

CliniDial: A Naturally Emerged Multimodal Dialogue Dataset for Team Reflection During Clinical Operation

Anonymous ACL submission

Abstract

In clinical operations, teamwork can be the crucial factor that determines the final outcome. Prior studies have shown that there can be 58% more deaths than expected due to insufficient collaboration. To understand how the team practices teamwork during the operation, we collected *CliniDial* from simulations of medical operations. *CliniDial* includes the audio data and its transcriptions, the simulated physiology signals of the patient manikins, and how the team operates from two camera angles. We annotated behavior codes following an existing framework to understand the teamwork process. Experimental results show that *CliniDial* poses significant challenges to the existing models.

1 Introduction

In clinical settings, teamwork is crucial for a successful operation, and effective team collaboration can improve the safety and well-being of the patients (Catchpole et al., 2008; Weaver et al., 2010; Schmutz et al., 2019; Rosen et al., 2018). Failures in teamwork and communication among healthcare providers are a major contributing factor to the estimated 250,000 preventable deaths that occur in the U.S. each year (Rosen et al., 2018; Makary and Daniel, 2016). Breakdowns in areas like leadership, situation awareness, decision-making and communication frequently underlie the many forms of preventable patient harm, including hospital infections, falls, diagnostic errors and surgical mistakes (Baker et al., 2005; Herzberg et al., 2019; Keers et al., 2013). There can be 58% more deaths than expected due to insufficient collaboration (Knaus et al., 1986). Motivated by these statistics, in this paper we model the communication between team members as well as the data in the operation room to detect the effective steps and interactions needed for a successful procedure.

To understand how teamwork unfolds in the operating room, we collected *CliniDial* from simula-

tions of medical operations. We collected the audio data, simulated physical signals from the patient manikins, as well as how the team operates from two camera angles. We then annotated behavior codes based on a team reflection behavior framework (Schmutz et al., 2021) to understand how the team members convey their objectives, strategies, and actions during the operation. We test various baseline methods with different setups. Experimental results show that *CliniDial* poses significant challenges to existing methods. In addition, we invite input from medical professionals to try to bridge the current NLP fields with the real-world applications they expect (Appendix E).

In summary, our contributions are two folds:

1. We present *CliniDial*, a naturally emerged multimodal dialogue dataset for team reflection during clinical operation.
2. We evaluate our dataset against various existing methods with different setups and provide a detailed analysis of their results. Our experimental results reveal that our dataset poses significant challenges to existing methods, urging methodology innovation in our NLP community.

2 How is *CliniDial* Different?

Our real-world setting distinguishes *CliniDial* from existing datasets in various aspects. First, there are significant **label imbalances** in the collected data. Such label imbalances are less common in conventional NLP datasets where researchers have some levels of control over the data distribution by data filtering or downsampling. However, since our dialogues occur naturally in the operation room, the interlocutors are not tasked to generate dialogues but rather to perform the clinical operation and take care of the “patient” as a team. We do not pose any constraints on how the team communicate, and we observe that the amount of majority class labels sig-

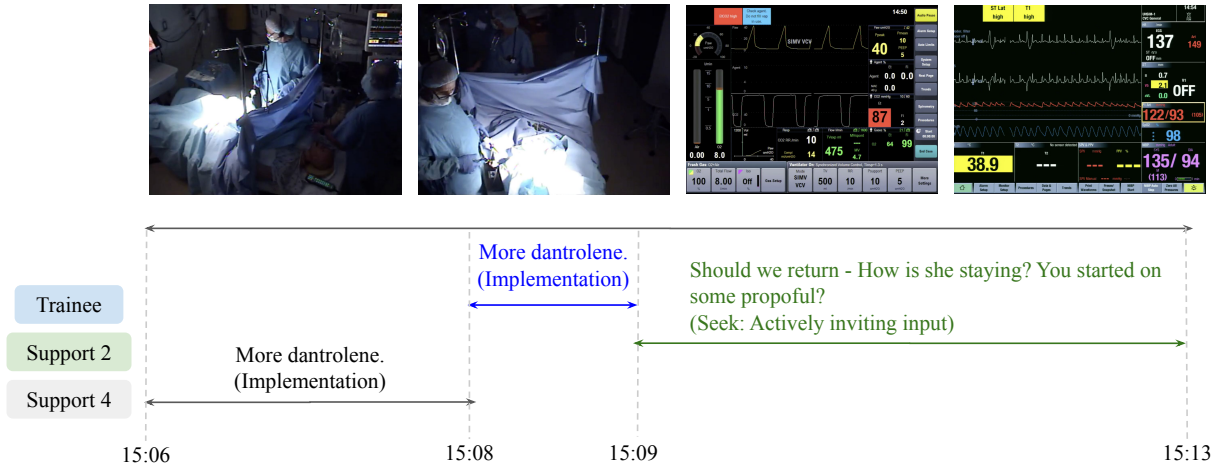


Figure 1: An example of the labeled dialogue in the simulated operation. Two cameras capture the scenes from two angles and two real-time monitoring systems provide the patient’s physiological signals. We only include the trainee and the two supports in this example, as they are the only three people speaking during this time frame.

079 significantly outmatches the minority class labels. Sec- 109
 080 ond, there are **rich and natural interactions** between 110
 081 the team members. Compared to the conven- 111
 082 tional dialogue benchmarks (Budzianowski et al., 112
 083 2018) which typically contain 30 turns at most, the 113
 084 dialogue in our collected dataset contains 311 turns 114
 085 on average. Third, there are **rich modalities** in the 115
 086 collected data. Compared to the conventional NLP 116
 087 datasets with text modality (Chen et al., 2021) or 117
 088 the conventional multimodal datasets which focus 118
 089 on vision and text modalities (Tapaswi et al., 2016; 119
 090 Lei et al., 2018; Castro et al., 2022), the data we 120
 091 collect includes not only the dialogue, but also the 121
 092 corresponding audio, the operation views from two 122
 093 camera angles, and the physiological signals from 123
 094 the “patient” aligned for each timestamp. 124

095 3 Dataset 125

096 3.1 Data Descriptions 126

097 **Scenarios.** A team of board certified anesthesiol- 127
 098 ogists together with support staff is tasked with the 128
 099 intraoperative management of a 36-year-old female 129
 100 who is undergoing a minimally invasive surgery 130
 101 ¹. This scenario takes place in a simulated oper- 131
 102 ating room where we present a mannequin as the 132
 103 female patient and simulate her physiological sig- 133
 104 nal changes from the backend. Specifically, the 134
 105 patient develops malignant hyperthermia (MH; a 135
 106 rare complication of general anesthesia that could 136
 107 develop in any patient) as the simulated scenario 137
 108 progresses. Many healthcare providers lack suffi-

¹The patient was diagnosed with acute cholangitis and is undergoing laparoscopic cholecystectomy

cient clinical exposure to MH, potentially hindering 109
 their ability to recognize, treat, and manage these 110
 rare but severe cases effectively (Isaak and Stiegler, 111
 2016). We want to stress that this is not a real oper- 112
 ation, and the intent is to train medical trainees in 113
 “near-life” surgical operations. 114

Roles. In the simulated operation, a confederate 115
 plays the role of the surgeon. The trainee who 116
 serves as the anesthesiologist is the main decision- 117
 maker ². The support participants are also trainees 118
 who support an anesthesiologist. Appendix B pro- 119
 vides additional details of the simulated operation 120
 and the roles of the team members. 121

122 3.2 Labels 122

Following Schmutz et al. (2021), we include three 123
 labels of “Seek”, “Evaluate” and “Plan” detailed 124
 in Table 3 in Appendix A. As our data is sourced 125
 from clinical operations, we are interested in not 126
 only how the teams engage in reflection or diag- 127
 nostic behaviors, but also how the team progresses 128
 from diagnostic actions to interventions or imple- 129
 mentation actions. Therefore, we assign an extra 130
 label “Implement” to such behaviors. Appendix A 131
 provides more details for each label. We describe 132
 the details of our annotation in Appendix C.2. 133

134 3.3 Dataset Statistics 134

Figure 1 provides an example of the annotated dia- 135
 logue in the simulated operation. Table 1 provides 136
 the statistics of our collected dataset. Table 2 pro- 137

²This is because malignant hyperthermia is a body’s adverse reaction to an anesthetic.

General	# Sessions	22
	# Participants / Session	6
Language	# Turns	6.5k
	# Words	49.9k
	# Turns / Session	311
	# Words / Session	2.3k
Others	Duration (min) / Session	19
	# Camera Angles	2
	# Physiological Signals	9

Table 1: Statistics of our collected dataset.

Label	None	Seek	Eval	Impl	Plan	All
Num	3.7k	1.3k	0.8k	0.6k	0.3k	6.9k

Table 2: Label distributions.

vides the label distributions. Appendix C provides more information for the dataset as well as the physiological signals included.

We apply ten-fold cross-validation on our dataset and report the average macro and micro F1 scores in the following setups. For each fold, we use 17, 2, and 3 sessions for training, validation, and testing, respectively.

4 Characteristic I: Imbalanced Class Distribution

Evaluation Setups. Here we constrain our study within the text domain and apply different methods to handle label imbalance. We directly tune a BERT base (Devlin et al., 2019) model to learn directly from the skewed data. In addition, we prompt the 8B and 70 B versions of the open-source Llama 3 model (abbreviated as Llama in figures) and closed-source GPT-4 and GPT-4o with and without demonstrations. Appendix D provides additional baseline models and their results.

Discussions. Figure 2 compares the F1 scores averaged by class (macro F1 scores) and F1 scores averaged by instances (micro F1 scores). For the tuning-based methods such as BERT_{base}, though it can achieve the highest micro F1 score, the macro F1 score remains much lower and is comparable to 0-shot or few-shot performances by GPT-4. This suggests that tuning-based methods bias the model to better learn the majority class, while the LLMs with a few demonstrations from each class do not suffer from the performance disparity between the macro and micro F1 scores. There is a significant performance boost for Llama 8B from 0-shot to 1-shot, suggesting even a single example can guide

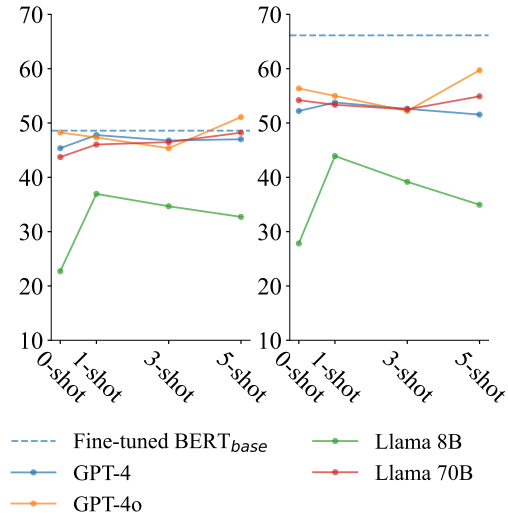


Figure 2: Comparison of macro F1 scores (F1 scores averaged by class, on the left) and micro F1 scores (F1 scores averaged by instances, on the right) versus number of demonstrations (number of shots). We compare both scores for the fine-tuned BERT_{base} model, 0-shot and few-shot prompting for LLMs.

smaller LLMs to better reason. In contrast, there is no significant performance boost if we increase the number of demonstrations in the few-shot setting. We hypothesize that since our dataset includes dialogues happening in the real world, there is a diverse forms of dialogue patterns. Therefore, a few demonstrations may be insufficient for the model to assess all the possible situations.

5 Characteristic II: Conversational Nature.

As shown in Figure 1, people are interacting with each other to actively communicate information in the operation process. Hence, an ideal model would leverage the context information of the interaction to better assess the current situation.

Evaluation Setups. We take the best performed closed-source LLM, GPT-4o, and the best performed open-source LLM, Llama 70B from Figure 2. We then prompt them with one turn both before and after the current round (context size of 3 in Figure 3) or two turns before and after the current turn (context size of 5 in Figure 3). In both situations, we report the performance by providing no demonstration (0-shot) or a single demonstration (1-shot).

Discussions. Figure 3 reports the performance comparison across different settings. For GPT-4o,

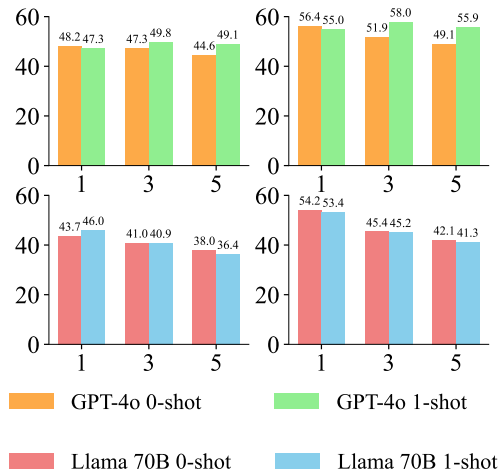


Figure 3: Comparison of macro F1 scores (F1 scores averaged by class, on the left) and micro F1 scores (F1 scores averaged by instances, on the right) versus the context size (x-axis). For instance, “3” on x-axis represents a context of size “3”, where we include one turn both before and after the current turn in our prompt to the LLM.

it may leverage the demonstration as well as the context information better, and acquires the best results when we feed it with one demonstration and context size of 3. In contrast, the demonstration and the context hurt Llama 3’s performance. One possible reason could be because of the long input prompt. On average, there is around 1,000 tokens per example if we feed the context information and one demonstration, while LLMs with a smaller context window size like Llama 3 may struggle with such long context information, similar to the findings by He et al. (2024).

6 Characteristic III: Multimodality Beyond Text and Vision.

Evaluation Setups. We evaluate the GPT-4o model, a multimodal end-to-end LLM with different modalities as the input, including feeding pure text (T), text and the operation video from two angles (+V), text and the physiology signals (+P). In addition, we try to let GPT-4o first verbalize what happens in the video, then pass the verbalized version of the information to GPT-4o together with the dialogue and instructions.

Discussion. From Figure 4, we can see that GPT-4o may still fail to leverage the visual or the physiological signals effectively. Moreover, when we verbalize the information for the physiological signals, GPT-4o suffers a 2% performance drop for

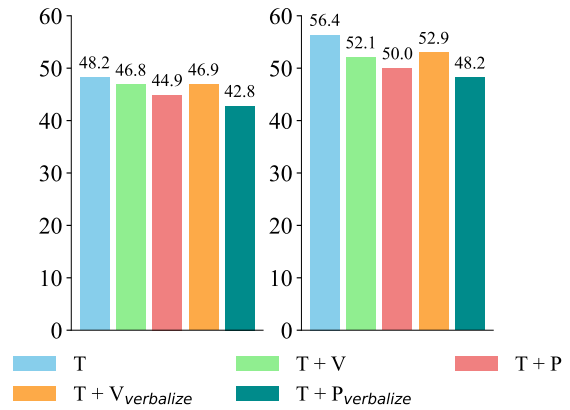


Figure 4: Comparison of macro F1 scores (F1 scores averaged by class, on the left) and micro F1 scores (F1 scores averaged by instances, on the right) when we pass in different modalities. “T” stands for text-only, “V”, “P” stand for visual signals and physiology signals, respectively. “T + V_{verbalize}” and “T + P_{verbalize}” stand for verbalizing the content by GPT-4o first, and then pass the text description with the other instructions to the GPT-4o model.

both macro and micro F1 scores. This indicates that though GPT-4o is capable handling a variety of vision tasks (Deng et al., 2024; Zhou et al., 2024b), reasoning over frames that require medical domain knowledge remains challenging. In addition, adding visual modality also hurts GPT-4o’s performance on our task. As we sample the frames corresponding to the timestamps when the dialogue happens, some cases may correspond to ten or more frames. For such cases, GPT-4o may struggle to leverage information from the video effectively, which aligns with the finding by Zhou et al. (2024a).

7 Conclusion

In this paper, we introduced *CliniDial*, a naturally emerged dialogue dataset from the clinical operation with distinguished characteristics from the existing benchmarks. Through experiments, we showed that existing methods do not work well on *CliniDial*. We hope the described characteristics in *CliniDial* could invite future effort to close the gap between NLP methods and the real-world applications.

Limitations

Due to the difficulty of setting up the environment and the data collection process, the dataset is collected mainly on 22 clinical operation sessions.

254	However, we note that there are 6.5k turns and	Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang.	305
255	49.9k words in total in <i>CliniDial</i> . We demonstrate	2021. DialogSum: A real-life scenario dialogue sum-	306
256	the distinct characteristics of our data from the ex-	marization dataset . In <i>Findings of the Association</i>	307
257	isting benchmarks and provide the performance of	<i>for Computational Linguistics: ACL-IJCNLP 2021</i> ,	308
258	popular NLP methods. However, due to the scope	pages 5062–5074, Online. Association for Computa-	309
259	of this study, we cannot evaluate every possible	tional Linguistics.	310
260	method and would like to invite future effort on		
261	a comprehensive evaluation of NLP methods on	Naihao Deng, Zhenjie Sun, Ruiqi He, Aman Sikka, Yu-	311
262	clinical data.	long Chen, Lin Ma, Yue Zhang, and Rada Mihalcea.	312
263		2024. Tables as texts or images: Evaluating the ta-	313
264		ble reasoning ability of llms and mllms . <i>Preprint</i> ,	314
265		arXiv:2402.12424.	315
266			
267	Ethics Statement	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	316
268		Kristina Toutanova. 2019. BERT: Pre-training of	317
269	We note that the study was approved by the Insti-	deep bidirectional transformers for language under-	318
270	tutional Review Board. Since the data from the	standing . In <i>Proceedings of the 2019 Conference of</i>	319
271	two cameras may reveal the identity of the team,	<i>the North American Chapter of the Association for</i>	320
272	we may not release the camera data. We are con-	<i>Computational Linguistics: Human Language Tech-</i>	321
273	sidering to release an anonymized version of the	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	322
274	dialogue transcription to facilitate future research	4171–4186, Minneapolis, Minnesota. Association for	323
275	on clinical NLP. We expect researchers to continue	Computational Linguistics.	324
276	building new algorithms and methods on top of this		
277	clinical dataset.	Qianyu He, Jie Zeng, Wenhao Huang, Lina Chen, Jin	325
278		Xiao, Qianxi He, Xunzhe Zhou, Jiaqing Liang, and	326
279		Yanghua Xiao. 2024. Can large language models	327
280		understand real-world complex instructions? In <i>Pro-</i>	328
281		<i>ceedings of the AAAI Conference on Artificial Intelli-</i>	329
282		<i>gence</i> , volume 38, pages 18188–18196.	330
283			
284	References	Simone Herzberg, Matt Hansen, Amanda Schoonover,	331
285		Barbara Skarica, James McNulty, Tabria Harrod,	332
286	Joseph Allen, Hayley Hung, Joann Keyton, Gabriel Mur-	Jonathan M Snowden, William Lambert, and Jeanne-	333
287	ray, Catharine Oertel, and Giovanna Varni. 2021. In-	Marie Guise. 2019. Association between measured	334
288	sights on group and team dynamics. In <i>Proceedings</i>	teamwork and medical errors: an observational	335
289	<i>of the 2021 International Conference on Multimodal</i>	study of prehospital care in the usa. <i>BMJ open</i> ,	336
290	<i>Interaction</i> , pages 855–856.	9(10):e025314.	337
291			
292	David P Baker, Sigrid Gustafson, Jeff Beaubien, Ed-	Hayley Hung, Litian Li, Jord Molhoek, and Jing Zhou.	338
293	uardo Salas, and Paul Barach. 2005. Medical team-	2024. The discontent with intent estimation in-the-	339
294	work and patient safety: the evidence-based relation.	wild: The case for unrealized intentions. In <i>Extended</i>	340
295	<i>AHRQ publication</i> , 5(53):1–64.	<i>Abstracts of the CHI Conference on Human Factors</i>	341
296		<i>in Computing Systems</i> , pages 1–9.	342
297			
298	Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang	Robert Scott Isaak and Marjorie Podraza Stiegler. 2016.	343
299	Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ra-	Review of crisis resource management (crm) princi-	344
300	madan, and Milica Gašić. 2018. MultiWOZ - a large-	ples in the setting of intraoperative malignant hyper-	345
301	scale multi-domain Wizard-of-Oz dataset for task-	thermia. <i>Journal of anesthesia</i> , 30:298–306.	346
302	oriented dialogue modelling . In <i>Proceedings of the</i>		
303	<i>2018 Conference on Empirical Methods in Natural</i>	Richard N Keers, Steven D Williams, Jonathan Cooke,	347
304	<i>Language Processing</i> , pages 5016–5026, Brussels,	and Darren M Ashcroft. 2013. Causes of medication	348
305	Belgium. Association for Computational Linguistics.	administration errors in hospitals: a systematic re-	349
306		view of quantitative and qualitative evidence. <i>Drug</i>	350
307		<i>safety</i> , 36:1045–1067.	351
308			
309		Florian Klonek, Fabiola Heike Gerpott, Nale Lehmann-	352
310		Willenbrock, and Sharon K Parker. 2019. Time to	353
311		go wild: How to conceptualize and measure process	354
312		dynamics in real teams with high-resolution. <i>Organi-</i>	355
313		<i>zational Psychology Review</i> , 9(4):245–275.	356
314			
315		William A Knaus, Elizabeth A Draper, Douglas P Wag-	357
316		ner, and Jack E Zimmerman. 1986. An evaluation of	358
317		outcome from intensive care in major medical centers.	359
318		<i>Annals of internal medicine</i> , 104(3):410–418.	360
319			
320			
321			
322			
323			
324			
325			
326			
327			
328			
329			
330			
331			
332			
333			
334			
335			
336			
337			
338			
339			
340			
341			
342			
343			
344			
345			
346			
347			
348			
349			
350			
351			
352			
353			
354			
355			
356			
357			
358			
359			
360			

361	Michaela Kolbe and Margarete Boos. 2019. Laborious but elaborate: the benefits of really studying team dynamics. <i>Frontiers in psychology</i> , 10:433269.	415
362		416
363		417
364	Nale Lehmann-Willenbrock and Hayley Hung. 2023. A multimodal social signal processing approach to team interactions. <i>Organizational Research Methods</i> , page 10944281231202741.	418
365		419
366		420
367		
368	Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg. 2018. TVQA: Localized, compositional video question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1369–1379, Brussels, Belgium. Association for Computational Linguistics.	421
369		422
370		423
371		424
372		425
373		
374	Martin A Makary and Michael Daniel. 2016. Medical error—the third leading cause of death in the us. <i>Bmj</i> , 353.	426
375		427
376		428
377	Michael A Rosen, Deborah DiazGranados, Aaron S Dietz, Lauren E Benishek, David Thompson, Peter J Pronovost, and Sallie J Weaver. 2018. Teamwork in healthcare: Key discoveries enabling safer, high-quality care. <i>American Psychologist</i> , 73(4):433.	429
378		430
379		
380		
381		
382	Jan B Schmutz, Zhike Lei, and Walter J Eppich. 2021. Reflection on the fly: development of the team reflection behavioral observation (turbo) system for acute care teams. <i>Academic Medicine</i> , 96(9):1337–1345.	
383		
384		
385		
386	Jan B Schmutz, Laurenz L Meier, and Tanja Manser. 2019. How effective is teamwork really? the relationship between teamwork and performance in healthcare teams: a systematic review and meta-analysis. <i>BMJ open</i> , 9(9):e028280.	
387		
388		
389		
390		
391	Chris Seiffert, Taghi M Khoshgoftaar, Jason Van Hulse, and Amri Napolitano. 2009. Rusboost: A hybrid approach to alleviating class imbalance. <i>IEEE transactions on systems, man, and cybernetics-part A: systems and humans</i> , 40(1):185–197.	
392		
393		
394		
395		
396	Christina Stevenson, Avneesh Bhangu, James J Jung, Aidan MacDonald, and Brodie Nolan. 2022. The development and measurement properties of the trauma non-technical skills (t-notechs) scale: a scoping review. <i>The American Journal of Surgery</i> , 224(4):1115–1125.	
397		
398		
399		
400		
401		
402	Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. Movieqa: Understanding stories in movies through question-answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 4631–4640.	
403		
404		
405		
406		
407		
408	Sallie J Weaver, Michael A Rosen, Deborah DiazGranados, Elizabeth H Lazzara, Rebecca Lyons, Eduardo Salas, Stephen A Knynch, Margie McKeever, Lee Adler, Mary Barker, et al. 2010. Does teamwork improve performance in the operating room? a multi-level evaluation. <i>The Joint Commission journal on quality and patient safety</i> , 36(3):133–142.	
409		
410		
411		
412		
413		
414		
	Yanxia Zhang, Jeffrey Olenick, Chu-Hsiang Chang, Steve WJ Kozlowski, and Hayley Hung. 2018. The i in team: Mining personal social interaction routine with topic models from long-term team data. In <i>23rd International Conference on Intelligent User Interfaces</i> , pages 421–426.	
	Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. 2024a. Mlvu: A comprehensive benchmark for multi-task long video understanding. <i>arXiv preprint arXiv:2406.04264</i> .	
	Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. 2024b. Image-of-thought prompting for visual reasoning refinement in multimodal large language models. <i>arXiv preprint arXiv:2405.13872</i> .	

Labels	Behavior Subcodes
Seek	Actively inviting input Expressing uncertainty
Evaluate	Stating a working hypothesis Recapping Explicitly assessing the situation Reasoning
Plan	Stating plans and priorities
Implement	Stating one's ongoing actions Designating tasks

Table 3: Behavior subcodes corresponding to each of our labels. We follow the definition from Schmutz et al. (2021) to determine the subcodes for “Seek”, “Evaluate”, and “Plan”. We add another label of “Implement” given the characteristics of our data source.

A Label Details

Table 3 provides an overview of the behavior subcodes for each label.

Seek includes:

- the action of actively inviting the team members to provide information and share ideas about the current event.
- expressing uncertainty with an implicit invitation to share information.

Evaluate includes:

- a clear formulation of a working hypothesis or diagnosis about the current situation.
- bringing together various pieces of information and providing a summary.
- providing an explicit judgment, giving value to a certain process, information, or strategy. This can be the process of evaluating information that has been gained through seeking information.
- explaining why certain things are more important, or why a specific behavior needs to be done.

Plan refers to laying out the course of action for the next few minutes that needs to contain at least two actions.

Implementation refers to stating the member is conducting the task or delegates a task to another team member.

B Scenario Details

The role of primary anesthesiologist was played by one of the course participants. The surgeon and secondary anesthesiologist (assistant) were played by other course participants. The role of surgeon served as a confederate along with the course instructors. The scenario begins with the primary anesthesiologist taking over the case from one of the course instructors. The patient is receiving general anesthesia and the procedure has already begun. The procedure is complicated by surgical difficulties resulting in the surgeon requesting additional muscle relaxants and increased insufflation pressures. There is also concern that the patient is developing sepsis given the significant gallbladder infection. The patient develops malignant hyperthermia (MH) as the simulated scenario progresses. The primary anesthesiologist must recognize this and begin appropriate treatment. Treatment algorithms for MH are well-known and broadly available (Hopkins et al., 2020; Rosenberg et al., 2020). Definitive treatment includes stopping the triggering agents, administering dantrolene, and supportive care.

C Dataset Information

The total number of anesthesiologists studied was 22; 15(68%) males and 7(32%) females. As part of the Maintenance of Certification in Anesthesiology (MOCA©), anesthesiologists who were board certified after 2000 were required to participate in a simulation course at a simulation center. The participants were board certified anesthesiologists who attended a simulation course at a midwestern academic medical center over a 5 year period. Date of initial certification was obtained from the American Board of Anesthesiologists (ABA) Physician Directory. The study was approved by the Institutional Review Board.

C.1 Physiological Signals

The physiological signals in our dataset include:

SpO2 refers to Peripheral Oxygen Saturation which measures the oxygen saturation level in the blood. Such signal is typically measured through a pulse oximeter.

ECG II refers to Electrocardiogram Lead II which represents the electrical activity of the heart as measured by electrodes placed on the body.

506 **APB** refers to Arterial Blood Pressure which rep- 551
507 represents the pressure exerted by blood on the walls 552
508 of the arteries during the cardiac cycle. 553

509 **HR** refers to Heart Rate which indicates the num- 554
510 ber of heartbeats per minute. 555

511 **NIBP** refers to Non-Invasive Blood Pressure 556
512 which measures blood pressure without the need to 557
513 insert instruments into the body. 558

514 **Temperature** represents the body’s temperature 560
515 and is often measured using a thermometer. 561

516 **Respiratory Waveform** represents the pattern of 562
517 inhalation and exhalation. 563

518 **CO₂** means Carbon Dioxide which typically 564
519 refers to end-tidal CO₂, which represents the con- 565
520 centration of carbon dioxide at the end of an ex- 566
521haled breath. 567

522 **IBP** refers to Invasive Blood Pressure which mea- 568
523 sures blood pressure using invasive techniques, typ- 569
524 ically involving a catheter inserted into an artery or 570
525 vein. 571

526 **C.2 Annotation Details**

527 Two researchers coded six out of 22 randomly se- 572
528 lected data files. The researchers discussed findings 573
529 and resolved discrepancies through the process of 574
530 social moderation. They achieved a Cohen’s kappa 575
531 score of 0.73. The two researchers then independ- 576
532 ently annotated the remaining dataset. 577

533 **D More Details about the Methods**

534 In addition to the methods in Section 4, we have 578
535 a majority vote baseline model which always pre- 579
536 dict the major class. As expected, it reaches a 580
537 decent micro F1 score (55.63) due to the class 581
538 imbalance, while a much lower macro F1 score
539 (14.01). In addition, we test two non-deep learning
540 methods such as RUSBoost (Seiffert et al., 2009)
541 and SMOTE (Chawla et al., 2002) algorithm which
542 is specifically designed to address class imbalance.
543 However, these pre-deep learning methods attains
544 24.21 and 32.32 macro F1 scores, much worse
545 than simply tuning BERT_{base} model or prompting
546 LLMs.

547 **E What Do Medical Professionals Expect** 548 **from NLP?**

549 We are also interested to see how the medical pro-
550 fessionals would view the results we get by em-

ploying these current NLP methods. Therefore, we 551
invite feedbacks from a medical professional who 552
has been working in the domain for over a decade. 553
Here are what we get: 554

1. They see a great opportunity to apply these 555
LLMs on behavioral evaluation in the medical 556
domain. They point out that the current evalua- 557
tion practices in medical domains have signifi- 558
cant limitations (Kolbe and Boos, 2019; Klonek 559
et al., 2019; Stevenson et al., 2022), which typ- 560
ically are labor-intensive and prone to personal 561
biases and errors. They expect NLPers to de- 562
velop consistent, reliable evaluation protocol to 563
give feedback to the healthcare professionals. 564

2. They expect a protocol that can take multimodal 565
input into consideration including the team dia- 566
logue, patient vitals, and procedure videos. We 567
note that this is one of the characteristics for 568
CliniDial. They also hope the NLP system could 569
pinpoint specific teamwork deficiencies in the 570
process. 571

3. They also point out the related NLP methods that 572
they find useful in their domain. For instance, in- 573
tent classification, dialogue summarization, and 574
multimodal reasoning works from NLP can pro- 575
vide quantifiable insights into teamwork dynam- 576
ics and communication patterns in multimodal 577
clinical data (Zhang et al., 2018; Allen et al., 578
2021; Lehmann-Willenbrock and Hung, 2023; 579
Hung et al., 2024). We note that *CliniDial* con- 580
tain rich conversational data. 581