# LLM-based ambiguity detection in natural language instructions for collaborative surgical robots

Ana Davila[1] and Jacinto Colan[2] and Yasuhisa Hasegawa[1]

*Abstract*— Ambiguity in natural language instructions poses significant risks in safety-critical human-robot interaction, particularly in domains such as surgery. To address this, we propose a framework that uses Large Language Models (LLMs) for ambiguity detection specifically designed for collaborative surgical scenarios. Our method employs an ensemble of LLM evaluators, each configured with distinct prompting techniques to identify linguistic, contextual, procedural, and critical ambiguities. A chain-of-thought evaluator is included to systematically analyze instruction structure for potential issues. Individual evaluator assessments are synthesized through conformal prediction, which yields non-conformity scores based on comparison to a labeled calibration dataset. Evaluating Llama 3.2 11B and Gemma 3 12B, we observed classification accuracy exceeding 60% in differentiating ambiguous from unambiguous surgical instructions. Our approach improves the safety and reliability of human-robot collaboration in surgery by offering a mechanism to identify potentially ambiguous instructions before robot action.

Fig. 1. Collaborative robot-assisted surgery requires seamless communication between surgeons and robotic assistants.

## I. INTRODUCTION

The integration of robotic systems into high-stakes environments, particularly surgical settings, has accelerated the need for reliable human-robot communication. Natural language instructions offer an intuitive interface between healthcare professionals and robotic assistants, but the inherent ambiguity of natural language presents significant challenges. In surgical contexts, where precision is essential, instruction ambiguity can lead to potentially dangerous misinterpretations, compromising patient safety and surgical outcomes [1]. These ambiguities manifest in multiple forms, including linguistic, contextual, and procedural ambiguities, each creating an opportunity for critical misunderstanding.

Traditional approaches to ambiguity management in human-robot instruction often rely on rule-based systems, predefined ontologies, or simple clarification dialogues [2]. While functional in controlled environments with limited instruction variability, these methods frequently struggle within the dynamic and complex nature of surgical scenarios. Surgical instructions often contain domain-specific terminology, rely on implicit knowledge, and carry high stakes for misinterpretation, exceeding the capabilities of
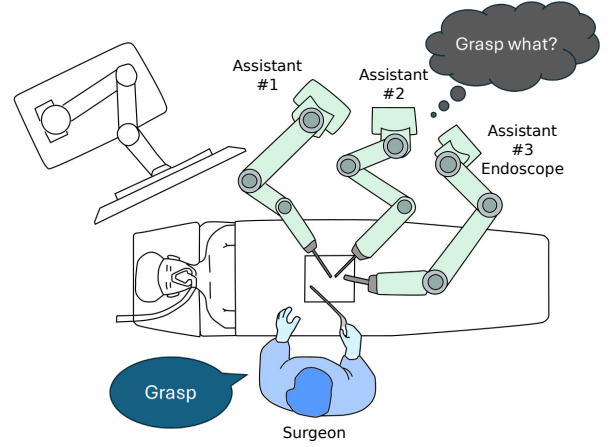
many conventional systems. Furthermore, these approaches may fail to capture the full spectrum of ambiguity types encountered in realistic operational settings.

Recent advancements in Large Language Models (LLMs) demonstrate remarkable capabilities in natural language understanding, context sensitivity, and nuanced reasoning. Their ability to process complex language and identify subtle semantic relationships suggests significant potential for overcoming the limitations of rule-based systems in ambiguity detection. While LLMs have been explored for enhancing robot task planning [3] and improving human-robot collaborative workflows [4], their specific application to systematic ambiguity detection in safety-critical surgical instructions remains relatively underexplored.

This paper introduces a novel framework for detecting ambiguity in natural language instructions for surgical robotic assistants as shown in Figure 1. Our approach leverages an ensemble of specialized LLM-based evaluators, each designed to identify distinct ambiguity types through targeted prompting strategies. By combining these diverse perspectives, we aim to comprehensively capture multiple dimensions of potential instruction ambiguity. To transform these evaluations into reliable ambiguity classifications, we implement a conformal prediction methodology that provides statistical guarantees about classification confidence. This approach allows a robotic system to quantify uncertainty in instruction interpretation and appropriately determine when clarification is necessary before execution.

The main contributions of this paper are:
- Development of an ensemble of LLM evaluators for

[1] Institutes of Innovation for Future Society, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8601, Japan

Correspondence: `davila.ana@robo.mein.nagoya-u.ac.jp`

[2] Department of Micro-Nano Mechanical Science and Engineering, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8603, Japan

surgical robotics to address various ambiguities.

- Use of conformal prediction to quantify ambiguity via calibrated scores.
- Demonstration of the approach's effectiveness with modern LLMs in recognizing ambiguous instructions.

## II. RELATED WORKS

### A. Ambiguity in Natural Language

Natural language ambiguity presents significant challenges for human-robot interaction systems, particularly in contexts requiring precise instruction interpretation. The problem of ambiguity has been well studied in computational linguistics and natural language processing (NLP). Foundational work in NLP categorizes ambiguities (e.g., lexical, syntactic, semantic, pragmatic) [5] and highlights challenges such as referential ambiguity, where context is needed to resolve which object an instruction refers to [6]. Recent work explores the use of LLMs themselves to handle ambiguity, for example, aligning LLM outputs with ambiguous inputs through self-assessment of uncertainty [7] or using structured taxonomies to guide instruction refinement [8]. These studies provide core concepts and LLM-based techniques relevant to the interpretation of potentially ambiguous robotic instructions.

### B. Ambiguity Detection in Robotic Systems

The robotics community has developed specialized approaches for handling unclear instructions across various domains. Wang et al. [9] present an Ask-when-Needed framework where LLMs evaluate instruction clarity before execution, seeking clarification only when necessary. Other studies integrate LLM with perception (vision-language models) to resolve object reference ambiguities during tasks such as manufacturing or collaborative assembly [10]. Interactive systems like SeeAsk strategically query users when visual grounding is uncertain during grasping tasks [11]. These approaches demonstrate the value of LLMs in identifying and sometimes resolving ambiguity in general robotic contexts, often through interaction or multi-modal sensing.

However, in the surgical context, the risks associated with ambiguity are amplified due to its high stakes, specialized terminology, dynamic environment, and need for extreme precision. Studies have highlighted the difficulty in translating qualitative voice commands (for example, "move more left") into precise robotic actions, especially with changing visual fields during procedures such as endoscopy [12]. Although several natural language interfaces for surgical robots have been developed [13], [14], [15], improving usability, they often lack dedicated mechanisms for systematic detection of ambiguities. Some research implicitly addresses ambiguity through multi-modal inputs like gesture control [16] or utilizes multi-modal LLMs for specific surgical sub-tasks like interpreting context for robotic assistance in blood suction [17]. However, a comprehensive framework specifically targeting the detection of diverse ambiguity types within natural language instructions for surgical robots remains an open area.
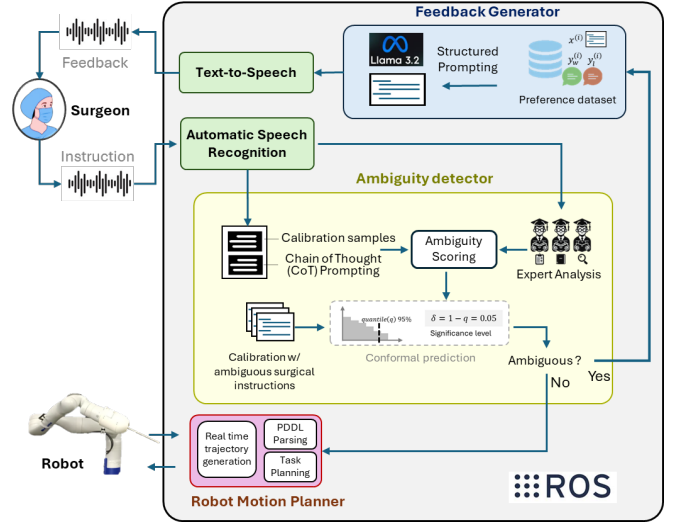


Fig. 2. Proposed framework for ambiguity detection

To move beyond simple detection towards reliable decision-making (i.e., deciding whether to execute or clarify), methods for quantifying uncertainty are essential. Conformal Prediction (CP) [18] has emerged as a promising statistical framework for providing rigorous confidence bounds on predictions, including those from LLMs. It has been applied in robotics to estimate uncertainty in learning [19], enabling systems to request clarification when confidence is low. Specific applications relevant to instruction ambiguity include ensuring probabilistically correct execution by identifying unsafe commands [20], aligning LLM planner uncertainty with task requirements across various ambiguity modes (e.g., spatial, numeric) [21], [22], and developing metrics to distinguish resolvable ambiguity from model hallucinations [23]. These works demonstrate the potential of CP to provide a principled basis for managing ambiguity by quantifying the model's certainty about its interpretation.

Although existing research addresses ambiguity from NLP foundations to general robotics, there remains a gap in systematically detecting the multiple facets of ambiguity inherent in natural language commands specifically within the high-stakes surgical context.

## III. METHODOLOGY

Our framework for detecting ambiguity in natural language instructions for surgical robots integrates an ensemble of Large Language Model (LLM)-based evaluators with a conformal prediction (CP) mechanism. Each evaluator is specialized to assess different facets of potential ambiguity. The CP framework then provides a statistically rigorous classification of the instruction's clarity based on the collective assessment of the evaluator ensemble. Figure 2 illustrates the overall architecture.

### A. Input Processing

The system accepts natural language requests, typically originating as spoken commands in a surgical setting. As
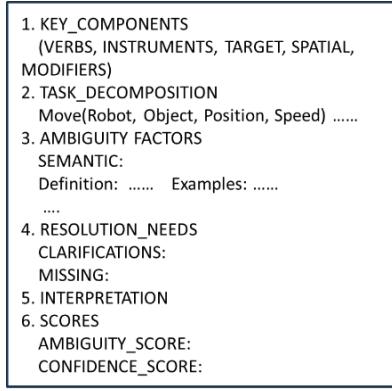
Fig. 3. Structure of CoT prompt



Fig. 4. Examples of instructions from the calibration dataset

described in our previous work [14], these commands are transcribed into text using a dedicated speech-to-text module integrated via ROS, making them suitable for processing by the subsequent LLM-based pipeline.

### B. LLM-Based Evaluator Ensemble

A central component is an ensemble of five distinct LLM-based evaluators designed to provide a multi-faceted assessment of instruction ambiguity. This ensemble approach leverages diverse perspectives for a more robust evaluation compared to a single, general assessment. Each evaluator utilizes a base LLM, prompted with specific instructions and illustrative examples, to generate an ambiguity score on a continuous scale from 0 (indicating complete clarity) to 10 (indicating high ambiguity).

*1) Chain-of-Thought (CoT) Evaluator:* One evaluator utilizes a Chain-of-Thought (CoT) prompting strategy [24]. As depicted in Figure 3, this involves guiding the LLM through a step-by-step reasoning process. The prompt is structured to first encourage the LLM to identify the key components within the instruction, such as the main verb, target objects, and any contextual elements. Following this identification, the LLM is prompted to decompose the task into a predefined list of potential robot actions relevant to surgical procedures. Subsequently, the definition and examples of multiple ambiguity factors are provided to the LLM to critically assess the potential for ambiguity based on the identified components and their relationships. These ambiguity factors include semantic, syntactic, pragmatic, and contextual types. We then request the LLM to identify potential clarifications that might be needed and any missing information required for unambiguous execution. Based on this comprehensive analysis, the LLM is asked to provide possible interpretations of the request and finally output an ambiguity score on a scale from 0 to 10, where a higher score indicates greater ambiguity. The final output of this evaluator is this ambiguity score reflecting the LLM's confidence in the instruction's clarity based on its step-by-step reasoning.

*2) Specialized Ambiguity Evaluators:* Four additional evaluators are prompt-tuned to focus on specific ambiguity types prevalent in surgical instructions:

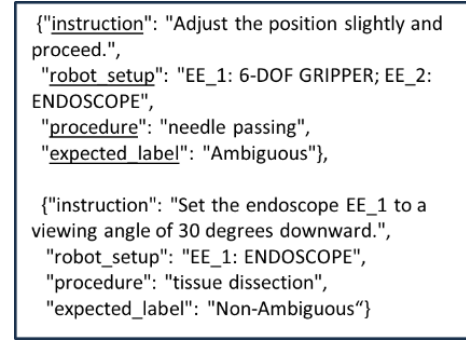- **Linguistic Evaluator:** Targets ambiguities originating from word meaning (lexical ambiguity, such as the interpretation of 'retract'), sentence structure (syntactic ambiguity, involving potentially unclear relations between phrases), or unresolved references (referential ambiguity, concerning pronouns like 'it').
- **Contextual Evaluator:** Assesses whether the instruction is adequately grounded in the implicit or explicit surgical context. It identifies instructions that depend on unclear assumptions or lack necessary situational details.
- **Procedural Evaluator:** Focuses on the clarity and completeness of the requested action sequence. This includes identifying underspecified parameters, exemplified by terms like 'slightly', unclear ordering of operations, or omissions of necessary steps.
- **Critical Safety Evaluator:** Specifically screens for ambiguities with the potential to induce safety risks upon misinterpretation. Such ambiguities could involve potential tissue damage, inappropriate instrument usage, or deviations from established surgical safety protocols.

### C. Calibration Dataset

The Conformal Prediction framework necessitates a calibration dataset. We constructed such a dataset comprising 40 natural language instructions representative of common surgical tasks, including endoscope navigation, instrument manipulation, and tissue handling [25], [26], [27], [28]. Each of these 40 requests was manually labeled as either "ambiguous" or "non-ambiguous". The labeling criteria considered factors such as the clarity of the action, the specificity of the target objects, the sufficiency of the context, and the potential for misinterpretation in a surgical setting. This labeled dataset serves as the basis for calibrating our conformal predictor.

### D. Conformal Prediction Framework

We employ conformal prediction [18] to translate the ambiguity scores from the LLM ensemble into classifications with statistical validity guarantees.

*1) Ensemble Score Aggregation:* For an input instruction $i$, let $s_{i,k}$ denote the ambiguity score from the $k$-th evaluator ($k = 1, \ldots, 5$). We compute the mean score $\mu_i = \frac{1}{5} \sum_{k=1}^{5} s_{i,k}$ and the sample variance $\sigma_i^2 = \frac{1}{4} \sum_{k=1}^{5} (s_{i,k} - \mu_i)^2$.

*2) Nonconformity Score Calculation:* A nonconformity measure quantifies the atypicality of a new instruction $i$ relative to the calibration examples associated with a candidate class $\delta \in$ Ambiguous, Non-ambiguous. The nonconformity score $NC_{\delta_i}$ for instruction $i$, under the hypothesis that its true class is $\delta$, is defined as:

$$NC_{\delta_i}(\mu_i, \sigma_i^2) = |\mu_i - \mu_{\text{cal},\delta}| + \beta \cdot \sigma_i^2 \qquad (1)$$

where:

- $\mu_{\text{cal},\delta}$ is the pre-computed mean of the average ambiguity scores ($\mu_j$) for all calibration examples $j$ belonging to class $\delta$.
- $\beta \geq 0$ is a hyperparameter balancing the contributions of the mean score deviation and the score variance. A larger $\beta$ assigns higher nonconformity to instructions where the evaluators exhibited significant disagreement, reflected by a large $\sigma_i^2$.

This score is computed twice for each new instruction $i$, yielding $NC_{\text{Amb}_i}$ and $NC_{\text{NonAmb}_i}$. A low $NC_{\delta_i}$ indicates that the instruction's score profile ($\mu_i, \sigma_i^2$) is characteristic of calibration examples from class $\delta$.

*3) P-value Calculation and Classification Decision:* Utilizing the distribution of nonconformity scores from the calibration set, we compute a p-value for each hypothesis class $\delta$ for the new instruction $i$. The p-value $p_{\delta_i}$ represents the proportion of calibration examples of class $\delta$ that are at least as non-conforming as instruction $i$:

$$p_{\delta_i} = \frac{|j \in \text{Cal}_\delta : NC_{\delta_j} \geq NC_{\delta_i}| + 1}{|\text{Cal}_\delta| + 1} \qquad (2)$$

Here, $\text{Cal}_\delta$ denotes the subset of calibration examples labeled $\delta$, and $NC_{\delta_j}$ is the nonconformity score for calibration example $j$ under the hypothesis $\delta$. The addition of 1 in the numerator and denominator provides a standard adjustment for finite calibration sets.

The instruction is classified based on these p-values and a pre-specified significance level $\alpha$.

$$\delta_{\text{out}} = \begin{cases} \text{Ambiguous} & \text{if } p_{\text{Amb}_i} > \alpha \text{ and } p_{\text{NonAmb}_i} \leq \alpha \\ \text{Non-ambiguous} & \text{if } p_{\text{Amb}_i} \leq \alpha \text{ and } p_{\text{NonAmb}_i} > \alpha \\ \text{Uncertain} & \text{if } p_{\text{Amb}_i} \leq \alpha \text{ and } p_{\text{NonAmb}_i} \leq \alpha \\ \text{Uncertain} & \text{if } p_{\text{Amb}_i} > \alpha \text{ and } p_{\text{NonAmb}_i} > \alpha \end{cases}$$
$$(3)$$

The primary classification outcomes are *Ambiguous* and *Non-ambiguous*.

In all other cases, where there is not a clear distinction based on the p-values, we classify the instruction as *Uncertain*, indicating that the system does not have sufficient statistical confidence to make a definitive classification at the chosen significance level. Both *Ambiguous* and *Uncertain* instructions can trigger a request for clarification from the human operator in a real-world surgical setting.

*E. Feedback Generation*

To provide users with actionable insights into the detected ambiguity, we implement a feedback generation module. This module utilizes a separate LLM instance, prompted to analyze the five evaluator scores ($s_{i,1}, \ldots, s_{i,5}$), identify the principal ambiguity factor from the CoT evaluator, and generate a one- or two-sentence message highlighting this primary issue and suggesting how the user might clarify their instruction, for instance: "Specify the location of the tissue to be cut". This facilitates an iterative clarification process between the user and the system.

## IV. EXPERIMENTAL VALIDATION

This section details the empirical evaluation of the proposed ambiguity detection framework using a curated dataset and standard classification metrics.

*A. Experimental Setup*

The evaluation utilized a dataset of 40 natural language instructions relevant to surgical robotics. This dataset was manually labeled and comprised 20 non-ambiguous instructions and 20 ambiguous instructions. The ambiguous set was specifically constructed to include 5 examples for each of four predominant ambiguity types: linguistic, contextual, procedural, and critical safety ambiguity. This structure allows for assessing the framework's ability to detect ambiguity overall and its sensitivity to different ambiguity manifestations.

We evaluated our framework's performance using two state-of-the-art Large Language Models (LLMs): Llama 3.2 11B and Gemma 3 12B. These models were accessed via the Transformers library from Hugging Face. All experiments were implemented in PyTorch 2.0 and conducted on a high-performance workstation featuring an AMD Ryzen Threadripper PRO 3995WX (2.7 GHz), 256 GB RAM, and an NVIDIA GeForce RTX A6000 GPU. For conformal prediction classification, the significance level ($\alpha$) was set to 0.1. The hyperparameter $\beta$, which weights the score variance in the nonconformity score calculation (Equation 1), was empirically determined to be 0.5 through preliminary tuning. Performance was assessed using standard classification metrics: accuracy, precision, recall, and F1-score.

*B. Quantitative Results*

The global performance in distinguishing between any ambiguous instruction and a non-ambiguous one is summarized by the confusion matrices in Figure 5.

Using Llama 3.2 11B, the framework achieved an overall accuracy of 70%. It correctly identified 16 out of 20 ambiguous instructions (Recall=0.80 for ambiguous class) and 12 out of 20 non-ambiguous instructions (Recall=0.60 for non-ambiguous class). Using Gemma 3 12B, the framework demonstrated superior performance, achieving an overall accuracy of 82.5%. This model correctly identified 17 out of 20 ambiguous instructions (Recall=0.85 for ambiguous class) and 16 out of 20 non-ambiguous instructions (Recall=0.80 for non-ambiguous class).
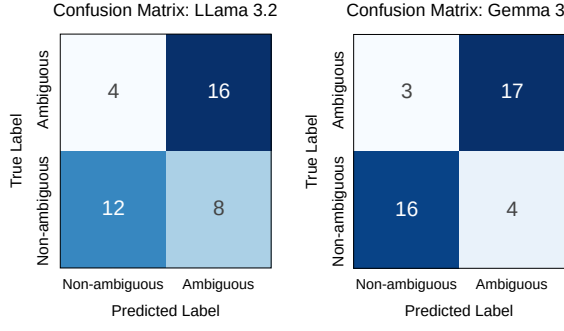
Fig. 5. Confusion matrices for ambiguity detection using the conformal prediction framework ($\alpha = 0.1$). Left: Llama 3.2 11B (70% accuracy). Right: Gemma 3 12B (82.5% accuracy).
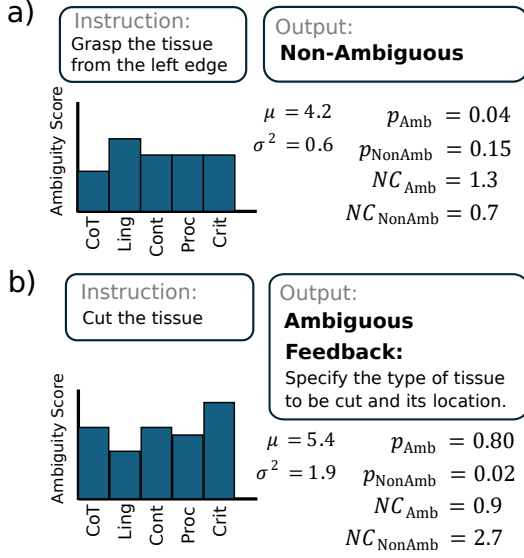


Fig. 6. Representative examples of the ambiguity detection process. (a) An instruction correctly classified as non-ambiguous. (b) An instruction correctly classified as ambiguous, with feedback generated based on the evaluators' outputs. The plots show individual evaluator scores and the resulting nonconformity scores (NCS).

To understand how effectively the framework identifies specific types of ambiguity when they are present, we analyzed performance metrics focused on each category within the ambiguous set. Tables I and II present these results. For each row, the metrics reflect the binary classification performance considering five pairs of samples, five ambiguous samples as true positives, and their corresponding five nonambiguous samples as true negatives. The "Total" row measures the performance across all 40 samples.

With Llama 3.2 11B (Table I), the framework correctly identified 80% (Recall = 0.80) of the procedural ambiguity examples. Its ability to identify other types was lower, correctly identifying 60% of linguistic and critical examples, and only 40% of contextual examples, although precision was high for contextual cases it did flag. With Gemma 3 12B (Table II), the framework achieved perfect recall (1.00) for linguistic ambiguities, identifying all 5 examples correctly. It also demonstrated a strong recall for contextual (0.80) and procedural (0.80) ambiguities. Performance on critical

| Ambiguity Type | F1-Score | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Linguistic | 0.60 | 0.60 | 0.60 | 0.60 |
| Contextual | 0.57 | 0.70 | 1.00 | 0.40 |
| Procedural | 0.80 | 0.80 | 0.80 | 0.80 |
| Critical | 0.67 | 0.70 | 0.75 | 0.60 |
| Total | 0.73 | 0.70 | 0.67 | 0.80 |

| Ambiguity Type | F1-Score | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Linguistic | 1.00 | 1.00 | 1.00 | 1.00 |
| Contextual | 0.73 | 0.70 | 0.67 | 0.80 |
| Procedural | 0.80 | 0.80 | 0.80 | 0.80 |
| Critical | 0.75 | 0.80 | 1.00 | 0.60 |
| Total | 0.83 | 0.83 | 0.81 | 0.85 |

ambiguities was lower (0.60 recall), although precision was perfect, meaning it made no false positive errors when classifying critical ambiguities within that subset. The overall F1 score across all 40 samples was 0.83.

### C. Qualitative Analysis

Figure 6 provides illustrative examples of the framework's classification process and feedback generation.

Example (a) shows an instruction correctly identified as *Non-ambiguous*, characterized by low scores from the evaluators and resulting p-values aligning with the non-ambiguous class hypothesis. Example (b) illustrates a case correctly identified as *Ambiguous*. Higher scores from specific evaluators, such as the *Contextual* one, contribute to p-values supporting the ambiguous classification. The feedback module then leverages the high-scoring evaluator type to generate a relevant clarification prompt for the user.

### D. Discussion

The experimental evaluation confirms the viability of the proposed LLM ensemble and conformal prediction framework for detecting ambiguity in surgical instructions. The results highlight the significant influence of the underlying LLM, with Gemma 3 12B achieving considerably higher overall accuracy and more consistent detection across different ambiguity types compared to Llama 3.2 11B in our setup. This underscores the importance of model selection for tasks requiring nuanced language understanding in specialized domains.

The analysis broken down by ambiguity type reveals varying levels of sensitivity. Both models, especially Gemma 3, showed proficiency in identifying linguistic and procedural ambiguities. This suggests the framework effectively captures issues related to word meaning, sentence structure, references, and action sequence clarity. Detecting contextual and critical safety ambiguities proved more challenging, as reflected in the recall scores for these categories, particularly for Llama 3.2. This difficulty likely stems from the inherent

reliance on implicit world knowledge, understanding of the ongoing surgical state, and safety constraints, which are harder for LLMs to grasp solely from the instruction text without richer context integration.

A primary limitation of this study is the small size of the evaluation dataset, particularly the limited number of examples for each specific ambiguity type. While indicative, the results regarding type-specific performance require validation on a larger, more diverse dataset to draw robust conclusions. Furthermore, the dataset construction and labeling inherently involve subjectivity, and the framework's performance depends on the calibration set reflecting the distribution of ambiguities encountered in practice. The choice of hyperparameters $(\alpha, \beta)$ and the specific nonconformity score formulation might also influence results and could be subject to further optimization.

## V. CONCLUSION

In this paper, we presented a novel approach for detecting ambiguity in natural language instructions for surgical robot assistants, employing an ensemble of specialized LLM-based evaluators specializing in different ambiguity types within a conformal prediction framework. Our experimental validation using Llama 3.2 11B and Gemma 3 12B demonstrated the effectiveness of this methodology, with Gemma 3 achieving an accuracy of 82.5% in distinguishing ambiguous from non-ambiguous instructions. Furthermore, we introduced a feedback generation mechanism to guide users in refining unclear requests. These results highlight the potential of our framework to enhance the safety and reliability of human-robot interaction in critical surgical environments. Future work will focus on expanding the calibration dataset, further refining the LLM evaluators and prompting strategies, and exploring the integration of this ambiguity detection system into a real-time surgical robotic platform.

## REFERENCES

[1] A. Satchidanand, J. Higginbotham, A. Bisantz, N. Aldhaam, A. El-sayed, I. Carr, A. A. Hussein, and K. Guru, ""put the what, where? cut here?!" challenges to coordinating attention in robot-assisted surgery: a microanalytic pilot study," *BMJ open*, vol. 11, no. 7, p. e046132, 2021.

[2] L. Villamar Gómez and J. Miura, "Ontology-based knowledge management with verbal interaction for command interpretation and execution by home service robots," *Robotics and Autonomous Systems*, vol. 140, p. 103763, 2021.

[3] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, "Progprompt: Generating situated robot task plans using large language models," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 11 523–11 530.

[4] C. Zhang, J. Chen, J. Li, Y. Peng, and Z. Mao, "Large language models for human–robot interaction: A review," *Biomimetic Intelligence and Robotics*, vol. 3, no. 4, p. 100131, 2023.

[5] B. Gleich, O. Creighton, and L. Kof, "Ambiguity detection: Towards a tool explaining ambiguity sources," in *Requirements Engineering: Foundation for Software Quality*, R. Wieringa and A. Persson, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 218–232.

[6] H. Yang, A. De Roeck, V. Gervasi, A. Willis, and B. Nuseibeh, "Analysing anaphoric ambiguity in natural language requirements," *Requirements engineering*, vol. 16, pp. 163–189, 2011.

[7] H. J. Kim, Y. Kim, C. Park, J. Kim, C. Park, K. M. Yoo, S. goo Lee, and T. Kim, "Aligning language models to explicitly handle ambiguity," 2024. [Online]. Available: https://arxiv.org/abs/2404.11972

[8] A. Niwa and H. Iso, "Ambignlg: Addressing task ambiguity in instruction for nlg," 2024. [Online]. Available: https://arxiv.org/abs/2402.17717

[9] W. Wang, J. Shi, Z. Ling, Y.-K. Chan, C. Wang, C. Lee, Y. Yuan, J. tse Huang, W. Jiao, and M. R. Lyu, "Learning to ask: When llm agents meet unclear instruction," 2025. [Online]. Available: https://arxiv.org/abs/2409.00557

[10] J. Fan and P. Zheng, "A vision-language-guided robotic action planning approach for ambiguity mitigation in human–robot collaborative manufacturing," *Journal of Manufacturing Systems*, vol. 74, pp. 1009–1018, 2024.

[11] Y. Mo, H. Zhang, and T. Kong, "Towards open-world interactive disambiguation for robotic grasping," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 8061–8067.

[12] K. Zinchenko, C.-Y. Wu, and K.-T. Song, "A study on speech recognition control for a surgical robot," *IEEE Transactions on Industrial Informatics*, vol. 13, no. 2, pp. 607–615, 2017.

[13] M. Elazzazi, L. Jawad, M. Hilfi, and A. Pandya, "A natural language interface for an autonomous camera control system on the da vinci surgical robot," *Robotics*, vol. 11, no. 2, 2022.

[14] A. Davila, J. Colan, and Y. Hasegawa, "Voice control interface for surgical robot assistants," in *2024 International Symposium on Micro-NanoMehatronics and Human Science (MHS)*, 2024, pp. 1–5.

[15] M. Moghani, L. Doorenbos, W. C.-H. Panitch, S. Huver, M. Azizian, K. Goldberg, and A. Garg, "Sufia: Language-guided augmented dexterity for robotic surgical assistants," in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2024, pp. 6969–6976.

[16] M. Jacob, Y.-T. Li, G. Akingba, and J. P. Wachs, "Gestonurse: a robotic surgical nurse for handling surgical instruments in the operating room," *Journal of Robotic Surgery*, vol. 6, pp. 53–63, 2012.

[17] S. Zargarzadeh, M. Mirzaei, Y. Ou, and M. Tavakoli, "From decision to action in surgical autonomy: Multi-modal large language models for robot-assisted blood suction," *IEEE Robotics and Automation Letters*, vol. 10, no. 3, pp. 2598–2605, 2025.

[18] A. N. Angelopoulos and S. Bates, "A gentle introduction to conformal prediction and distribution-free uncertainty quantification," 2022. [Online]. Available: https://arxiv.org/abs/2107.07511

[19] C. Celemin and J. Kober, "Knowledge-and ambiguity-aware robot learning from corrective and evaluative feedback," *Neural Computing and Applications*, vol. 35, no. 23, pp. 16 821–16 839, 2023.

[20] J. Wang, G. He, and Y. Kantaros, "Probabilistically correct language-based multi-robot planning using conformal prediction," *IEEE Robotics and Automation Letters*, vol. 10, no. 1, pp. 160–167, 2025.

[21] A. Z. Ren, A. Dixit, A. Bodrova, S. Singh, S. Tu, N. Brown, P. Xu, L. Takayama, F. Xia, J. Varley, Z. Xu, D. Sadigh, A. Zeng, and A. Majumdar, "Robots that ask for help: Uncertainty alignment for large language model planners," 2023. [Online]. Available: https://arxiv.org/abs/2307.01928

[22] K. Liang, Z. Zhang, and J. F. Fisac, "Introspective planning: Aligning robots'uncertainty with inherent task ambiguity," in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, Eds., vol. 37. Curran Associates, Inc., 2024, pp. 71 998–72 031.

[23] J. F. M. Jr. and D. Manocha, "Lap, using action feasibility for improved uncertainty alignment of large language model planners," 2024. [Online]. Available: https://arxiv.org/abs/2403.13198

[24] J. Wei, X. Wang, D. Schuurmans, M. Bosma, b. ichter, F. Xia, E. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 24 824–24 837.

[25] K. Fozilov, J. Colan, A. Davila, K. Misawa, J. Qiu, Y. Hayashi, K. Mori, and Y. Hasegawa, "Endoscope automation framework with hierarchical control and interactive perception for multi-tool tracking in minimally invasive surgery," *Sensors*, vol. 23, no. 24, p. 9865, 2023.

[26] Y. Yamada, J. Colan, A. Davila, and Y. Hasegawa, "Task segmentation based on transition state clustering for surgical robot assistance," in *2023 8th International Conference on Control and Robotics Engineering (ICCRE)*, 2023, pp. 260–264.

[27] S. Liu, J. Colan, Y. Zhu, T. Kobayashi, K. Misawa, M. Takeuchi, and Y. Hasegawa, "Latent regression based model predictive control for

tissue triangulation," *Advanced Robotics*, vol. 37, no. 24, pp. 1552–1565, 2024.

[28] Y. Yamada, J. Colan, A. Davila, and Y. Hasegawa, "Multimodal semi-supervised learning for online recognition of multi-granularity surgical workflows," *International Journal of Computer Assisted Radiology and Surgery*, vol. 19, no. 6, pp. 1075–1083, 2024.