

# TRIAGEAGENT: Towards Better Multi-Agents Collaborations for Large Language Model-Based Clinical Triage

Anonymous EMNLP submission

## Abstract

The global escalation in emergency department patient visits poses significant challenges to efficient clinical management, particularly in clinical triage. Traditionally managed by human professionals, clinical triage is susceptible to substantial variability and high workloads. Although large language models (LLMs) demonstrate promising reasoning and understanding capabilities, directly applying them to clinical triage remains challenging due to the complex and dynamic nature of the clinical triage task. To address these issues, we introduce TRIAGEAGENT, a novel heterogeneous multi-agent framework designed to enhance collaborative decision-making in clinical triage. TRIAGEAGENT leverages LLMs for role-playing, incorporating self-confidence and early-stopping mechanisms in multi-round discussions to improve document reasoning and classification precision for triage tasks. In addition, TRIAGEAGENT employs the medical Emergency Severity Index (ESI) handbook through a retrieval-augmented generation (RAG) approach to provide precise clinical knowledge and integrates both coarse- and fine-grained ESI-level predictions in the decision-making process. Extensive experiments demonstrate that TRIAGEAGENT outperforms state-of-the-art LLM-based methods on three clinical triage test sets. Furthermore, we have released the first public benchmark dataset for clinical triage with corresponding ESI levels and human expert performance for comparison.

## 1 Introduction

Emergency Departments (EDs) play a crucial role in the healthcare system by continuously assessing and prioritizing patients based on urgency and severity. This process, known as **clinical triage**, utilizes the **Emergency Severity Index (ESI)** as a standardized guide for decisions on rapid medical intervention, which is vital for prioritizing treatment and allocating resources. However, the

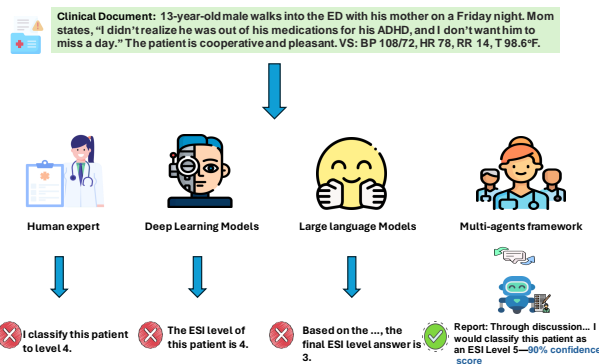


Figure 1: Workflow illustrations of clinical triage task.

growing number of patients poses significant challenges to the rapid and precise classification of cases, which is crucial for accurate ESI categorization. Currently, hospitals rely on human experts to review clinical notes and determine case urgency (as illustrated in Figure 1). Although effective, this manual method is time-consuming, labor-intensive, and burdensome for clinical staff. The increasing patient volume and complex triage process often lead to staff fatigue, diminishing accuracy and efficiency, and raising the risk of inconsistent classification or misdiagnosis.

Consequently, there is a high demand for AI methods to automate ESI classification. Traditional deep learning (DL) models (as illustrated in Figure 1) (Kojima et al., 2023; Yao et al., 2021; Sánchez-Salmerón et al., 2022) have assisted in clinical triage but often fall short due to the complex and dynamic nature of the task, which requires extensive labeled data and real-time adaptation. LLMs such as GPT (Kojima et al., 2023; OpenAI et al., 2024), Med-PaLM (Chowdhery et al., 2022), and Llama (Touvron et al., 2023) offer promising solutions with advanced text understanding capabilities, reducing time costs and errors by quickly interpreting and categorizing clinical documents. Additionally, LLMs can leverage external tools, such as knowledge base APIs (Qin et al., 2023; Zhuang

071 et al., 2023), to enhance domain-specific knowl- 122  
072 edge, adaptability, speed, and accuracy. However, 123  
073 even with these tools, a single LLM may struggle 124  
074 with the complexities of emergency scenarios, such 125  
075 as diverse patient presentations and the need for im- 126  
076 mediate multi-disciplinary coordination (Chenais 127  
077 et al., 2023). Therefore, a multi-agent approach is 128  
078 necessary to improve performance by distributing 129  
079 tasks and assigning specialized roles to agents with 130  
080 diverse expertise. 131

081 Recent advancements have significantly en- 132  
082 hanced multi-agent systems in areas such as rea- 133  
083 soning (Wang et al., 2023b), sophisticated plan- 134  
084 ning (Yao et al., 2023; Sun et al., 2023), and 135  
085 memory (Wang et al., 2023a). These improve- 136  
086 ments enable multi-agent LLMs to analyze medi- 137  
087 cal data, formulate treatment plans, and recall pa- 138  
088 tient histories more effectively (Tang et al., 2024). 139  
089 However, directly deploying multi-agent LLMs 140  
090 in clinical triage remains challenging due to their 141  
091 lack of optimization for triage-specific decision- 142  
092 making, resulting in accuracy levels around 60% 143  
093 (as illustrated in Table 2). This suboptimal per- 144  
094 formance stems from the intricate coordination re- 145  
095 quired among agents and the need to design archi- 146  
096 tectures that better utilize LLMs. Additionally, the 147  
097 lack of open-source benchmark datasets hampers 148  
098 comparisons with human experts, thereby affecting 149  
099 the practical credibility of multi-agent LLMs. 150

100 We identify four major challenges hindering clin- 151  
101 ical triage effectiveness. The first is **data process-** 152  
102 **ing speed**. The increasing patient volume neces- 153  
103 sitates rapid data processing and decision-making, 154  
104 yet traditional methods often face delays due to con- 155  
105 tinuous data updates and extensive preprocessing, 156  
106 hindering timely clinical interventions. The sec- 157  
107 ond is **diversity in clinical documents**. Variability 158  
108 in patient histories, symptoms, writing styles and 159  
109 terminologies complicates the classification task. 160  
110 The third is **contextual understanding and inter-** 161  
111 **pretability**. The complexity of clinical contexts 162  
112 makes it challenging for models to accurately inter- 163  
113 pret information and provide transparent, evidence- 164  
114 based results, leading to a lack of clinician trust. 165  
115 The fourth is **boundaries between different ESI** 166  
116 **levels**. Precise differentiation and classification of 167  
117 ESI levels (1-5) are essential to avert critical medi- 168  
118 cal errors. For instance, distinguishing between a 169  
119 level 1 (most urgent) and a level 2 (less urgent) case 170  
120 is vital, as it determines whether a patient receives 171  
121 immediate medical intervention or encounters a

prolonged wait. Nevertheless, the subtle nuances 122  
between different ESI levels present considerable 123  
challenges to the clinical triage process. 124

To address these challenges, we propose 125  
TRIAGEAGENT, a novel heterogeneous multi-agent 126  
collaboration framework for clinical triage that 127  
leverages LLM-based agents enhanced with ex- 128  
ternal tools and embedded medical knowledge. 129  
TRIAGEAGENT enables effective information ex- 130  
change and reliable interactions among agents, 131  
making the framework ideal for zero-shot docu- 132  
ment classification and handling complex tasks 133  
without prior demonstrations. The **key innovations** 134  
include retrieval-augmented generation for context- 135  
relevant evidence, a confidence score-based mech- 136  
anism for precise decision-making, and an early 137  
stopping mechanism to improve time efficiency. 138  
These features enhance contextual understanding, 139  
interpretability, and precision in ESI-level classifi- 140  
cation. Additionally, the framework supports real- 141  
time decision-making and dynamic, interactive de- 142  
bates among agents, refining information for more 143  
accurate triage and improving the timeliness of 144  
clinical interventions. 145

146 Experiments on three ESI clinical triage test sets 147  
demonstrate that TRIAGEAGENT significantly im- 148  
proves zero-shot performance with GPT-3.5-turbo 149  
and GPT-4, reducing discordance rates by up to 150  
10.84% and 18.42%, respectively. Additionally, 151  
we are the first to publicly release a clinical triage 152  
dataset that includes clinical notes alongside ESI 153  
levels and human expert performance, providing a 154  
valuable resource for academic research and clini- 155  
cal practice. By setting new benchmarks, we aim to 156  
advance the field of clinical triage in both academic 157  
and practical applications. We will publicly release 158  
our code and dataset once the paper is published. 159

## 160 2 Related Work 161

### 162 2.1 LLM Applications in Healthcare 163

164 Large Language Models (LLMs) have recently ex- 165  
166 perience significant advancements across various 167  
168 fields, particularly in healthcare (Ling et al., 2024; 169  
170 Bi et al., 2024; Nori et al., 2023; Bao et al., 2023). 171  
These models are increasingly employed in medical 172  
applications, including text-based diagnostics (Ma 173  
et al., 2024), genetic analysis (Bi et al., 2024), phar- 174  
maceutical applications (Liu et al., 2023), and med- 175  
ical summary generation (Shaib et al., 2023). Cur- 176  
rent research on LLMs in healthcare primarily fo- 177  
cuses on integrating external tools to enhance clini-

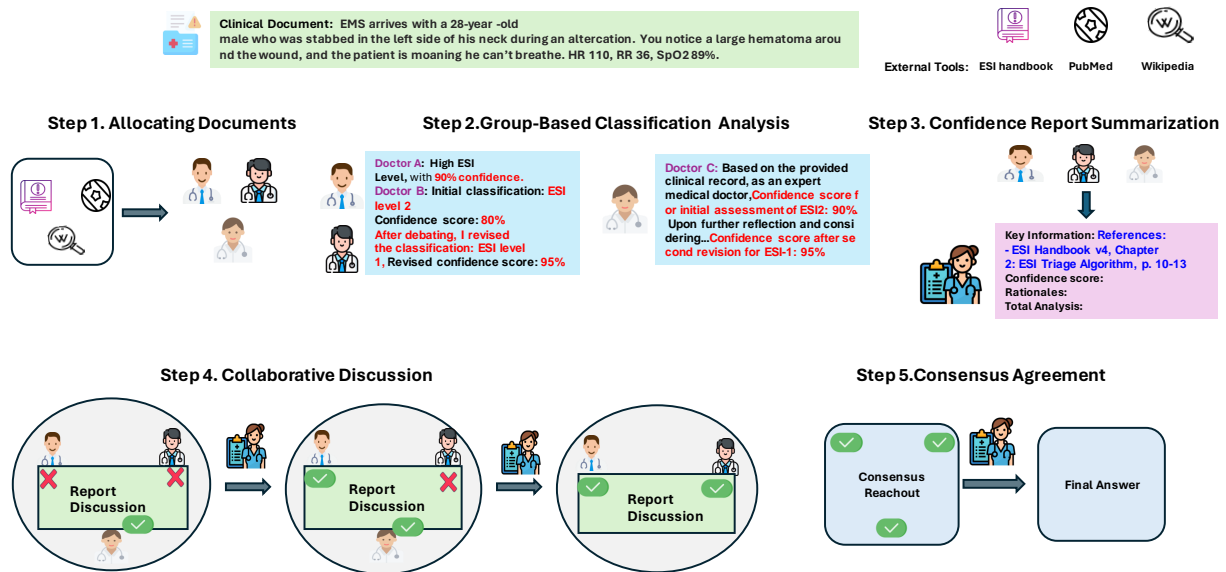


Figure 2: A diagram of our proposed heterogeneous multi-agent collaboration framework, TRIAGEAGENT, illustrates the reasoning process through five stages when given a clinical document as input. The five stages include 1) allocating documents, 2) group-based classification analysis, 3) confidence report summarization, 4) collaborative discussion, and 5) consensus agreement.

cal insights and refining models through instruction tuning. For instance, GeneGPT (Jin et al., 2023) leverages Web APIs from the National Center for Biotechnology Information (NCBI) to access diverse biomedical information and then employs GPT models for reasoning tasks. Additionally, the methods in (Zhang et al., 2024b; Singhal et al., 2022; Oniani et al., 2024; Kang et al., 2023) utilize instruction tuning combined with prompt design to adapt LLMs for various healthcare tasks, including decision support, medical question answering, and disease diagnosis.

## 2.2 LLM Applications in Auto-Agents

Research in both academia and industry has focused on autonomous agents trained in isolated, self-contained environments with limited knowledge bases (Wang et al., 2024a). Significant advancements have been achieved in deploying LLM-based agents capable of independently sensing and decision-making, as detailed in (Yao et al., 2023; Xie et al., 2023; Zhou et al., 2023). The trend has shifted towards collaborative multi-agent systems, which enhance the capabilities of LLM agents through iterative feedback and teamwork, as discussed in (Xi et al., 2023; Wang et al., 2024b; Li et al., 2023). These systems emulate human learning and decision-making processes, involving agents assuming specific roles (Wang et al., 2024b; Hong et al., 2023) and engaging in effective communication (Qian et al., 2023; Wu et al., 2023; Li

et al., 2023). Recent research has also explored improving agent performance through adversarial tactics such as debates (Du et al., 2023; Liang et al., 2023; Xiong et al., 2023) and negotiations (Fu et al., 2023), with innovative frameworks where agents interact competitively (Liang et al., 2023) or negotiate roles as buyers and sellers (Fu et al., 2023).

## 3 Methodology

This section presents the TRIAGEAGENT framework, which simulates teamwork and problem-solving in diagnosing and planning treatment for ED cases. We explore the heterogeneous structure formulation of the framework, as introduced in Appendix A. We describe the five stages of TRIAGEAGENT’s operation stages, illustrated in Figure 2 and detailed in the following procedures: First, **allocating documents**, where patient clinical records are assigned to expert agents to initiate ESI discussions. Second, **group-based classification analysis**, where agents are divided into two groups to conduct coarse and fine-grained classification. Third, **confidence report summarization**, where a summarized report is generated based on previous analyses, including classification results and confidence scores. Fourth, **collaborative discussion**, where agents engage in discussions over the summarized report, iteratively refining it according to key information, rationales, and confidence scores. Fifth, **consensus agreement**, resulting in an ultimate, precise, and thoroughly validated revised

report, highlighting the importance of collaborative decision-making.

### 3.1 Document Allocation

Given a patient’s narrative clinical note  $P = \{P_1, P_2, \dots\}$ , a clinical natural language query  $q$ , and a structured ESI handbook with level references  $R = \{R_1, R_2, \dots\}$ , this stage assigns documents to agents and initiates expert discussions, as illustrated in Figure 2 Step 1. For more details about query, see Appendix B.

### 3.2 Group-Based Document Classification

The classification stage is summarized as a function  $f : (P, R, q) \rightarrow C$ , where  $C$  represents the set of hierarchically structured ESI-level categories. To improve classification efficiency and accuracy, agents are divided into two groups, employing a direct and coarse-to-fine-grained classification, respectively.

**Coarse-to-Fine-Grained Classification** This group comprises two agents. The first agent,  $A_1$ , performs an initial coarse classification of the patient record  $P_i$  into two broad categories: high (levels 1, 2, and 3) or low (levels 3, 4, and 5), with level 3 included in both. This can be represented as  $A_1 : P_i \rightarrow \{C_{\text{high}}, C_{\text{low}}\}$ . The second agent,  $A_2$ , then refines these broad categories into detailed ESI levels: if classified as  $C_{\text{high}}$ , the second agent selects from  $\{1, 2, 3\}$ ; if  $C_{\text{low}}$ , it selects from  $\{3, 4, 5\}$ . This can be represented as:  $A_2 : \{C_{\text{high}}, C_{\text{low}}\} \rightarrow \{\{1, 2, 3\}, \{3, 4, 5\}\}$ . This two-step process reduces misclassifications and enhances precision. As illustrated in the left of step 2 in Figure 2, the first agent,  $A_1$  (referred to as Doctor A), initially assigns a high ESI level with 90% confidence.  $A_2$  (referred to as Doctor B) then refines the high-level category to ESI level 2 with an 80% confidence score. After rounds of discussion, Doctor B revises the classification to ESI level 1, achieving a revised confidence score of 95%.

**Direct Fine-Grained Classification** This group consists of a single agent,  $A_3$ , who directly assigns ESI levels:  $A_3 : P_i \rightarrow \{1, 2, 3, 4, 5\}$ . As illustrated in the right of step 2 in Figure 2,  $A_3$  (referred to as Doctor C) initially assigns ESI level 2 with 90% confidence and later refines the classification to level 1 with 95% confidence after further consideration.

By combining the two strategies described above, TRIAGEAGENT effectively addresses the chal-

lenges of distinguishing boundaries between ESI levels and enhances decision-making efficiency and accuracy through this collaborative effort.

### 3.3 Confidence Report Summarization

In this stage, the summarizer agent  $A_s$  summarizes previous document classification results, including confidence scores, rationales, and supporting evidence from  $(A_1, A_2, A_3)$ . This step consolidates the findings and uses each agent’s analysis report to construct the summary prompt  $Prompt_{rs}$ , ensuring a well-supported and reliable decision. Additionally, debates among the agents are incorporated into the summarizer’s prompt. The summarizer then generates a synthesized report by extracting key information and analyzing the previous classifications provided by the agents. This process can be mathematically formulated as:  $Repo = LLM(P, R, r_{rs}, Prompt_{rs})$ , where  $Repo$  represents the synthesized report,  $P$  denotes the patient’s clinical notes,  $R$  refers to the ESI handbook references,  $r_{rs}$  is the role of the summarizer, and  $Prompt_{rs}$  is the guideline prompt for the summarizer, including analysis reports from  $(A_1, A_2, A_3)$ . The synthesized report is structured as follows:  $Repo = [key\ information; confidence\ score; rationale; consolidated\ analysis]$ . This approach effectively combines insights from multiple agents, ensuring that the triage decision is based on comprehensive and validated information. Consequently, this method enhances the accuracy interpretability and reliability of the clinical triage process. As depicted in Step 3 of Figure 2, the key information includes references from the ESI Handbook (e.g., ESI Handbook v4, Chapter 2: ESI Triage Algorithm, p. 10-13). The summarized report consolidates rationales and total analysis, ensuring that all relevant information is considered.

### 3.4 Collaborative Discussion

In this stage, agents engage in multiple rounds of discussions based on the synthesized summary report  $Repo$  to refine their individual classifications. Unlike the commonly-used voting mechanism (Tang et al., 2024), TRIAGEAGENT critically reflects on the classification results, reasoning, and confidence scores, incorporating peer-provided evidence. Each agent  $A_i$  starts with an initial classification result  $C_i$  and confidence score  $S_i$  follows the following process:  $A_i$  reviews the classification results  $C_j$ , reasoning  $R_j$ , and confidence scores  $S_j$  from every other agent  $A_j (j \neq i)$ . If  $A_i$  is



persuaded by  $A_j$ 's reasoning or finds  $S_j > S_i$ , it updates  $C_i$  to  $C_j$  with explanations. Conversely, if  $A_i$  rejects  $A_j$ 's reasoning or finds  $S_j$  lower or equal to  $S_i$ , it justifies keeping  $C_i$ . This iterative process continues until agents reach a preliminary consensus or the early stopping mechanism is triggered. In Step 4 of Figure 2, agents participate in collaborative discussions to resolve discrepancies and refine the report.

### 3.5 Early-stopping Mechanism

To enhance the efficiency of group chat discussions, we implement an early-stopping mechanism inspired by Byzantine Consensus theory (Castro and Liskov, 1999). This approach requires at least  $3p + 1$  agents to handle  $p$  faulty agents in a single communication round. Additionally, our termination mechanism draws inspiration from advancements in LLMs fine-tuned with Reinforcement Learning from Human Feedback (RLHF), allowing consensus after several debate rounds (Du et al., 2023; Ouyang et al., 2022). The mechanism terminates communication when agents consistently confirm their reasoning with high confidence, thereby reducing unnecessary computations. It operates under two conditions: the first is repetition of high-confidence answers by a single agent: if an agent repeatedly provides the same answer with high confidence, that agent triggers early-stopping and exits the group discussions. The second is repetition of high-confidence answers by multiple agents: if all agents consistently provide the same answer with high confidence, the dialogue is terminated. This dynamic, real-time stopping condition enhances the traditional theory's efficiency, ensuring efficient and accurate consensus in group discussions. These conditions collectively foster an adaptive termination criterion, prioritizing efficiency and accuracy in reaching conclusions (Yin et al., 2023). In our case, we apply the early-stopping mechanism to each round of discussion of the TRIAGEAGENT.

### 3.6 Consensus Agreement

Finally, agents reach a formal consensus by integrating the refined answers, reasoning, and confidence scores from the collaborative discussion stage. This stage ensures all agents agree on a single classification, leveraging their combined domain knowledge to validate the final decision. This collaborative process ensures the final decision is robust and well-supported by comprehensive analysis. As shown in Step 5 of Figure 2, the final con-

Table 1: Statistics of the clinical triage dataset

Dataset	Training	Test-1	Test-2	Test-3
# of Docs	218	72	72	72

sensus is reached and the definitive classification is provided after all agents agree on the outcome.

## 4 Experimental Setup

**Dataset** We construct a clinical triage dataset by collecting patient cases from the publicly available Emergency Severity Index (ESI) Handbook v4 (esi, Accessed: 2024-04-06). This dataset is designed to evaluate machine learning models and methods for categorizing ESI levels in medical documents. To our knowledge, this is the first publicly released clinical triage dataset that includes clinical notes and corresponding ESI levels necessary for triage tasks, serving as a benchmark for evaluating our framework's effectiveness. Since the patient cases are sourced from the official ESI Handbook, no de-identification is needed. The dataset is divided into a training set and three test sets (test-1, test-2, and test-3). The training set contains 218 cases with the following distribution across ESI levels: ESI-1 (14), ESI-2 (92), ESI-3 (65), ESI-4 (22), and ESI-5 (25). Each test set contains 72 cases, maintaining fixed proportions of ESI levels: ESI-1 (12), ESI-2 (20), ESI-3 (13), ESI-4 (12), and ESI-5 (15). The dataset statistics are provided in Table 1. For more details, see Appendix C.

**Implementation** We use GPT-3.5-Turbo (OpenAI, 2024) and GPT-4 (OpenAI et al., 2024) from OpenAI for as our base models for the zero-shot experiments. Our TRIAGEAGENT framework utilizes the publicly open-sourced Autogen framework from Microsoft<sup>1</sup>. The temperature is 0.9,  $top_k$  is 1.0, and the cache seed is 42. The maximum number of iterations is 12, and the frequency penalty is 0.1. For SCTR and SCONS, we perform 8 iterations with a temperature of 0.9.

**Baselines** The performance of the TRIAGEAGENT framework is evaluated against several state-of-the-art baselines, including methods that employ LLM-based planning, tool usage, and retrieval-augmented generation.

- Chain-of-thought (Kojima et al., 2023) integrates step-by-step reasoning into the prompt for LLMs. We implemented CoT on our clinical triage dataset as a baseline method for comparison.

<sup>1</sup><https://microsoft.github.io/autogen/>

Table 2: Performance comparison of TRIAGEAGENT with baseline methods on the clinical triage dataset. The reported performance (%) in this table are averaged values from the three test sets in the dataset (Table 1). The highest performance is highlighted in **bold**. A lower total discordance represents a higher performance. CoT denotes chain-of-thought prompting, SCTR denotes self-contrast prompting, SCons denotes self-consistency prompting, and EoT denotes the exchange-of-thought prompting method.

Supervision	Method	Total Discordance	UnderTriage	Significant UnderTriage	OverTriage	Significant OverTriage
<b>GPT-3.5</b>						
-Supervised	Vanilla	39.18	21.76	15.28	18.06	6.85
-Zero-shot	MedAgent (w/Handbook)	39.58	<b>5.56</b>	5.56	34.03	15.97
	CoT (1-Agnt)	41.40	16.70	12.50	24.70	8.33
	SCTR (1-Agnt)	39.35	11.57	8.33	27.78	11.58
	SCons (1-Agnt)	36.11	17.59	8.33	<b>15.74</b>	7.87
	EoT (4-Agnt)	36.81	11.57	6.94	15.28	9.72
	SCons (4-Agnt)	34.72	5.56	<b>4.63</b>	27.31	10.65
	SCons (4-Agnt) (w/Handbook)	31.02	7.41	7.41	23.61	6.94
	SCons (4-Agnt)+Confidence (w/Handbook)	32.87	6.02	6.02	26.85	10.65
	<b>TRIAGEAGENT (Vanilla)</b>	34.72	<b>5.56</b>	<b>4.63</b>	27.31	10.56
	<b>TRIAGEAGENT (w/Handbook)</b>	31.02	7.87	7.41	22.69	<b>5.56</b>
	<b>TRIAGEAGENT (Ours)</b>	<b>30.56</b>	6.94	6.48	24.54	9.72
<b>GPT-4</b>						
-Supervised	Vanilla	23.50	8.10	6.94	14.80	8.33
	Vanilla (w/Handbook)	22.68	9.50	5.70	<b>7.10</b>	<b>1.90</b>
-Zero-shot	MedAgent (w/Handbook)	30.56	4.17	3.24	25.93	18.52
	CoT (1-Agnt)	37.40	14.30	8.33	23.30	10.64
	EoT (4-Agnt)	29.86	9.03	5.56	20.83	12.50
	SCons (4-Agnt)	29.63	11.11	7.87	18.06	8.33
	SCons (4-Agnt) (w/Handbook)	23.61	5.09	3.70	18.52	9.26
	SCons (4-Agnt)+ Confidence (w/Handbook)	23.61	5.09	3.70	18.52	9.26
	<b>TRIAGEAGENT (Vanilla)</b>	29.63	11.11	7.87	18.06	8.33
	<b>TRIAGEAGENT (w/Handbook)</b>	23.61	5.09	3.70	18.52	9.26
	<b>TRIAGEAGENT (Ours)</b>	<b>18.98</b>	<b>2.30</b>	<b>2.80</b>	17.10	8.80
<b>Human Eval</b>	<b>Human Experts</b>	31.43	12.80	8.61	18.60	10.50

- Self-contrast (Zhang et al., 2024a) improves stability and accuracy by contrasting different solving perspectives and summarizing discrepancies. We applied self-contrast on our clinical triage dataset to analyze and reconcile conflicting classification results as a baseline for comparison.
- Self-consistency (Wang et al., 2023b) enhances zero-shot and few-shot CoT by generating predominant responses through multiple chain samplings. We implemented self-consistency on our clinical triage dataset to generate multiple response chains and select the most consistent answers as a baseline for comparison.
- Exchange-of-thought (Yin et al., 2023) enables cross-model communication and problem-solving integration. We implemented EoT on our clinical triage dataset to facilitate communication between various agents as a baseline for comparison.
- MedAgent (Tang et al., 2024) is a role-playing collaboration framework for medical scenarios

using LLMs. We utilized MedAgent’s multi-agent framework on our clinical triage dataset as a baseline for comparison.

- TRIAGEAGENT utilizes dynamically updated confidence scores from various reasoning perspectives supported by external evidence to enhance the performance. It selects the top  $K$  most confident answers for critical assessment. Agents employ retrieval-augmented generation to refine choices, ensuring minimal discordance and maximal coherence. This iterative process integrates multiple methodologies and specialized knowledge retrieval to improve decision accuracy and reliability.

**Evaluation Protocol** Our primary evaluation metric is the *total discordance rate*, which measures the percentage of incorrectly predicted queries. This metric is critical as it provides a comprehensive overview of the accuracy; a lower value indicates better performance. Additionally, we evaluate the *undertriage rate*, *overtriage rate*,

467 *significant undertriage rate*, and *significant over-*  
468 *triage rate* to assess the model’s performance in  
469 specific areas of clinical urgency categorization.  
470 While these metrics are important, the *total dis-*  
471 *cordance rate* remains the primary measure of ac-  
472 curacy. Detailed definitions of the five ESI levels  
473 (I-V) and the evaluation metrics are provided in  
474 Appendix D.

## 475 5 Results

### 476 5.1 Main Results

477 We evaluate the performance of TRIAGEAGENT  
478 by averaging the results across the three test sets.  
479 The backbone LLMs in our experiments include  
480 GPT-3.5, GPT-4, Llama-2, and Llama-3. Table 2  
481 presents the main results with a better performance  
482 from the GPT-3.5 and GPT-4 models. Detailed re-  
483 sults for Llama-2-7B and Llama-3-8B are provided  
484 in Appendix E.

485 **Performance Comparison with Baselines** The  
486 performance of our framework compared to state-  
487 of-the-art (SOTA) methods is presented in Ta-  
488 ble 2. The TRIAGEAGENT framework surpasses  
489 traditional prompt engineering methods, includ-  
490 ing *supervised learning*, *self-contrast*, and *self-*  
491 *consistency*, with improvements of 9.25%, 8.79%,  
492 and 5.55%, respectively. It also outperforms SOTA  
493 multi-agent frameworks in a zero-shot setting using  
494 the GPT-3.5-turbo model, exceeding *MedAgents*  
495 by 9.02% and *EoT* by 6.25%.

496 **Comparison with CoT Methods** Performance  
497 can decline when employing overly complex Chain  
498 of Thought (CoT) methods. Simply stacking  
499 prompts without a clear, logical sequence can re-  
500 sult in hallucinations—erroneous outputs caused  
501 by insufficient document comprehension and mis-  
502 understanding of medical terminologies. However,  
503 our approach, which integrates multi-agent role-  
504 playing with confidence assessments, effectively  
505 addresses these issues and demonstrates its poten-  
506 tial as a more robust method for applying LLMs in  
507 clinical triage.

508 **Comparison with Single-Agent Methods** Meth-  
509 ods such as *CoT*, *self-contrast* and *supervised learn-*  
510 *ing* lack crucial interactions among multiple LLMs.  
511 This absence inhibits these methods from adap-  
512 tively refining their responses, leading to subopti-  
513 mal performance in triage question-answering sce-  
514 narios. Consequently, their discordance rates aver-  
515 age around 38.95% on our clinical triage dataset

when using GPT-3.5 (as illustrated in Table 2), high-  
lighting the need for improvement.

### 5.2 Ablation Study

Our ablation study analyzes team optimization and  
external resource optimization. After determining  
the optimal structure, the TRIAGEAGENT frame-  
work simulates multi-role team collaboration, en-  
abling agents to acquire the necessary capabilities  
to effectively accomplish triage tasks.

**Team Optimization** Our heterogeneous frame-  
work employs multiple agents, each with a specific  
role, to achieve optimal outcomes. We explored var-  
ious configurations to enhance team performance.  
By adjusting the number of role-specific agents, we  
found that four agents provided the most optimized  
structure for performance on the ESI triage dataset.  
Figure 3(a) details this optimal configuration and  
demonstrates how adjusting agent roles and num-  
bers enhances overall system performance.

**External Resource Optimization** We evaluate  
the ESI Handbook, PubMed, and Wikipedia for su-  
pervising our model (Figure 3(c)). The ESI Hand-  
book proved to be the most effective, significantly  
enhancing model performance with its targeted clin-  
ical diagnostic guidelines. In contrast, PubMed  
and Wikipedia require processing extensive addi-  
tional data. We also analyze the optimal number of  
agents using the ESI Handbook to maximize task  
efficiency (Figure 3(b)). Increasing the number  
of agents improved model effectiveness without  
significantly changing overall accuracy, enhancing  
the contextual understanding and interpretability of  
medical texts while optimizing resource utilization.

### 5.3 Case Study

**Error Analysis** Based on our findings, we con-  
duct an expert evaluation to identify key limita-  
tions and common issues in our model. As shown  
in Figure 4, we categorize these errors into four  
major types. The first type of error is **lack of doc-**  
**ument understanding**. This type of error arises  
from insufficient medical knowledge or incorrect  
linking to ESI levels, leading to misjudgments of  
clinical severity. The second type of error is **mis-**  
**retrieval of domain knowledge**. Errors result from  
inaccurately retrieving irrelevant or mismatched in-  
formation will compromise triage accuracy. The  
third type of error is **confidence-based consistency**  
**errors**. This type of errors is caused by confi-  
dently providing contradictory responses or failing

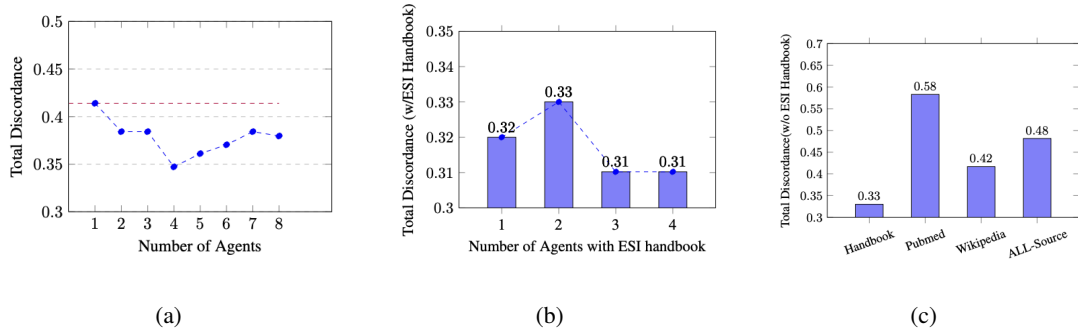


Figure 3: Ablation study results. (a), (b) and (c) show the impact of agent structure composition, optimizing the agent team with the ESI handbook and optimizing the agent team with various external resources, respectively. The lower *Total Discordance* value in this figure represents better performance.

Table 3: Time efficiency on multiple agents (average seconds per test case)

Model	Test-1	Test-2	Test-3
CoT (GPT-3.5)	0.17	0.19	0.17
CoT (GPT-4)	0.16	0.19	0.17
Self-Consistency (GPT-3.5)	0.21	0.23	0.24
Three-agents (GPT-3.5)	0.55	1.01	0.58
Three-agents (GPT-4)	0.57	1.00	0.59
Four-agents (GPT-3.5)	1.28	1.36	1.41
Four-agents (GPT-4)	1.31	1.56	1.53
Five-agents (GPT-3.5)	1.55	1.56	1.49
Six-agents (GPT-3.5)	2.11	2.05	2.13
TriagAgent (GPT-3.5)	1.30	1.45	1.52
TriagAgent (GPT-4)	1.31	1.43	1.50

Table 4: Cost efficiency on ESI datasets classification

Model	Performance	#API Calls
Two-agent (GPT-3.5)	38.42%	324
Three-agent (GPT-3.5)	38.42%	486
Four-agents (GPT-3.5)	30.56%	648
Four-agents (GPT-4)	18.98%	604
CoT (GPT-4)	37.40%	216

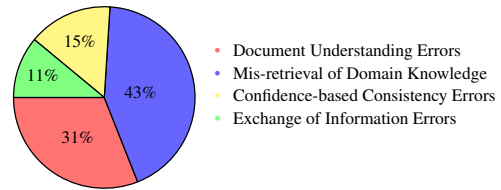


Figure 4: Ratio of different categories of error cases.

**Cost Analysis** Operational efficiency, particularly API token usage, is crucial for our framework. Table 4 compares the performance and costs of different configurations. While TRIAGEAGENT requires more API calls than a single-agent setup, feedback from clinical departments suggests these costs are reasonable. Our system improves efficiency by 12.54% compared to human experts manually classifying documents and achieves 18.54% higher performance than LLMs using CoT prompts operated by human experts.

to reach a consensus, often due to internal inefficiencies or flaws in the early stopping mechanism, resulting in falsely assured incorrect outcomes. The last type of errors is **exchange of information errors**. This type of errors results from incorrect data transfer between agents, disrupting logical sequences and leading to erroneous conclusions. See Appendix F for more details.

**Time Analysis** Time efficiency is crucial in emergency clinical triage. Table 3 compares the performance and time costs of different methods. Our multi-agent architecture incurs a slight but negligible increase in time costs compared to a single-agent model. Feedback from emergency department experts indicates that our method does not significantly affect overall time efficiency but greatly enhances decision accuracy, which is crucial for better patient outcomes and resource allocation.

## 6 Conclusions

This paper introduces a novel heterogeneous multi-agent framework, TRIAGEAGENT for clinical triage, utilizing LLM-based role-playing agents in a multi-stage group chat setting. This zero-shot, training-free, and interpretable framework comprises five significant stages. Experiments on clinical triage datasets demonstrate our framework significantly outperforms zero-shot baselines and experienced professionals. Case studies and human evaluations highlight areas for improvement, such as reducing document understanding errors and knowledge misretrieval. Future research can enhance the framework’s efficiency by improving document comprehension and correcting domain knowledge retrieval errors. Upon acceptance, we will release our dataset as open source.



611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658

## Ethics Statement

This research adhered to the highest ethical standards and best practices, as outlined in the ACL Code of Ethics. All experiments were conducted using publicly available datasets, thereby avoiding concerns related to privacy, confidentiality, or personal information. The datasets used are fully anonymized and have been vetted to ensure compliance with ethical guidelines. Additionally, we have carefully considered the broader impacts and potential applications of our work, ensuring that it does not inadvertently cause harm or misuse. Consequently, we believe this research is free from ethical issues.

## Limitations

In this paper, we introduce a heterogeneous multi-agent collaboration framework called TRIAGEAGENT. Despite our efforts, the framework faces limitations inherent to the healthcare industry.

**Limited Expert Evaluation** Our research is limited by the involvement of three human experts, restricting the scope and depth of expert evaluation, which may impact the generalizability of our findings to broader clinical settings.

**Cross-Institution Collaboration** The complexity of emergency departments requires advanced triage systems to address patient conditions within the same urgency levels. This need arises from varying institutional conditions, protocols, and patient demographics.

**Workflow Integration** Our model provides a final triage decision, but actual triage often involves multiple decision-making stages. Thus, it's crucial to evaluate how well our system integrates into existing workflows and complements human-led emergency care.

**Privacy** Deploying our framework necessitates strict privacy measures and clinical worker training. Processing clinical notes can expose sensitive information, making compliance with HIPAA and GDPR.

**Time and Cost Efficiency** TRIAGEAGENT aims to enhance clinical triage by automating initial patient assessment, reducing manual workload, waiting times, burnout and operational costs. However, using external technologies like OpenAI's API adds expenses and dependencies. A detailed cost-benefit

analysis is essential for sustainability and economic viability.

## References

Accessed: 2024-04-06. Emergency severity index (esi): A triage tool for emergency department. Available online. [https://media.emscimprovement.center/documents/ESI\\_Handbook2125.pdf](https://media.emscimprovement.center/documents/ESI_Handbook2125.pdf).

Zhijie Bao, Wei Chen, Shengze Xiao, Kuang Ren, Jiaao Wu, Cheng Zhong, Jiajie Peng, Xuanjing Huang, and Zhongyu Wei. 2023. [Disc-medllm: Bridging general large language models and real-world medical consultation](#).

Zhenyu Bi, Sajib Acharjee Dip, Daniel Hajialigol, Sindhura Kommu, Hanwen Liu, Meng Lu, and Xuan Wang. 2024. [Ai for biomedicine in the era of large language models](#).

Miguel Castro and Barbara Liskov. 1999. Practical byzantine fault tolerance. In *Proceedings of the Third Symposium on Operating Systems Design and Implementation*, pages 173–186, USA. USENIX Association.

G Chenais, E Lagarde, and C Gil-Jardiné. 2023. [Artificial intelligence in emergency medicine: Viewpoint of current applications and foreseeable opportunities and challenges](#). *J Med Internet Res*, 25:e40031.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. [Improving factuality and reasoning in language models through multiagent debate](#).

Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. [Improving language model negotiation with self-play and in-context learning from ai feedback](#).

659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713

714	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2023. <a href="#">Metagpt: Meta programming for a multi-agent collaborative framework.</a>	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach,	767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829
720	Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. 2023. <a href="#">Genegpt: Augmenting large language models with domain tools for improved access to biomedical information.</a>		
724	Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. 2023. <a href="#">Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks.</a>		
728	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. <a href="#">Large language models are zero-shot reasoners.</a>		
731	Yuan Li, Yixuan Zhang, and Lichao Sun. 2023. <a href="#">Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents.</a>		
735	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2023. <a href="#">Encouraging divergent thinking in large language models through multi-agent debate.</a>		
739	Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy Chowdhury, Yun Li, Hejie Cui, Xuchao Zhang, Tianjiao Zhao, Amit Panalkar, Dhagash Mehta, Stefano Pasquali, Wei Cheng, Haoyu Wang, Yanchi Liu, Zhengzhang Chen, Haifeng Chen, Chris White, Quanquan Gu, Jian Pei, Carl Yang, and Liang Zhao. 2024. <a href="#">Domain specialization as the key to make large language models disruptive: A comprehensive survey.</a>		
748	Zhengliang Liu, Zihao Wu, Mengxuan Hu, Bokai Zhao, Lin Zhao, Tianyi Zhang, Haixing Dai, Xianyan Chen, Ye Shen, Sheng Li, Brian Murray, Tianming Liu, and Andrea Sikora. 2023. <a href="#">Pharmacygpt: The ai pharmacist.</a>		
753	Danqing Ma, Meng Wang, Ao Xiang, Zongqing Qi, and Qin Yang. 2024. <a href="#">Transformer-based classification outcome prediction for multimodal stroke treatment.</a>		
756	Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. <a href="#">Capabilities of gpt-4 on medical challenge problems.</a>		
759	David Oniani, Xizhi Wu, Shyam Visweswaran, Sumit Kapoor, Shravan Kooragayalu, Katelyn Polanska, and Yanshan Wang. 2024. <a href="#">Enhancing large language models for clinical decision support by incorporating clinical practice guidelines.</a>		
764	OpenAI. 2024. <a href="#">Gpt-3.5 turbo documentation.</a> <a href="https://platform.openai.com/docs/models/gpt-3-5-turbo">https://platform.openai.com/docs/models/gpt-3-5-turbo</a> . Accessed: 2024-06-15.		

830	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Ra-	889
831	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	jkumar, Joelle Barral, Christopher Semturs, Alan	890
832	Girish Sastry, Heather Schmidt, David Schnurr, John	Karthikesalingam, and Vivek Natarajan. 2022. <a href="#">Large</a>	891
833	Schulman, Daniel Selsam, Kyla Sheppard, Toki	<a href="#">language models encode clinical knowledge.</a>	892
834	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav		
835	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai,	893
836	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	and Chao Zhang. 2023. <a href="#">Adaplanner: Adaptive plan-</a>	894
837	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	<a href="#">ning from feedback with language models.</a>	895
838	lipe Petroski Such, Natalie Summers, Ilya Sutskever,		
839	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	Rocío Sánchez-Salmerón et al. 2022. <a href="#">Machine learning</a>	896
840	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	<a href="#">methods applied to triage in emergency services: A</a>	897
841	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	<a href="#">systematic review.</a> <i>International Emergency Nursing</i> ,	898
842	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	60:101109.	899
843	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,		
844	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming	900
845	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and	901
846	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	Mark Gerstein. 2024. <a href="#">Medagents: Large language</a>	902
847	Clemens Winter, Samuel Wolrich, Hannah Wong,	<a href="#">models as collaborators for zero-shot medical reason-</a>	903
848	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	<a href="#">ing.</a>	904
849	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-		
850	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	905
851	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	906
852	Zheng, Juntang Zhuang, William Zhuk, and Barret	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	907
853	Zoph. 2024. <a href="#">Gpt-4 technical report.</a>	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	908
		Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	909
854	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	910
855	roll L. Wainwright, Pamela Mishkin, Chong Zhang,	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	911
856	Sandhini Agarwal, Katarina Slama, Alex Ray, John	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	912
857	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	913
858	Maddie Simens, Amanda Askell, Peter Welinder,	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	914
859	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	915
860	<a href="#">Training language models to follow instructions with</a>	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	916
861	<a href="#">human feedback.</a>	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	917
		bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	918
862	Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	919
863	Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	920
864	Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023.	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	921
865	<a href="#">Communicative agents for software development.</a>	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	922
		Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	923
866	Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen,	Melanie Kambadur, Sharan Narang, Aurelien Rod-	924
867	Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang,	riguez, Robert Stojnic, Sergey Edunov, and Thomas	925
868	Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su,	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-</a>	926
869	Huadong Wang, Cheng Qian, Runchu Tian, Kunlun	<a href="#">tuned chat models.</a>	927
870	Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen		
871	Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi,	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	928
872	Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong,	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	929
873	Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan,	Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei,	930
874	Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng	and Ji-Rong Wen. 2024a. <a href="#">A survey on large language</a>	931
875	Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and	<a href="#">model based autonomous agents.</a>	932
876	Maosong Sun. 2023. <a href="#">Tool learning with foundation</a>		
877	<a href="#">models.</a>	Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu,	933
		Xifeng Yan, Jianfeng Gao, and Furu Wei. 2023a.	934
878	Chantal Shaib, Millicent L. Li, Sebastian Joseph, Iain J.	<a href="#">Augmenting language models with long-term mem-</a>	935
879	Marshall, Junyi Jessy Li, and Byron C. Wallace. 2023.	<a href="#">ory.</a>	936
880	<a href="#">Summarizing, simplifying, and synthesizing medical</a>	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc	937
881	<a href="#">evidence using gpt-3 (with varying success).</a>	Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery,	938
		and Denny Zhou. 2023b. <a href="#">Self-consistency improves</a>	939
882	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mah-	<a href="#">chain of thought reasoning in language models.</a>	940
883	davi, Jason Wei, Hyung Won Chung, Nathan Scales,		
884	Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl,	Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao	941
885	Perry Payne, Martin Seneviratne, Paul Gamble, Chris	Ge, Furu Wei, and Heng Ji. 2024b. <a href="#">Unleashing the</a>	942
886	Kelly, Nathaneal Scharli, Aakanksha Chowdhery,	<a href="#">emergent cognitive synergy in large language mod-</a>	943
887	Philip Mansfield, Blaise Agueria y Arcas, Dale Web-	<a href="#">els: A task-solving agent through multi-persona self-</a>	944
888	ster, Greg S. Corrado, Yossi Matias, Katherine Chou,	<a href="#">collaboration.</a>	945



946 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,  
947 Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,  
948 Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadal-  
949 lah, Ryen W White, Doug Burger, and Chi Wang.  
950 2023. [Autogen: Enabling next-gen llm applications  
951 via multi-agent conversation.](#)

952 Ziheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen  
953 Ding, Boyang Hong, Ming Zhang, Junzhe Wang,  
954 Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan,  
955 Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran  
956 Wang, Changhao Jiang, Yicheng Zou, Xiangyang  
957 Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng,  
958 Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan  
959 Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui.  
960 2023. [The rise and potential of large language model  
961 based agents: A survey.](#)

962 Tianbao Xie, Fan Zhou, Zhoujun Cheng, Peng Shi, Lu-  
963 oxuan Weng, Yitao Liu, Toh Jing Hua, Junning Zhao,  
964 Qian Liu, Che Liu, Leo Z. Liu, Yiheng Xu, Hongjin  
965 Su, Dongchan Shin, Caiming Xiong, and Tao Yu.  
966 2023. [Openagents: An open platform for language  
967 agents in the wild.](#)

968 Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing  
969 Qin. 2023. [Examining inter-consistency of large lan-  
970 guage models collaboration: An in-depth analysis via  
971 debate.](#)

972 LH Yao, KC Leung, CL Tsai, CH Huang, and LC Fu.  
973 2021. [A novel deep learning-based system for triage  
974 in the emergency department using electronic medi-  
975 cal records: Retrospective cohort study. \*J Med Inter-  
976 net Res\*, 23\(12\):e27008.](#)

977 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak  
978 Shafraan, Karthik Narasimhan, and Yuan Cao. 2023.  
979 [React: Synergizing reasoning and acting in language  
980 models.](#)

981 Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng  
982 Guo, Junqi Dai, Xuanjing Huang, and Xipeng Qiu.  
983 2023. [Exchange-of-thought: Enhancing large lan-  
984 guage model capabilities through cross-model com-  
985 munication.](#)

986 Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying  
987 Peng, Jun Wang, Yueting Zhuang, and Weiming Lu.  
988 2024a. [Self-contrast: Better reflection through incon-  
989 sistent solving perspectives.](#)

990 Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang  
991 Chen, Zekun Li, and Linda Ruth Petzold. 2024b.  
992 [Alpacare:instruction-tuned large language models for  
993 medical application.](#)

994 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou,  
995 Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue  
996 Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Gra-  
997 ham Neubig. 2023. [Webarena: A realistic web envi-  
998 ronment for building autonomous agents.](#)

999 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun,  
1000 and Chao Zhang. 2023. [Toolqa: A dataset for llm  
1001 question answering with external tools.](#)



## A Model Architecture Comparison with Previous Work

We compare TRIAGEAGENT with representative previous works, as shown in Figure 5. We refer to our system as a **heterogeneous multi-agent** framework because each agent represents a different role, creating a role-playing heterogeneous structure. This structure simulates real-world collaboration among different roles, with each agent (or role) having its unique function and task, thereby improving the overall system’s decision-making efficiency and accuracy. The changing colors of nodes in the figure illustrate our dynamic heterogeneous nature, representing agents’ evolving perspectives based on different rounds of discussions. In the second row, nodes represent agents at different time steps, arrows indicate the edges, and colors signify the roles of the agents.

**Node** A node represents an agent at a specific time step, each with a unique role in a heterogeneous role-playing structure. This setup simulates real-world collaboration, enhancing decision-making efficiency and accuracy by assigning distinct functions and tasks to each agent. Additionally, the changing node colors in the figure highlight our framework’s dynamic nature, illustrating how agents’ perspectives evolve as they assimilate and process new information through subsequent discussion rounds.

**Edge** Edges represent the communication channels between nodes during multi-agent collaboration, illustrating how information flows between agents and through the system. In our LLM-agent-based feed-forward network, these edges show how agents share information across different time steps to generate the final answer for the task query.

## B Query Details

below:

query  $q$ : What is the ESI level of the following clinic record? Please give me a final unique answer after a second revision of your first proposed answer. You can learn from the Emergency Severity Index handbook v4. Double-check the ESI handbook and ask yourself again(two-round self-check)when you are sure about this ESI level classification before you give me the answer. Then classify the following medical record according to ESI level, candidate answers are ESI-1,2,3,4,5.

Here is the record:

## C Dataset Construction

The dataset was constructed by extracting relevant patient cases from the ESI handbook v4, focusing on a comprehensive range of clinical scenarios. Each case was carefully reviewed and labeled by professional human experts to ensure accuracy. Each case was carefully reviewed and labeled by professional human experts to ensure accuracy. The dataset is divided into a training set and three test sets (test-1, test-2, and test-3), with the training set containing 180 cases and each test set containing 72 cases. We maintained fixed proportions of ESI levels in the test sets as follows: ESI-1 (12), ESI-2 (20), ESI-3 (13), ESI-4 (12), and ESI-5 (15).

For the training scenarios, we have a total of 218 cases with the following distribution across ESI levels: ESI-1 (14), ESI-2 (92), ESI-3 (65), ESI-4 (22), and ESI-5 (25). The explanations in the training dataset are manually annotated by human experts. Detailed proportions of each ESI level in the training and test sets are provided in Table 1.

Each training and test set includes cases with detailed clinical notes and corresponding ESI levels. Although the test sets do not contain explicit explanations for each label, they have been accurately annotated by human experts based on the clinical information provided. The lack of explanations in the test set is due to the scarcity of medical data and the high cost of manual labeling. This highlights the advantage of our framework in reducing human labor costs and improving the efficiency of medical text classification.

### C.1 Recruitment and Payment

The human experts assisting us in the classification of medical documents are professional experts who voluntarily participated in our comparative study as collaborators. We did not provide them with any additional payments or benefits.

### C.2 Instructions Given to Participants

We invited professional human experts to serve as human annotators. The full text of the instructions given to participants is: "This study aims to evaluate the effectiveness of a generative AI model in predicting ESI levels and compare its performance with traditional nurse triage. Your participation will help us understand the potential of AI in augment-

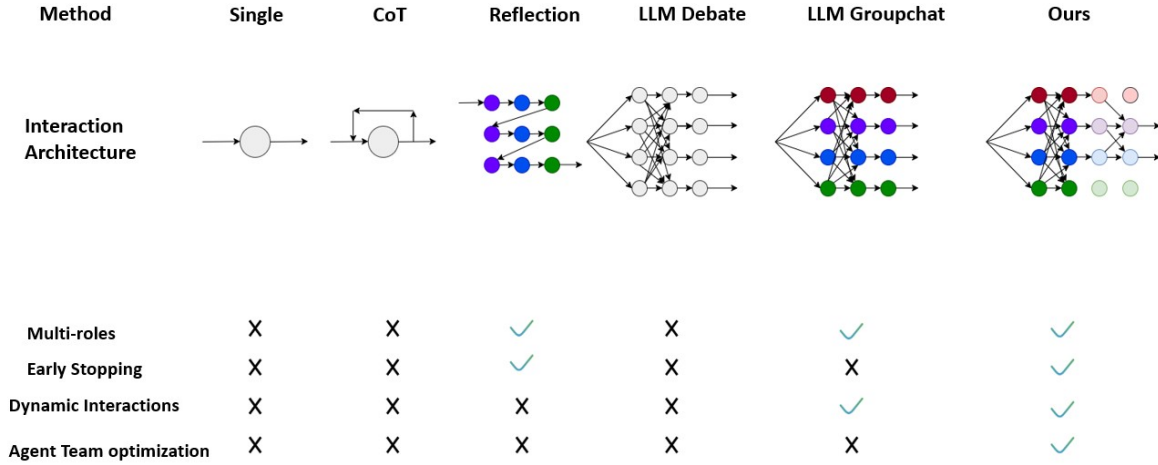


Figure 5: Topology Structure Comparisons of previous methods

ing emergency department workflows and improving patient care. All test results are de-identified and will only be used for the purposes of this research study. Please do not look up answers or use any additional resources to complete the test as that can negatively impact the validity of this study."

## D ESI Hierarchy and Evaluation Metrics

### D.1 Hierarchy of ESI levels

The ESI (Emergency Severity Index) levels classify medical events by urgency, from 1 (most urgent) to 5 (least urgent). Medical personnel can refer to the medical event in the ESI handbook with their medical experience to quickly categorize the current emergency medical event according to the ESI classification.

- ESI-1: Most Urgent; Immediate life-saving intervention
- ESI-2: High urgency; Potentially life-threatening, prompt attention necessary.
- ESI-3: Urgent; Requires multiple resources but not immediately life-threatening.
- ESI-4: Less urgent; Requires one resource; not immediately life-threatening.
- ESI-5: Least Urgent; No resources needed immediately; wait time is acceptable.

### D.2 Evaluation Metrics

The *Total discordance* is calculated as the ratio of the total number of misclassified texts to the total number of texts, representing the overall error rate of the model. This metric is given by the formula:

$$\text{Total discordance} = \frac{\text{Total Misclassifications}}{\text{Total number of texts}} \quad (1)$$

where *Total misclassifications* is the number of queries incorrectly classified by the model, and *Total number of texts* is the total number of queries analyzed. The *Undertriage* rate is defined as the fraction of instances where the predicted label is greater than the true label, which is calculated as:

$$\text{Undertriage} = \frac{\text{Number of predictions True\_labels}}{\text{Total number of texts}} \quad (2)$$

Similarly, the *Overtriage* rate is defined as the fraction of instances where the predicted label is less than the true label:

$$\text{Overtriage} = \frac{\text{Number of predictions True\_labels}}{\text{Total number of texts}} \quad (3)$$

The *Significant Undertriage* rate captures the scenarios where the true label is critical (ESI 1 or 2) but the prediction underestimates the urgency (predicted as 3, 4, or 5):

$$\text{Significant Undertriage} = \frac{\text{Predicted-3, 4, or 5}}{\text{Total number of texts}} \quad (4)$$

Lastly, the *Significant Overtriage* rate addresses the instances where less critical labels (true label

Method	Zero-shot(CoT)
<b>Total discordance</b>	<b>56.25%</b>
<b>Undertriage</b>	3.48%
<b>Significant undertriage</b>	3.48%
<b>Overtriage</b>	52.78%
<b>Significant overtriage</b>	35.42%
Method	Fine-Tuned
<b>Total discordance</b>	<b>27.78%</b>
<b>Undertriage</b>	10.65%
<b>Significant undertriage</b>	7.64%
<b>Overtriage</b>	17.13%
<b>Significant overtriage</b>	6.94%

Table 5: Performance of LLaMa-3-8B on our clinical triage task.

is 2, 3, or 4) are overestimated by the model (predicted as 1):

$$\text{Significant Overtriage} = \frac{\text{Predicted-3,4,or 5}}{\text{Total number of texts}} \quad (5)$$

## E Performance of Different Backbone LLMs

**Using Different Language Models** In our experiments, we utilized GPT-3.5-turbo<sup>2</sup>, GPT-4<sup>3</sup>, their combination, and LLaMa-2-7B<sup>4</sup> and Llama-3-8B<sup>5</sup>. When GPT-3.5-turbo and GPT-4 engaged in interactive debates on AI-Triage tasks, we observed performance improvements of 31.88% for GPT-3.5-turbo and 24.6% for GPT-4. In a zero-shot scenario, LLaMa-2-7B performed 17.025% worse than LLaMa-3-8B. We also fine-tuned LLaMa-3-8B using training scenarios from the ESI handbook to enhance its domain knowledge. Compared to GPT-4, which achieved a performance improvement of 23.50% on the same training dataset, LLaMa-3-8B showed a performance improvement of 27.78%, making it 4.28% worse than GPT-4. Detailed results for LLaMa-3 are presented in Table 5, demonstrating that its performance on clinical triage tasks was inferior to that of the GPT models. For details about the performance of LLaMa-3 on our task, see Table 5

## F Case Study

To better illustrate the issue, we have excerpted parts of the complete conversation as key examples

<sup>2</sup><https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates/>

<sup>3</sup><https://openai.com/index/gpt-4/>

<sup>4</sup><https://llama.meta.com/llama2/>

<sup>5</sup><https://llama.meta.com/llama3/>

to more effectively explain the corresponding types of errors.

### •Type-1: Lack of Document Understanding

*Admin (to chat\_manager):*

What is the ESI level of the following clinic record? Please give me a final unique answer after a second revision of your first proposed answer. The Candidate answers: ESI-1,2,3,4,5. Double-check the clinic notes and related medical knowledge and ask yourself again (two-round self-check) when you are sure about this ESL level classification before you give me the answer. Here is the record:“A healthy 10-year-old male is brought to the emergency department by his mother, who reports that her son has not moved his bowels for a week. He is complaining of 7/10 generalized abdominal pain, nausea, and lack of appetite. Vital signs: BP 107/66, HR 75, RR 20, T 98.6°F, SpO2 99%”.

*Agent\_1 (to chat\_manager):*

Based on the clinical record provided, the most appropriate Emergency Severity Index (ESI) level for this patient would be ESI level 3. This decision is based on the patient’s presenting complaints of abdominal pain, nausea, and lack of appetite, along with the vital signs indicating stable blood pressure, heart rate, respiratory rate, temperature, and oxygen saturation. While the patient’s symptoms are concerning and require evaluation, there are no immediate life-threatening issues or unstable vital signs that would warrant a higher ESI level.

Confidence score: 70%.

*Agent\_2 (to chat\_manager):*

Based on a thorough review of the clinical record and consideration of the patient’s presentation, vital signs, and symptoms, I would revise the ESI level for this patient to ESI level 2. The patient’s complaints of abdominal pain, nausea, lack of appetite, and constipation for a week warrant urgent evaluation, but the stable vital signs indicate that there are no immediate life-threatening issues. However, the patient’s symptoms require timely assessment and treatment, justifying a higher acuity level than initially considered.

Confidence score: 85%.

*Agent\_3 (to chat\_manager):*

1227	Based on the information provided in the clinical	55-year-old female tells you. No previous medical	1278
1228	record, a more appropriate Emergency Severity	history and vital signs are within normal limits.	1279
1229	Index (ESI) level for this patient would be ESI		1280
1230	level 2. The patient is a 10-year-old male with	<hr/>	1281
1231	concerning symptoms of abdominal pain, nausea,	<i>Agent_1 (to chat_manager):</i>	1282
1232	lack of appetite, and constipation for a week.	Based on the information provided and the	1283
1233	While the vital signs are stable, the duration of	Emergency Severity Index handbook, this patient's	1284
1234	symptoms and the need for urgent evaluation	clinic note would be categorized as Low ESI level.	1285
1235	justify a higher acuity level. This patient may	The patient's condition is not life-threatening but	1286
1236	require timely intervention and assessment to	requires medical intervention, which aligns with a	1287
1237	address the underlying cause of his symptoms,	lower ESI level.	1288
1238	making ESI level 2 more appropriate in this case.	References: - ESI Handbook v4, Chapter 2: ESI	1289
1239	Confidence score: 90%.	Triage Algorithm, p.10-13.	1290
1240		<hr/>	1291
1241	<i>Result_summarizer (to chat_manager):</i>	<i>Agent_2 (to chat_manager):</i>	1292
1242	After synthesizing the inputs from the Agent_1,	The scenario describes a 55-year-old female who	1293
1243	Agent_2, and Agent_3, the final unique answer for	has a foreign object (a bug) in her ear, causing	1294
1244	the ESI level of the patient in the clinical record	discomfort and a buzzing noise. She has no	1295
1245	would be ESI level 2. This decision is based on the	previous medical history and her vital signs are	1296
1246	patient being a 10-year-old male presenting with	within normal limits. This is a non-life-threatening	1297
1247	concerning symptoms of abdominal pain, nausea,	condition, but it does require medical intervention	1298
1248	lack of appetite, and constipation for a week.	to remove the foreign object and alleviate the	1299
1249	While the vital signs are stable, the duration and	patient's discomfort.	1300
1250	nature of the symptoms warrant urgent evaluation	The Emergency Severity Index (ESI) is a five-level	1301
1251	and potential intervention, justifying a higher	emergency department triage algorithm that	1302
1252	acuity level.	categorizes patients by both acuity and resource	1303
1253	Confidence score: 92%.	needs. ESI level 1 is for patients who need	1304
1254		immediate life-saving intervention, while ESI level	1305
1255		5 is for patients with non-urgent conditions.	1306
1256	<b>•Type-2: Mis-retrieval of Domain Knowledge</b>	Given the patient's condition, she would likely be	1307
1257		classified as ESI level 4. This level is for patients	1308
1258	This example mis-retrieved the domain knowl-	with conditions that could potentially escalate to a	1309
1259	edge from <i>ESI Handbook v4, Chapter 2: ESI</i>	more serious problem requiring emergency inter-	1310
1260	<i>Triage Algorithm, p. 10-13.</i> Then this example	vention if left untreated, but are not life-threatening.	1311
1261	was erroneously categorized as needing urgent	The patient's condition is not life-threatening, but	1312
1262	intervention.	it is causing her significant discomfort and could	1313
1263		potentially lead to complications such as infection	1314
1264	<i>User_proxy (to chat_manager):</i>	if not addressed.	1315
1265	What is the ESI level of the following clinic	The most alarming aspect of this issue is the	1316
1266	record according to the Emergency Severity	potential for complications if the foreign object is	1317
1267	Index handbook? Please give me a final unique	not removed promptly. The use of a Q-tip could	1318
1268	answer after a second revision from your first	potentially push the bug further into the ear canal,	1319
1269	proposed answer. The Candidate answers:	causing more damage. It's also noteworthy that	1320
1270	ESI-1,2,3,4,5. Double Check the clinic notes	the patient's vital signs are within normal limits,	1321
1271	and related medical knowledge and ask yourself	indicating that she is not in immediate danger.	1322
1272	again (two-round self-check) when you are sure	After a second revision, I still believe that the ESI	1323
1273	about this ESL level classification before you	level for this patient would be ESI level 4. The	1324
1274	give me the answer. Here is the record:" This	patient's condition is not life-threatening, but it	1325
1275	sounds really strange. A bug flew into my right	does require medical intervention.	1326
1276	ear while I was gardening. I tried to get it out	Confidence score: 90%.	1327
1277	by using a Q-tip. I just don't know what else to	<hr/>	1328
	do, but this buzzing noise is driving me crazy," a	<i>Agent_3 (to chat_manager):</i>	



1329	After reviewing the patient’s condition and the	wrist,” reports a 58-year-old female with a history	1380
1330	Emergency Severity Index handbook, I would	of migraines. There is no obvious deformity. Vital	1381
1331	classify this patient as ESI level 4 with a confi-	signs are within normal limits, and she rates her	1382
1332	dence score of 90%. The patient’s condition is not	pain as 5/10.	1383
1333	life-threatening but requires medical intervention		1384
1334	to remove the foreign object from her ear and alle-	<hr/>	1385
1335	viate her discomfort. This aligns with the criteria	<i>Response 1 by chat_manager (initial response):</i>	1386
1336	for ESI level 4, which is for conditions that could	Based on the patient’s description and clinical	1387
1337	potentially escalate to a more serious problem if	information provided, I recommend classifying	1388
1338	left untreated but are not life-threatening.	this patient as ESI level 4. The injury mechanism	1389
1339		and the pain level, along with normal vital signs,	1390
1340	<hr/>	suggest a non-life-threatening condition.	1391
1341	<i>Agent_2 (to chat_manager):</i>	Confidence score: 90%.	1392
1342	Based on the patient’s clinical record and the		1393
1343	Emergency Severity Index handbook, I would	<hr/>	1394
1344	classify this patient as ESI level 4 with a confi-	<i>Response 2 by chat_manager (second response</i>	1395
1345	dence score of 90%. The patient’s condition is not	<i>after self-check):</i>	1396
1346	life-threatening but requires medical intervention	Upon further reflection and review of the medical	1397
1347	to remove the foreign object from her ear and alle-	knowledge, I now assess the patient’s condition	1398
1348	viate her discomfort. This aligns with the criteria	as ESI level 3. Considering the possibility of a	1399
1349	for ESI level 4, which is for conditions that could	hidden injury due to the nature of the fall, a more	1400
1350	potentially escalate to a more serious problem if	cautious approach is warranted.	1401
1351	left untreated but are not life-threatening.	Confidence score: 92%.	1402
1352			1403
1353	<hr/>	<b>•Type-4: Exchange of Information Errors</b>	1404
1354	<i>Result_summarizer (to chat_manager):</i>		1405
1355	After synthesizing the inputs from the Agent_1,	<hr/>	1406
1356	Agent_2, and Agent_3, the consensus is that	<i>User_proxy (to chat_manager):</i>	1407
1357	the patient should be classified as ESI level	What is the ESI level of the following clinical	1408
1358	4. This level is appropriate for patients with	record according to the Emergency Severity In-	1409
1359	conditions that could potentially escalate to a	dex handbook? Please provide a final answer after	1410
1360	more serious problem if left untreated but are not	a second revision from your initial proposed	1411
1361	life-threatening. The patient’s condition of having	answer. The candidate answers are: ESI-1, 2, 3, 4, 5.	1412
1362	a bug in her ear is causing significant discomfort	Double-check the clinic notes and related medical	1413
1363	and requires medical intervention, but it is not	knowledge before finalizing your decision. Here is	1414
1364	life-threatening. Therefore, the final decision is	the record: “A bug flew into my right ear while I	1415
1365	ESI level 4.	was gardening. I tried to get it out by using a Q-tip.	1416
1366		I just don’t know what else to do, but this buzzing	1417
1367	<b>•Type-3: Confidence-based Consistency Er-</b>	noise is driving me crazy,” a 55-year-old female	1418
1368	<b>rors</b>	tells you. No previous medical history and vital	1419
1369	This error might be caused by the early-stopping	signs are within normal limits.	1420
1370	mechanism and confidence-score passing process.		1421
1371	<hr/>	<i>Agent_1 (to chat_manager):</i>	1422
1372	<i>Admin (to chat_manager):</i>	Based on the information provided and the	1423
1373	What is the ESI level of the following clinic	Emergency Severity Index handbook, this patient’s	1424
1374	record? Please give me a final unique answer	clinic note would be categorized as a low ESI level.	1425
1375	after a second revision from your first proposed	The patient’s condition is not life-threatening	1426
1376	answer. The Candidate answers: ESI-1,2,3,4,5.	but requires medical intervention, aligning with	1427
1377	Double Check the clinic notes and related medical	a lower ESI level.(e.g., ESI level 5), with a	1428
1378	knowledge and ask yourself again(two-round	confidence score of 80%.	1429
1379	self-check)when you are sure about this ESL level		1430
	classification before you give me the answer.Here	<hr/>	
	is the record: “I slipped on the ice, and I hurt my	<i>Agent_2 (to chat_manager):</i>	
		The scenario describes a 55-year-old female	
		with a foreign object (a bug) in her ear, causing	

1431 discomfort and a buzzing noise. She has no  
1432 previous medical history, and her vital signs are  
1433 within normal limits. This is a non-life-threatening  
1434 condition but requires medical intervention to  
1435 remove the foreign object and alleviate the  
1436 patient's discomfort.

1437 Given the patient's condition, she would likely be  
1438 classified as ESI level 4. This level is for patients  
1439 with conditions that could potentially escalate to  
1440 a more serious problem if left untreated but are  
1441 not life-threatening. After a second revision, I still  
1442 believe that the ESI level for this patient is ESI  
1443 level 4, with a confidence score of 90%.

---

1445 *Agent\_3 (to chat\_manager):*

1446 After reviewing the patient's condition and the  
1447 Emergency Severity Index handbook, I would  
1448 classify this patient as ESI level 4 with a confi-  
1449 dence score of 90%. The patient's condition is not  
1450 life-threatening but requires medical intervention  
1451 to remove the foreign object from her ear and  
1452 alleviate her discomfort. This aligns with the  
1453 criteria for ESI level 4.

---

1455 *Result\_summarizer (to chat\_manager):*

1456 After synthesizing the inputs from the Agent\_1,  
1457 Agent\_2, and Agent\_3, the consensus is that  
1458 the patient should be classified as ESI level  
1459 4. This level is appropriate for patients with  
1460 conditions that could potentially escalate to a  
1461 more serious problem if left untreated but are not  
1462 life-threatening. The patient's condition of having  
1463 a bug in her ear is causing significant discomfort  
1464 and requires medical intervention, but it is not  
1465 life-threatening. Therefore, the final decision is  
1466 ESI level 4, with a consensus confidence score of  
1467 90%.

---

1469 Explanation: In this example, the inconsistent  
1470 classification by Agent\_1 ("low ESI level" as  
1471 ESI level 5) compared to the other agents (ESI  
1472 level 4) could lead to confusion and incorrect  
1473 final decisions, representing an Exchange of  
1474 Information Error.