

TCAndon-Router: Adaptive Reasoning Router for Multi-Agent Collaboration

Jiuzhou Zhao*, Chunrong Chen*, Chenqi Qiao*, Lebin Zheng*,
Minqi Han*, Yanchi Liu* Yongzhou Xu* Xiaochuan Xu* Min Zhang*

*Tencent Cloud Andon

{joskazhao, charentchen, chenqiqiao, lebinzheng}@tencent.com
{minqihan, yanncyliu, alanxu, xxcxu, alexzmzhang}@tencent.com

Abstract

Multi-Agent Systems(MAS) have become a powerful paradigm for building high performance intelligent applications. Within these systems, the router responsible for determining which expert agents should handle a given query plays a crucial role in overall performance. Existing routing strategies generally fall into two categories: performance routing, which balances latency and cost across models of different sizes, and task routing, which assigns queries to domain-specific experts to improve accuracy. In real-world enterprise applications, task routing is more suitable; however, most existing approaches rely on static single-label decisions, which introduce two major limitations: (i) difficulty in seamlessly integrating new agents as business domains expand, and (ii) routing conflicts caused by overlapping agent capabilities, ultimately degrading accuracy and robustness. To address these challenges, we propose TCAndon-Router(TCAR): an adaptive reasoning router for multi-agent collaboration. Unlike traditional routers, TCAR supports dynamic agent onboarding and first generates a natural-language reasoning chain before predicting a set of candidate agents capable of handling the query. In addition, we design a collaborative execution pipeline in which selected agents independently produce responses, which are then aggregated and refined into a single high-quality response by a dedicated Refining Agent. Experiments on public datasets and real enterprise data demonstrate that TCAR significantly improves routing accuracy, reduces routing conflicts, and remains robust in ambiguous scenarios. We have released TCAR at <https://huggingface.co/tencent/TCAndon-Router> to support future research on explainable and collaborative multi-agent routing.

1 Introduction

As large language models (LLMs) continue to improve in agent-based scenarios, multi-agent systems (MAS) have become increasingly mature. Their core idea is to decompose complex problems into smaller sub-tasks, each of which is handled by a specialized expert agent [1]. Representative applications include Vibe Coding [2] and Deep Research [3]. Today, large-scale enterprise systems—such as those in finance, healthcare, education, and cloud computing are also increasingly relying on multi-agent architectures to deliver accurate and efficient problem solving capabilities. Despite their broad applicability, MAS still face numerous challenges, including agent collaboration [4], routing [5], and security [6]. As these ecosystems become more complex, routing becomes particularly crucial: it determines which agent should handle each specific task and therefore sets the lower bound for the overall performance of the MAS.

Current routing strategies generally fall into two paradigms. The first is performance based routing, as depicted in 1(a) which dynamically selects LLMs of different sizes based on the estimated difficulty of a query to balance precision and cost[7, 8, 9]. The second is task based routing, as shown in 1(b) which assigns user queries to domain-specific expert agents to achieve higher task precision[5]. Performance based routing focuses on computational efficiency and model selection, whereas task based routing serves as the backbone of enterprise multi-agent systems by ensuring that queries are handled by domain experts rather than general-purpose models. Despite its strong performance, task based routing faces a critical limitation in real-world enterprise environments: **agent conflicts**. Unlike clean academic datasets, enterprise scenarios often involve overlapping agent responsibilities and multi-intent queries. For example, in cloud computing, an issue such as "website latency" can simultaneously be related to CDN configuration, public network quality, or application-layer bottlenecks. Traditional task based routers typically rely on single-label classification, which forces a single choice even when multiple experts

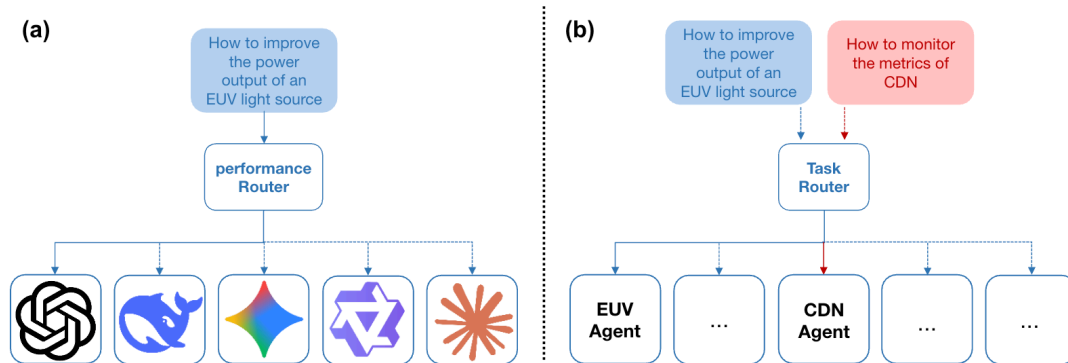


Figure 1: (a) performance router, which balances latency and cost across models of different sizes. (b) task router, which assigns queries to domain-specific experts to improve accuracy.

are appropriate. This leads to routing errors, brittle behavior under ambiguous queries, and reduced system reliability. Furthermore, most routers support only static routing, making them difficult to adapt to dynamic category changes in enterprise environments. The lack of explicit reasoning during routing also limits interpretability and robustness.

To address these challenges, we propose TCAR, adaptive reasoning router for multi-agent collaboration. Unlike existing routers that can only output a single label, TCAR generates structured natural language reasoning [10] and identifies all expert agents that may be relevant to the current problem. This design transforms routing from a traditional black-box classification task into a transparent and interpretable reasoning process, enabling the system to articulate when and why multiple agents might be applicable. Beyond the agent selection mechanism, we further introduce a novel collaborative execution pipeline. Each selected expert agent first produces a domain-specific response, after which a Refining Agent [11] aggregates and reconciles these outputs to form a logically coherent final answer. This architecture aligns with recent advances in collaborative multi-agent LLM systems, where domain experts provide complementary perspectives and a coordinator synthesizes the final decision [12, 13]. This paradigm also mirrors human organizational workflows: multiple subject-matter experts contribute their insights, and a senior analyst ultimately integrates them into a unified conclusion.

We evaluate TCAR on public benchmarks (CLINC150 [14], HWU64 [15], MINDS14 [16], and SGD [17]) as well as on a private dataset from real-world cloud computing scenarios. Experimental results demonstrate that combining multi-agent routing with a refinement mechanism significantly improves accuracy, reduces routing conflicts, and enhances system robustness when handling ambiguous or cross-domain queries.

The main contributions of this work are as follows:

- A reasoning-centric multi-agent routing system. We introduce TCAR, the first domain router that outputs natural-language decision rationales and supports flexible multi-agent selection.
- Conflict resolution through collaborative execution: We propose an innovative "multi-agent generation and refinement" pipeline that mitigates agent conflicts and effectively integrates complementary domain expertise.
- An open-source collaborative routing framework. We release the model and associated resources to advance research in interpretable multi-agent routing and reasoning-driven LLM collaboration.

2 Related Work

Performance based routing selects among models of different sizes to balance cost and latency. Hybrid LLM[18] proposes a model based on the BERT architecture to assess the difficulty of a query and make routing decisions. Self-REF[19] uses an LLM to output a confidence score, For low confidence, it uses an LLM or rejects the answer, while for high confidence, it uses an SLM.

Task based routing assigns queries to domain-specific expert agents to achieve higher accuracy. This work focuses on task based routing, whose methodological landscape can be grouped into several major

directions. Early approaches largely employed BERT-style encoders [20], formulating routing as a fixed single-label classification problem grounded in a predefined intent space and model structure. To support adding new agents without retraining, some studies shifted toward KNN-based methods, such as computing similarity between query and agent representations[21], or using query-agent dual-tower models or fusion-based cross-encoders [22] to dynamically compute relevance scores. [9]Train a binary classifier to determine whether the model performs excellently on a specific task. More recently, LLM-based routing methods have emerged—for example, ArchRouter [5] leverages human preference to perform agent selection. AgentRouter[23] Transform MAS into a knowledge graph, and decide which agent to select by predicting the edges through the model. Confidence-Driven LLM Router[24] decide which agent is better based on the confidence of the responses output by different agents, thereby obtaining labeled preference data, which is then used to train the routing model. However, these methods still follow the paradigm of producing a single agent output, lacking explicit and structured reasoning explanations. As a result, they struggle to handle the complexity and interpretability demands in real enterprise environments where multiple agents may simultaneously be relevant.

Multi-agent collaboration frameworks have also advanced rapidly in recent years. Representative systems include AutoGen [25], CAMEL [26], and CrewAI [27], which enable multiple agents to jointly solve complex tasks through mechanisms such as message passing, role-playing, and multi-turn dialogues. However, these frameworks typically rely on manually specified participant agents or simple rule-based selection. Their focus lies in how agents collaborate to execute a task rather than which agents should participate. Consequently, they lack a systematic mechanism for identifying the relevant subset of expert agents from a larger pool, and they pay limited attention to how the outputs of multiple agents should be integrated into a coherent final answer.

3 Method

We present TCAR, an adaptive reasoning router for multi-agent collaboration

3.1 Problem Formulation

A multi-agent system consists of a set of expert agents with different domain specializations. Let $A = \{a_1, a_2, \dots, a_N\}$ denote the set of all available expert agents. In real systems, each expert agent is not hard-coded into the router through a fixed identifier; instead, it exposes its capabilities and responsibilities to the routing model through a natural-language description. We denote this description as d_i , and define a natural-language "encoding function" g such that $g(a_i) = d_i$. The set of all descriptions is $D = \{d_1, d_2, \dots, d_N\}$. Given a user query q , the router does not operate directly on q alone. Instead, it receives a full prompt formed by concatenating the routing instruction, the user query, and all agent descriptions. Formally, we define a prompt construction function:

$$\text{Prompt}(q, A) = \Phi(\text{ins}, q, D)$$

where Φ denotes ordered text concatenation and ins is the router instruction. The routing task is then to output both the reasoning C_q for agent selection and the final selected agent set A_q :

$$\mathcal{R} : \text{Prompt}(q, A) \rightarrow (C_q, A_q)$$

3.2 Dynamic Agent Addition

Traditional methods typically bind a router \mathcal{R} to a fixed agent set A , when A changes \mathcal{R} must be retrained. In TCAR, adding a new agent requires only appending a_{new} to the existing set, yielding $A' = A \cup \{a_{\text{new}}\}$ which produces a new description set $D' = D \cup \{d_{\text{new}}\}$. The routing process then becomes $(C_q, A_q) = \mathcal{R}(\text{Prompt}(q, D'))$, $A_q \subseteq A$, meaning that the router \mathcal{R} itself does not need to be modified or retrained to support newly added agents. This property greatly improves efficiency and usability in real-world enterprise environments.

3.3 Multi-Agent Conflicts

Selecting Multiple Agents. Traditional task routing typically formulates routing as a single-label classification problem: $\mathcal{R}(\text{Prompt}(q, A)) = a_i \in A$ meaning that each query can only be assigned to a single agent. This formulation implicitly assumes strict and mutually exclusive domain boundaries, as well

as a clearly defined intent space—assumptions that rarely hold in real enterprise environments. When multiple agents are simultaneously suitable for handling a query due to overlapping domain responsibilities or cross-domain problem characteristics, an agent conflict occurs, formally defined as: $|A_q| > 1$. Existing single-label routers cannot represent this structural property; they are forced to pick only one agent among multiple plausible candidates. This leads to higher routing error rates, brittle behavior under cross-domain or ambiguous queries, and a lack of interpretability in routing decisions—ultimately undermining system stability and trustworthiness. To address this, we redefine routing as selecting a subset of agents from the full agent set: $A_q \subseteq A, |A_q| \geq 1$. The objective is no longer to choose a single “correct” agent, but to identify all experts that may be relevant to the query.

Reasoning-Based Agent Selection. To effectively model multi-agent conflict scenarios, we reformulate routing from a “direct label prediction” task into a two-stage “reason-then-select” process. Specifically, instead of outputting only a discrete agent label, the router generates both a natural-language reasoning chain C_q and the corresponding agent subset A_q :

$$(C_q, A_q) = \mathcal{R}(\text{Prompt}(q, A))$$

The reasoning C_q is required to explicitly analyze the potential causes of the user query, the relevant technical stack, and the responsibility boundaries of each agent. This produces a textual explanation of why the selected agents are relevant to the current query. This design provides two key advantages. First, the reasoning process encourages fine-grained semantic alignment between the query and agent descriptions, leading to more stable selection of the appropriate agent subset in cases of overlapping responsibilities or semantic ambiguity—rather than relying on brittle short-text matching. Second, explicit reasoning enhances interpretability: operations engineers can inspect C_q to diagnose routing errors, refine agent descriptions, or adjust routing strategies, forming a closed feedback loop of “data \rightarrow reasoning \rightarrow policy”. Experimental results show that, compared with black-box single-label or multi-label routers, incorporating natural-language reasoning significantly improves routing robustness under cross-domain, weakly-specified, and noisy queries.

Multi-Agent Collaborative Execution. Once the router is allowed to output multiple agents, potential inter-agent conflicts no longer need to be “compressed” into a single decision during routing. Instead, they are explicitly preserved and leveraged downstream through collaborative execution. Given the routing result $A_q = \{a_{i_1}, \dots, a_{i_k}\}$, $k = |A_q|$, the system invokes all selected candidate agents in parallel, with each agent producing a response to the same query $r_{i_j} = a_{i_j}(q)$. On top of these responses, we design a subsequent collaboration and aggregation mechanism. On the one hand, different agents provide partial yet specialized answers from their own perspectives, helping to cover diverse potential root causes in cross-domain queries. On the other hand, a downstream aggregation module—the Refining Agent—filters, aligns, and fuses $\{r_{i_j}\}$ into a single consistent final response. Formally, this collaborative process is defined as:

$$y_q = \mathcal{F}(\{r_{i_j}\}_{j=1}^k)$$

where \mathcal{F} denotes the overall operator of multi-agent collaboration and aggregation. Through the design of multi-agent selection + collaborative execution + aggregated fusion, the system no longer relies on a one-shot single-point decision to eliminate conflicts. Instead, conflicts are explicitly transformed into multi-perspective evidence that is unified at the answer level, thereby improving final response accuracy and robustness while maintaining high coverage.

3.4 Training TCAR

To endow the router \mathcal{R} with both multi-agent selection capabilities and interpretable reasoning, we adopt a two-stage[28] training strategy. We first apply supervised fine-tuning (SFT)[29] to teach the model the basic “reason-then-select” pattern. We then further enhance the quality of agent selection and the stability of reasoning through reward-based reinforcement learning using DAPO [30].

Supervised Fine-Tuning. In the first stage, we train the model using labeled data with standard causal language modeling. SFT equips the model with the following capabilities: (i) performing semantic alignment between the query and agent descriptions, thereby improving its generalization ability when dynamically onboarding new agents; and (ii) generating structured outputs, including the reasoning chain C_q and the candidate agent set A_q . Instead of adopting the conventional `<think>` tag used in “think”-style models, we introduce a dedicated `<reason>` tag to ensure compatibility with both instruction-following models and think models. The SFT stage is optimized using the standard autoregressive

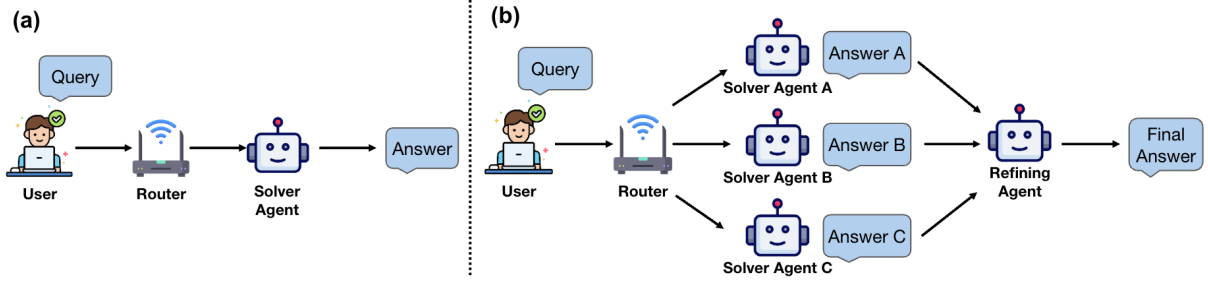


Figure 2: (a) Solo Agent: the router outputs a single agent, and that agent directly responds to the user’s query. (b) Refining Agent: the router outputs a set of candidate agents, and the refining agent integrates their outputs into a final answer.

language modeling loss:

$$\mathcal{L}_{\text{SFT}}(\theta) = - \sum_{t=1}^T \log p_{\theta}(y_t \mid y_{<t}, q, A),$$

$$\{y_1, \dots, y_T\} = (C_q, A_q).$$

However, SFT also has its limitations: the model may overfit to annotation templates or generate reasoning in a mechanical manner, making it necessary to further improve both the quality of reasoning and the accuracy of agent selection.

Reinforcement Learning. To improve the robustness of routing decisions and the soundness of the generated reasoning, we introduce a reinforcement learning stage on top of SFT. Here, we adopt the DAPO method to train the model, with the training objective formulated as:

$$\mathcal{J}_{\text{DAPO}}(\theta) = \mathbb{E}_{(q,a) \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot \mid q)} \left[\frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}}) \hat{A}_{i,t} \right) \right]$$

s.t. $0 < |\{o_i \mid \text{is_equivalent}(a, o_i)\}| < G$.

For reward design[31], since the SFT stage already enables the model to produce well-formatted outputs with a reasonable reasoning length C_q , we apply the reward only to the final agent set A_q . Because A_q is a set that may contain one or multiple agents, we decompose the routing reward into two parts. Let A^* denote the gold agent set and A_q the predicted candidate agent set. The first reward focuses on whether each element in A_q belongs to A^* (analogous to precision), while the second reward evaluates whether A_q covers all correct agents (analogous to recall or a coverage constraint). Formally, the rewards are defined as follows:

First, we define a local reward based on the degree of set overlap, which measures the proportion of agents in the predicted set A_q that are correct. This reward encourages the model to reduce irrelevant or incorrect agent predictions.

$$R_1(A_q, A^*) = \frac{|A_q \cap A^*|}{\max(1, |A_q|)},$$

Second, we introduce a coverage reward to penalize missing correct agents. The model receives an additional reward if and only if its predicted set A_q fully contains the ground-truth set A^* ; otherwise, this reward is zero.

$$R_2(A_q, A^*) = \frac{|A_q \cap A^*|}{\max(1, |A^*|)},$$

To prevent the model from generating excessively long outputs or endlessly listing agents, we introduce an additional length-penalty term[32]. Combining all three components, we define the final reward as:

$$R(A_q, A^*) = \alpha R_1(A_q, A^*) + (1 - \alpha) R_2(A_q, A^*) - \beta \max(|A_q| - |A^*|, 0)$$

where α and β are weighting hyperparameters. The coefficient α balances the trade-off between the correctness of the predicted set (precision-like) and its coverage of the ground-truth set (recall-like), while β controls the strength of the length-penalty term. In practice, we directly incorporate this reward into the policy optimization objective during the reinforcement learning stage, encouraging the router to avoid over-predicting irrelevant agents while still covering all truly relevant experts in multi-agent scenarios.

4 Experiments

We present the experimental results and analyses.

4.1 Experimental Setup

This section introduces the datasets, baselines, evaluation metrics, and implementation details used to evaluate TCAR. All experiments focus on the core task of multi-agent domain routing, aiming to thoroughly assess TCAR’s effectiveness under domain conflicts, ambiguous intents, and real-world enterprise scenarios.

Datasets. We conduct evaluations on four public intent-classification datasets and one internal enterprise dataset:

- **CLINC150** [14]: Contains 150 intent categories covering everyday conversational scenarios. The task boundaries are clear, but the large number of classes leads to long agent descriptions, posing challenges for LLMs.
- **HWU64** [15]: Includes 64 intents with relatively high cross-domain ambiguity.
- **MINDS14** [16]: A multilingual dataset spanning 14 domains, used to test router stability under cross-lingual semantic transfer.
- **SGD** [17]: A multi-turn dialogue dataset designed to evaluate routing performance in multi-turn conversational settings.
- **QCloud**: Additionally, we construct a real-world enterprise dataset from Tencent Cloud, covering domains such as networking, cloud storage, CDN, security, databases, and application operations. Compared with public datasets, it exhibits stronger domain coupling, higher intent ambiguity, and more frequent agent conflicts, making it a crucial benchmark for assessing the practicality of multi-agent routing systems.

Evaluation Metrics. To comprehensively evaluate routing performance, we adopt the following metrics:

- For single-agent datasets, we primarily report **accuracy**. For multi-agent datasets, we compute **F1** to measure whether the model successfully and accurately identifies all agents capable of handling a given query.
- **End-to-End Task Success Rate**, which evaluates the overall response quality after multi-agent collaboration and Refining Agent processing, offering a system-level assessment of routing capability.

Baselines. Since static label routing is difficult to deploy in most real-world enterprise scenarios, we focus on comparing against dynamic routing approaches. We evaluate a variety of both proprietary and open-source large language models, including GPT-5.1[33], Claude-4.5[34], DeepSeek-v3.1[35], ArcRouter[5], as well as the Qwen3[36] family, which also serves as our base model for training.

Implementation Details. We experimented with several model sizes from the Qwen3 family and ultimately selected Qwen3-4B-Instruct-2507 for open-sourcing, achieving a strong balance between effectiveness and efficiency. To ensure training data compatibility and avoid reliance on model-specific reasoning formats, we did not use Qwen3’s native `<think>` tag. Instead, we introduced a unified `<reason>` tag for generating reasoning chains. Full-parameter SFT was conducted using the ms-swift framework[37] with data parallelism across 8 GPUs, Adam optimizer, batch size of 256, training for 1 epoch, an initial learning rate of $2e-5$, and a warmup ratio of 0.1. During the RL stage, we use the SFT model with a

temperature of 1 and perform 8 rollouts per training sample. If the rollout consistency exceeds $\tau > 0.6$, we consider the sample to be already well-learned during SFT—essentially a low-entropy token [38]. Such samples tend to cause near-zero variance when computing advantages during RL, leading to ineffective optimization. Therefore, following the principles of [30], we remove these samples from the RL training set to reduce unnecessary dynamic sampling and improve training efficiency. We also experimented with initializing the model using Slerp[39, 40], which allows the model to achieve higher pass@k performance at the early stage of RL training. We aim for the model in the RL phase to transition from achieving high pass@k to achieving high pass@1[41]. Thereby preventing the model from suffering entropy collapse[42]. Both stages of training (SFT and RL) are implemented using the ms-swift framework.

4.2 Main Results

Table 1 presents the performance of all models across the five datasets (CLINC150, HWU64, MINDS14, SGD, QCloud). Overall, general-purpose LLMs such as GPT-5.1 and DeepSeek-v3.1 outperform most small-scale open-source models by a considerable margin. Despite having only 4B parameters, our TCAR achieves state-of-the-art performance on all datasets except CLINC150. CLINC150 contains 150 intent categories, resulting in extremely long prompts (averaging 18k tokens), which pose challenges for smaller models that struggle with ultra-long sequence processing. On multi-turn (SGD) and multilingual (MINDS14) datasets, TCAR demonstrates clear and substantial advantages. On the QCloud dataset which features high agent conflict frequency and real enterprise ambiguity—TCAR surpasses even current leading general-purpose LLMs in F1, highlighting its robustness in multi-agent routing scenarios.

Models	CLINC150	HWU64	MINDS14	SGD	QCloud
GPT-5.1	93.84	85.59	95.59	73.90	92.93
Claude-Sonnet-4.5	94.21	87.40	96.20	76.02	91.45
DeepSeek-v3.1-terminus	88.29	88.10	95.72	79.70	92.98
ArcRouter	62.98	69.33	91.79	65.59	-
Qwen3-Embedding-4B	57.21	54.27	94.12	37.02	-
Qwen3-4B-Instruct-2507	70.12	80.29	90.08	58.74	80.81
TCAR(4B)	91.25	91.63	96.70	91.58	93.98

Table 1: Comparison of model performance across datasets. For the QCloud dataset, the metric is reported as F1, while for all other datasets the metric is Accuracy.

We further evaluate the effectiveness of the downstream Refining Agent when TCAR outputs multiple candidate agents. As shown in Figure 3, the Refining Agent implemented using DeepSeek-v3.1 as the underlying LLM—significantly enhances the quality of response in our QCloud data set. For consultation-type queries, a single agent is often sufficient to produce high quality responses. However, for troubleshooting-type queries, an individual agent typically fails to provide complete coverage, whereas the Refining Agent integrates and consolidates responses from multiple agents, resulting in more comprehensive and reliable answers.

4.3 Ablation Studies

To systematically analyze the contribution of each component in TCAR, we conduct ablation studies along three dimensions: reasoning capability, training strategy, and the source of RL initialization.

Effect of Reasoning This experiment evaluates the impact of explicit reasoning chains[10] on routing performance.

- **Without reasoning:** The model performs direct classification prediction and outputs only A_q .
- **With reasoning:** The model first generates C_q as an explicit and interpretable reasoning chain and then outputs the set of agents A_q .

We evaluated both settings on the public datasets and the QCloud dataset. The results 2 show that even without reasoning, the model achieves strong performance after fine-tuning; however, incorporating reasoning consistently leads to further improvements. This observation suggests that reasoning enhances the model’s generalization ability.

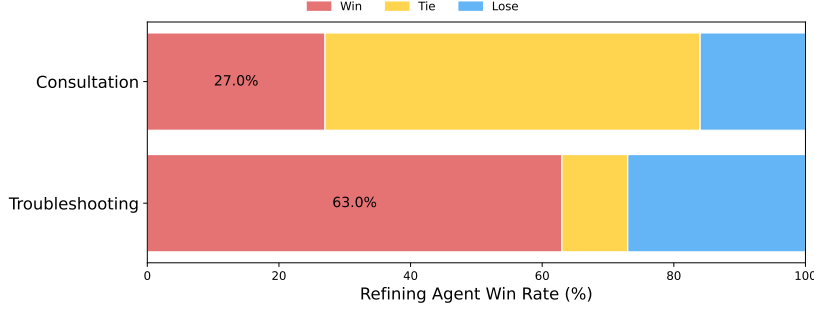


Figure 3: In scenarios where consultation-type and troubleshooting-type queries create conflicts among candidate agents, we compare two strategies: randomly selecting a single agent to respond versus allowing multiple agents to answer followed by aggregation by the Refining Agent. The results show that the Refining Agent achieves a significantly higher human preference win rate, particularly for troubleshooting-type queries.

Reasoning Setting	CLINC150	HWU64	MINDS14	SGD	QCloud
TCAR (without Reasoning)	89.65	87.83	96.45	84.96	91.33
TCAR (with Reasoning)	91.25	91.63	96.70	91.58	93.98

Table 2: Comparison of the performance with and without reasoning output

SFT-Only vs. SFT+RL This experiment evaluates the improvements brought by applying DAPO during the RL stage.

- **SFT-only**: The model is trained using supervised fine-tuning only.
- **SFT+RL**: After SFT, we further filter the training set and retain only high-entropy tokens for RL.

We compare the performance of the SFT-only model with the model further optimized through RL. Experimental results show that RL leads to consistent and significant improvements. In particular, on the QCloud dataset, RL enhances recall while maintaining high precision, effectively mitigating the common issue of SFT-only models being overly conservative and outputting only a single agent.

Training Method	CLINC150	HWU64	MINDS14	SGD	QCloud
SFT-Only	91.32	90.30	96.69	86.30	93.77
SFT+RL	91.25	91.63	96.70	91.58	93.98

Table 3: Comparison between SFT-only training and SFT followed by RL.

RL Initialization: SFT-Only vs. SFT+Slerp Since the initialization point of RL has a substantial impact on the final performance, we compare two different SFT models used for RL initialization:

- **SFT-Only**: RL initialized from the SFT model.
- **SFT+Slerp**: RL initialized from the model obtained via SFT-Slerp merging.

Slerp is a model-merging technique [39, 40] that effectively enhances model generalization. We trained five different SFT models on five distinct datasets and then merged them using Slerp. The resulting merged model exhibits higher exploration capability in the early stages of RL, reflected by a higher pass@k score. We aim for the model in the RL phase to transition from achieving high pass@k to achieving high pass@1[41]. We conducted RL training using both the SFT model and the Slerp-merged model. As shown in Table 4, the Slerp-based model does not exhibit a significant improvement in the final metrics. However, from Figure 4b, we observe that although the reward curves of the two initialization strategies are similar, the Slerp model maintains a higher entropy, indicating stronger exploration. This

RL Initialization	CLINC150	HWU64	MINDS14	SGD	QCloud
SFT-Only	91.25	91.63	96.70	91.58	93.97
SFT+Slerp	90.43	91.45	96.82	91.53	93.20

Table 4: Comparison between initializing the RL stage with the SFT model and initializing it with the Slerp-merged model.

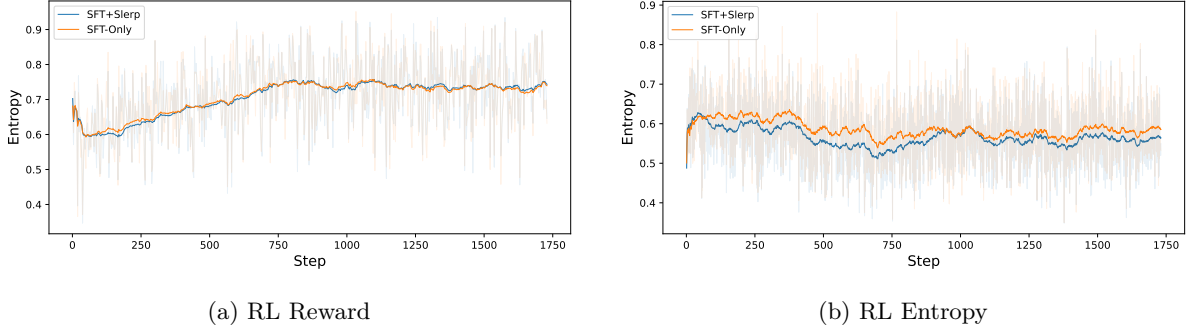


Figure 4: Comparing RL initialized from the SFT model and from the Slerp-merged model, we observe that both approaches achieve almost the same reward throughout training. In contrast, the Slerp-initialized model maintains higher entropy, suggesting that it explores a broader solution space.

enhanced exploration leads to better generalization and provides greater potential for further fine-tuning based on the merged model.

Summary of Ablation Findings

Based on the ablation results, we summarize the following conclusions:

- Reasoning capability is a key factor for improving routing interpretability and robustness.
- RL is an essential training stage that enhances recall while maintaining high precision.
- The choice of RL initialization point is crucial; Slerp-based models offer stronger generalization and are less prone to entropy collapse during RL.

5 Analysis

After validating the overall effectiveness of TCAR through the main results and ablation studies, we further analyze its behavioral characteristics and internal mechanisms from multiple perspectives. This section provides an in-depth examination of error types, reasoning quality, multi-agent collaboration effects, and system efficiency, offering a more comprehensive understanding of the model’s strengths and limitations.

5.1 Error Analysis

To systematically understand the major sources of routing errors, we categorize the mistakes and observe that TCAR’s errors primarily fall into three types: (1) **Incomplete problem descriptions**. User queries may lack key information. For example, in “What should I do if my webpage loads slowly?”, although TCAR’s reasoning often infers that multiple agents are needed, the lack of essential context may still lead to deviations; see Appendix B. (2) **Insufficient understanding of highly domain-specific scenarios**. In Tencent Cloud use cases, the model sometimes struggles with understanding specialized terminology, leading to incorrect routing decisions; see Appendix B. (3) **Mismatch between reasoning and final prediction**. Some cases exhibit reasonable reasoning steps but produce agent outputs inconsistent with the reasoning itself. How to further improve the alignment and consistency of the reasoning chain remains an open problem.

5.2 Multi-Agent Collaboration Analysis

Multi-agent collaboration effectively supplements information from different perspectives. In Tencent Cloud troubleshooting scenarios, multiple agents are often required to jointly diagnose and resolve issues, as a single agent is typically insufficient. Moreover, the Refining Agent demonstrates strong capability in reconciling conflicting responses. With RAG support, the Refining Agent compares and filters multiple candidate answers using the knowledge base, enabling conflicting opinions to converge. For instance, in troubleshooting scenarios localization tasks, partial correctness from different agents can be combined to form a complete solution—something a single agent cannot achieve.

5.3 Efficiency and Cost Analysis

We further evaluate the deployment cost of TCAR in online environments. (1) **The number of selected agents is far below the theoretical maximum.** Although multi-agent outputs are supported, the actual average number of routed agents is only 1.37, indicating that most queries can be solved by a single agent and that the model naturally favors concise predictions. (2) **The latency impact of the reasoning chain is well controlled.** For a 4B model deployment requirements remain modest; the average reasoning length is under 100 tokens. When deployed on GPUs, the average latency is less than 1 second, which is entirely acceptable. (3) **The incremental cost of downstream collaboration is small.** Because the number of selected agents is low, downstream execution incurs an average overhead equivalent to only 0.37 additional agent calls—far below any risk of combinatorial explosion.

6 Limitations

Although TCAR demonstrates stable and significant performance improvements across multiple public datasets and real-world enterprise scenarios, several limitations remain that warrant attention in future work.

Dependence on the quality of agent descriptions. TCAR’s routing reasoning relies on manually written agent descriptions, which serve as the semantic foundation for constructing reasoning chains. When these descriptions are overly brief, ambiguous, or fail to fully represent the true scope of an agent’s responsibilities, the resulting reasoning chain may deviate from the intended semantics, ultimately affecting routing accuracy. Although reinforcement learning improves the model’s robustness to noise in descriptions, the overall quality of these descriptions remains a critical factor influencing performance.

Challenges in long-tail knowledge and domain transfer. For low-frequency or highly domain-specific scenarios (e.g., rare APIs or special network configurations), the model may still exhibit instability due to insufficient training samples. Furthermore, when the instruction format is heavily modified, the model’s instruction-following capability may degrade. If the instructions are substantially altered, an additional round of fine-tuning on domain-specific private data becomes a more reasonable strategy.

Overall, these limitations do not undermine TCAR’s effectiveness in mainstream routing tasks or real enterprise deployments, but they highlight promising directions for future research aimed at improving robustness, generalization, and performance under extreme or specialized conditions.

7 Conclusion

This paper addresses the prevalent challenges in real-world multi-agent enterprise systems, including agent conflicts, cross-domain intent ambiguity, and limited interpretability. We propose **TCAR**, a reasoning-centric multi-agent routing framework. By explicitly generating a natural-language reasoning chain, TCAR extends traditional single-agent classification into a flexible subset prediction over all potentially relevant agents, significantly enhancing the coverage and interpretability of routing decisions. Building on this, we further design a downstream multi-agent collaborative execution mechanism and a Refining Agent that consolidates multiple agent outputs, enabling the system to harness complementary expertise across complex or cross-domain tasks. Experiments on multiple public datasets and large-scale real ITSM data from Tencent Cloud demonstrate that TCAR achieves substantial improvements in task accuracy, recall under conflict-prone scenarios, and the final answer quality resulting from multi-agent collaboration. Ablation studies further validate the effectiveness of reasoning chains, reinforcement

learning, and the SLERP-based model fusion strategy. Meanwhile, our analyses reveal key factors influencing model performance, including reasoning-chain structure, agent description quality, and the sparsity of long-tail domain knowledge. In summary, TCAR highlights the feasibility and practical value of a reasoning-driven, multi-agent collaborative routing paradigm in large-scale enterprise environments. For future work, we plan to explore structured reasoning-chain constraints, more efficient collaboration protocols, and low-cost expansion strategies for emerging business domains to further enhance robustness and generalization in complex real-world scenarios.

References

- [1] Ali Dorri, Salil S. Kanhere, and Raja Jurdak. Multi-agent systems: A survey. *IEEE Access*, 6:28573–28593, 2018.
- [2] Christian Meske, Tobias Hermanns, Esther von der Weiden, Kai-Uwe Loser, and Thorsten Berger. Vibe coding as a reconfiguration of intent mediation in software development: Definition, implications, and research agenda. *arXiv preprint arXiv:2507.21928*, 2025.
- [3] OpenAI. Introducing deep research. <https://openai.com/index/introducing-deep-research/>, 2025.
- [4] Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’Sullivan, and Hoang D Nguyen. Multi-agent collaboration mechanisms: A survey of llms. *arXiv preprint arXiv:2501.06322*, 2025.
- [5] Co Tran, Salman Paracha, Adil Hafeez, and Shuguang Chen. Arch-router: Aligning llm routing with human preferences. *arXiv preprint arXiv:2506.16655*, 2025.
- [6] Zhixin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. Agent-safetybench: Evaluating the safety of llm agents. *arXiv preprint arXiv:2412.14470*, 2024.
- [7] OpenAI. Introducing gpt-5. <https://openai.com/index/introducing-gpt-5/>, 2025. Accessed: 2025-11-25.
- [8] Wittawat Jitkrittum, Harikrishna Narasimhan, Ankit Singh Rawat, Jeevesh Juneja, Congchao Wang, Zifeng Wang, Alec Go, Chen-Yu Lee, Pradeep Shenoy, Rina Panigrahy, et al. Universal model routing for efficient llm inference. *arXiv preprint arXiv:2502.08773*, 2025.
- [9] Tal Shnitzer, Anthony Ou, Mírian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. Large language model routing with benchmark datasets. *arXiv preprint arXiv:2309.15789*, 2023.
- [10] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc., 2022.
- [11] Soyeong Jeong, Aparna Elangovan, Emine Yilmaz, and Oleg Rokhlenko. Adaptive multi-agent response refinement in conversational systems. *arXiv preprint arXiv:2511.08319*, 2025.
- [12] Florian Grötschla, Luis Müller, Jan Tönshoff, Mikhail Galkin, and Bryan Perozzi. Agentsnet: Coordination and collaborative reasoning in multi-agent llms. *arXiv preprint arXiv:2507.08616*, 2025.
- [13] Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of intelligent llm agents. *arXiv preprint arXiv:2306.03314*, 2023.
- [14] Stefan Larson, Anish Mahendran, Joseph J Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K Kummerfeld, Kevin Leach, Michael A Laurenzano, Lingjia Tang, et al. An evaluation dataset for intent classification and out-of-scope prediction. *arXiv preprint arXiv:1909.02027*, 2019.
- [15] Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. Benchmarking natural language understanding services for building conversational agents. In *Increasing naturalness and flexibility in spoken dialogue interaction: 10th international workshop on spoken dialogue systems*, pages 165–183. Springer, 2021.
- [16] Daniela Gerz, Pei-Hao Su, Razvan Kuztos, Avishek Mondal, Michal Lis, Eshan Singhal, Nikola Mrksic, Tsung-Hsien Wen, and Ivan Vulic. Multilingual and cross-lingual intent detection from spoken data. *CoRR*, abs/2104.08524, 2021.
- [17] Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. Schema-guided dialogue state tracking task at dstc8. *arXiv preprint arXiv:2002.01359*, 2020.

- [18] Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Ruhle, Laks VS Lakshmanan, and Ahmed Hassan Awadallah. Hybrid llm: Cost-efficient and quality-aware query routing. *arXiv preprint arXiv:2404.14618*, 2024.
- [19] Yu-Neng Chuang, Prathusha Kameswara Sarma, Parikshit Gopalan, John Boccio, Sara Bolouki, Xia Hu, and Helen Zhou. Learning to route llms with confidence tokens. *arXiv preprint arXiv:2410.13284*, 2024.
- [20] Toby Simonds, Kemal Kurniawan, and Jey Han Lau. Modem: Mixture of domain expert models. *arXiv preprint arXiv:2410.07490*, 2024.
- [21] Deepak Babu Piskala, Vijay Raajaa, Sachin Mishra, and Bruno Bozza. Dynamic llm routing and selection based on user preferences: Balancing performance, cost, and ethics. *arXiv preprint arXiv:2502.16696*, 2025.
- [22] Justin Chiu and Keiji Shinzato. Cross-encoder data annotation for bi-encoder based product matching. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 161–168, 2022.
- [23] Zheyuan Zhang, Kaiwen Shi, Zhengqing Yuan, Zehong Wang, Tianyi Ma, Keerthiram Murugesan, Vincent Galassi, Chuxu Zhang, and Yanfang Ye. Agentrouter: A knowledge-graph-guided llm router for collaborative multi-agent question answering. *arXiv preprint arXiv:2510.05445*, 2025.
- [24] Tuo Zhang, Asal Mehradfar, Dimitrios Dimitriadis, and Salman Avestimehr. Leveraging uncertainty estimation for efficient llm routing. *arXiv preprint arXiv:2502.11021*, 2025.
- [25] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- [26] Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [27] João Moura and CrewAI contributors. Crewai: A python framework for multi-agent automation. <https://github.com/crewAIInc/crewAI>, 2025. Version X.Y.Z (please replace with actual version), open-source.
- [28] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [29] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- [30] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- [31] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [32] Hongze Tan, Jianfei Pan, Jinghao Lin, Tao Chen, Zhihang Zheng, Zhihao Tang, and Haihua Yang. Gtpo and grpo-s: Token and sequence-level reward shaping with policy entropy. *arXiv preprint arXiv:2508.04349*, 2025.
- [33] OpenAI. Gpt-5.1. <https://openai.com/index/gpt-5-1/>, 2025. Accessed: 2025-01-10.
- [34] Anthropic. Claude sonnet 4.5. <https://www.anthropic.com/news/claude-sonnet-4-5>, 2025. Accessed: 2025-01-10.
- [35] DeepSeek-AI. Deepseek-v3 technical report, 2024.

- [36] Qwen Team. Qwen3 technical report, 2025.
- [37] Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. Swift:a scalable lightweight infrastructure for fine-tuning, 2024.
- [38] Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025.
- [39] Ken Shoemake. Animating rotation with quaternion curves. In *Proceedings of the 12th annual conference on Computer graphics and interactive techniques*, pages 245–254, 1985.
- [40] mlabonne. Merge large language models with mergekit. <https://huggingface.co/blog/mlabonne/merge-models>, 2025. Accessed: 2025-12-09.
- [41] Christian Walder and Deep Karkhanis. Pass@ k policy optimization: Solving harder reinforcement learning problems. *arXiv preprint arXiv:2505.15201*, 2025.
- [42] Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.

A prompt

TCAR Prompt:

You are a Router Agent. Your job is to determine which Agents should handle the user's issue based on the chat history. Your goal is to select the appropriate Agents to ensure the problem can be solved efficiently and accurately. You may select up to 3 Agents.

Agent List
{agents}

Notes
+ The Description is each Agent's capability description. When selecting an Agent, you must ensure that its capabilities clearly allow it to resolve the user's problem.
+ If no Agent can reasonably handle the problem, output "oos" in the <ID> tag.

Output Format
Output in XML format. If multiple Agents are needed, output multiple <ID> lines.
<Reason>The reason for choosing the Agent(s)</Reason>
<ID>Agent ID</ID>

Chat History

Refining Agent Prompt:

You are a Refining Agent. I will provide you with a user query along with the answers generated by different agents. Your task is to compare and integrate these answers and produce the most helpful and accurate final response for the user.

Requirements
Please output one final answer that satisfies the following criteria:
+ Fully compare the content from different agents, and combine the strengths of each answer to produce the response that best solves the user's problem. Avoid repetition, errors, and irrelevant information.
+ If there are conflicts among the answers, you must present all conflicting points in the final response and explicitly explain the conflict.
+ If the answers are complementary, you should integrate them into a unified, coherent response.
+ If all answers fail to provide useful information or all refuse to answer the query, you must output a polite refusal to answer.
+ Output only the final answer — do not include any other commentary.

Query: {query}
Answers: {answers}
Final Answer:

B Badcase

Incomplete problem descriptions

Input Query: What should I do if a webpage loads very slowly?
TCAR Output:
<Reason> The user reports that "the webpage loads very slowly," which requires diagnosis from both the network path and acceleration layers. The Network Agent can analyze public/private network latency, packet loss, jitter, and path quality. The CDN Agent can examine CDN access latency, origin fetch behavior, and distribution policies. If EdgeOne is being used, the EO Agent should investigate origin latency, edge node issues, and configuration problems. </Reason>
<ID>network_agent</ID>
<ID>cdn_agent</ID>
<ID>eo_agent</ID>

Insufficient understanding of highly domain-specific scenarios

Agent List:

```
<Agent>
  <Name>CKafka Message Queue Agent</Name>
  <Description>
    Query CKafka instance information and diagnose CKafka-related issues.</Description>
  <ID>ckafka_agent</ID>
</Agent>

<Agent>
  <Name>RocketMQ Message Queue Agent</Name>
  <Description>
    Analyze abnormal RocketMQ topic production rates and query configuration changes (
      scaling/reconfiguration).
  </Description>
  <ID>rocketMQ_agent</ID>
</Agent>
```

Input Query: Configuration change for CKafka message queue

TCAR Output:

<Reason> The user reports a “message queue configuration change.” The RocketMQ Message Queue Agent has the capability to handle configuration change queries and can effectively address the user’s issue. </Reason>

<ID>rocketMQ_agent</ID>