BehaviorSFT: Behavioral Token Conditioning for Clinical Agents Across the Proactivity Spectrum

Anonymous ACL submission

Abstract

Large Language Models (LLMs) as clinical agents require careful behavioral adaptation. While adept at reactive tasks (e.g., diagnosis reasoning), LLMs often struggle with proactive engagement, like unprompted identification of critical missing information or risks. We introduce **BehaviorBench**, a comprehensive dataset to evaluate agent behaviors across a clinical assistance spectrum, ranging from reactive query responses to proactive interventions (e.g., clarifying ambiguities, flagging overlooked critical data). Our BehaviorBench experiments reveal LLMs' inconsistent proactivity. To address this, we propose BehaviorSFT, a novel training strategy using behavioral tokens to explicitly condition LLMs for dynamic behavioral selection along this spectrum. BehaviorSFT boosts performance, achieving up to 97.3% overall Macro F1 on BehaviorBench and improving proactive task scores (e.g., from 95.0% to 96.5% for Qwen2.5-7B-Ins). Crucially, blind clinician evaluations confirmed BehaviorSFT-trained agents exhibit more realistic clinical behavior, striking a superior balance between helpful proactivity (e.g., timely, relevant suggestions) and necessary restraint (e.g., avoiding over-intervention) versus standard fine-tuning or explicit instructed agents.¹

1 Introduction

002

012

017

021

030

040

As Large Language Models (LLMs) transition from experimental systems to deployed agents in clinical environments, a critical question emerges: "*when* should these systems act *reactively* or *proactively* (Fauscette, 2024)?." Unlike general-purpose AI agents, healthcare agents can operate in high-stakes environments where both action and inaction carry significant consequences (Kim et al., 2025). We define *reactive* behaviors as those where the agent responds only to explicit queries with precisely the involve volunteering additional information, raising concerns, or suggesting actions beyond what was directly solicited. Importantly, proactivity in clinical contexts extends beyond merely asking clarifying questions, a common, but limited, focus in existing NLP research (Li et al., 2024; Hu et al., 2024). While question-asking represents one dimension of proactivity, our work encompasses a broader spectrum - including unsolicited intervention, critical evaluation, and recommendation. These behaviors align closely with the "Appraisal" phase of Evidence-Based Medicine (EBM) (Denby, 2008), where clinicians actively assess available information, identify information gaps, and determine appropriate next steps. An agent that remains strictly reactive may fail to raise an alert when problems are observed with critical lab values or medication contraindications (Walter Costa et al., 2021; Wright et al., 2018), potentially compromising patient safety (McCoy et al., 2014). In contrast, an excessively proactive system that frequently interrupts with unsolicited recommendations risks contributing to alert fatigue, interruption of workflow, and potential rejection by healthcare professionals (Sutton et al., 2020). This trade-off between reactive and proactive behaviors forms the core challenge addressed in this paper. The appropriate balance between these modalities varies dramatically based on clinical context, urgency, risk levels, and the specific healthcare roles being augmented, demanding adaptive behavior policy rather than a fixed mode, especially as systems achieve higher levels of autonomy (Figure 4).

information requested, while proactive behaviors

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

To systematically discuss how an agent's reactive and proactive stance should adapt with its increasing capabilities, we adapt the SAE Levels of Driving Automation (SAE, 2021) into a six-level taxonomy for healthcare AI agent autonomy. This framework detailed in Table 8 in Appendix helps to illustrate a key principle: as an AI agent ascends

¹Project Page: https://behavior-adaptation.github.io/



Figure 1: Six representative tasks from BEHAVIORBENCH, showcasing the spectrum of agent behaviors in clinical settings. The figure illustrates (a-c, f) proactive tasks where the LLM agent identifies issues or offers insights without direct prompting, and (b, d, e) reactive tasks responding to explicit clinician queries.

these autonomy levels, its capacity and responsibility to engage in sophisticated proactive behaviors, rather than merely reactive ones, become increasingly critical.

The autonomy level taxonomy highlights that effective healthcare AI, particularly for achieving Level 3 (Conditional Proactive Assistance) and above, must move beyond simple reactive responses (Levels 1-2). As AI autonomy increases, the nature of clinician responsibility evolves, shifting from direct task execution to supervision, validation of AI-driven insights, and management of exceptions. Our work, therefore, focuses on enabling AI agents to learn and exhibit the adapted spectrum of reactive and proactive behaviors crucial for safe and effective operation at these higher levels of conditional and collaborative automation. BEHAVIORBENCH is designed to evaluate these capabilities across this spectrum, and **BehaviorSFT** aims to train agents to achieve this behavioral adaptability, particularly for robust performance at Levels 2 and 3, with an eye towards future capabilities at Level 4.

Effectively adapting *which* of these behaviors is appropriate, and *when*, is essential for clinical AI systems that can safely operate at increasing levels of autonomy. In this work, we ask *what* proactivity means for healthcare AI and how we build systems that are appropriately behaving? To 110 this end, we propose a novel six-level taxonomy 111 for healthcare AI autonomy that maps progression 112 from human-controlled to autonomous operation. 113 We trace the evolution from early reactive systems 114 (Tu et al., 2024; Han et al., 2023) to more recent de-115 velopments like MediQ (Li et al., 2024) and AIME 116 (McDuff et al., 2025; Tu et al., 2024), which in-117 corporate proactive elements while demonstrating 118 the critical interplay between proactivity and ur-119 gency. Our benchmark was curated from real med-120 ical cases sourced from New England Journal of 121 Medicine (NEJM) clinical case reports (Brinkmann 122 et al., 2024). We employed a LLM (Gemini-2.5 123 Flash) to meticulously ground these cases in their 124 factual details and then reformat them into multi-125 turn, multi-clinician-patient conversational scenar-126 ios, integrating multi-modal inputs such as text, 127 images, and tabular data. Indeed, we propose 128 this LLM-assisted methodology for converting ex-129 isting static clinical datasets into rich, reactive-130 proactive benchmark scenarios as a key contribu-131 tion of our work. Additionally, we present a novel 132 training methodology, BehaviorSFT, which em-133 ploys explicit behavioral tokens to condition LLM 134 responses along the reactive-proactive spectrum. 135

136Our approach demonstrates significant improve-137ments, achieving up to 97.3% overall Macro F1 on138BehaviorBench (compared to 96.7% for general139SFT) with particularly notable gains in proactive140tasks (from 95.0% to 96.5%). The primary contri-141butions are:

142

143

144

145

146

147

149

150

151

152

153

154

155

156

157

158

159

160

161

162

164

165

166

- We introduce BEHAVIORBENCH, an evaluation dataset that assesses LLM capabilities across both reactive and proactive tasks in healthcare contexts.
 - We provide detailed analysis of recent LLMs' performance on BEHAVIORBENCH, revealing significant variability in contextual awareness and appropriate behavioral adaptation.
 - 3. We propose BehaviorSFT, a new fine-tuning strategy that leverages behavioral tokens to guide LLMs in dynamically adapting their responses along the reactive-proactive tasks.

2 BEHAVIORBENCH

We introduce BEHAVIORBENCH, a novel dataset specifically designed to assess agent capabilities across the reactive-proactive tasks. Derived from real clinical cases, BEHAVIORBENCH comprises of 6,876 real-world clinical case scenarios from which we derived a total of 142,496 tasks distributed across the 13 distinct task categories. This framework provides a more granular analysis of an agent's ability to discern context and modulate its behavior accordingly, moving beyond standard metrics, such as accuracy, that are solely based on reactive responses. Detailed dataset statistics can be found in the Appendix D.

To ensure that the generated tasks effectively probe clinical reasoning, we construct the dataset in 169 a two-step process. First, we carefully prompt the 170 LLM (see Appendix G) generating the tasks to use 171 detailed summary from real-world clinical cases, 172 including patient history, diagnostics, conversation 173 snippets, and final diagnoses. This ensures that the 174 questions, answers, and rationales reflect genuine 175 clinical context instead of relying on pseudolabels 176 generated without any realistic groundings. All 178 draft tasks then underwent several back-and-forth revision cycles with two physicians, who reviewed 179 any hallucinations and confirmed each scenario's practical plausibility for N=10 cases. Then, to evaluate the agent's proactive capabilities, we augment 182

the base scenarios by intentionally introducing subtle challenges, such as hypothetical scenarios with probable clinical errors, conflicting data points (e.g. modifying numerical values slightly between reports, or presenting exam findings seemingly at odds with imaging), and omitted information expected by clinical standards. The resulting reactiveproactive tasks are as follows:

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

227

228

229

Reactive Tasks evaluates whether the agents can handle information when requested directly.

- 1. fact_retrieval: Finds specific facts mentioned in the text (e.g., "What was the patient's initial temperature?").
- 2. timeline_sequence: Puts events in order using clear time references (e.g., tracing how lung exam findings changed between the initial presentation and Turn N, based on provided descriptions from those time points).
- ddx_reasoning: Explains the reasoning for a possible diagnosis using only the evidence given (e.g., identifying findings prior to Turn M, such as specific X-ray descriptions and sputum results, that suggested bronchopneumonia over simple lobar pneumonia).
- 4. treatment_decision: Connects a doctor's thinking or action to the stated reason or data supporting it (e.g., evaluating a specific diagnostic leaning mentioned in Turn *K* based only on the evidence explicitly available at that time, like sputum results).

Balanced Tasks are initiated by specific, provided information but demand a more significant cognitive step involving deeper thinking, such as multi-step inference, synthesis of multiple data points, or evaluating the impact of new information on existing understanding.

- 1. reasoning_differential_evolution: Compares the patient's situation at two different times and explains how the doctor's assessment should change because of new information (e.g., asking how the list of possible diagnoses should shift from Timepoint A to Timepoint B considering newly available sputum culture results and vital signs).
- 2. integrity_missing_turn_inference: Figures out what was likely said in a missing part of a conversation based on what came

Table 1: Comparison of Public Medical Benchmarks. Modality codes: t=text, i=image, b=tabular/structured data. \checkmark indicates that the benchmark natively supports the evaluation dimension; \varkappa indicates it does not.

Benchmark	Size	Modality	Behavior Evaluation	Sequential Eval.	Dialogue Interaction	Multiple Roles
MedQA (Jin et al., 2021)	1,273	t	×	×	×	X
MedMCQA (Pal et al., 2022)	6,100	t	×	×	×	×
MultiMedQA (Singhal et al., 2023)	13,115	t	×	×	×	×
MediQ (Li et al., 2024)	1,273	t	×	\checkmark	\checkmark	\checkmark
MediQ-AskDocs (Li et al., 2025)	17,000	t	×	\checkmark	\checkmark	\checkmark
ClinicBench (Chen et al., 2024)	11,000	t	×	×	×	×
MedChain (Liu et al., 2024)	12,163	t+i	×	\checkmark	\checkmark	\checkmark
MedAgentBench (Jiang et al., 2025)	300	t+b	×	\checkmark	\checkmark	\checkmark
HealthBench (Arora et al., 2025)	5,000	t	×	×	\checkmark	×
BEHAVIORBENCH (Ours)	142,496	t+i+b	\checkmark	\checkmark	\checkmark	\checkmark

before and after (e.g., "Turn N orders a test, Turn N + M discusses the result. What likely happened in Turn N + K, where 0 < K < M?").

Proactive Tasks require the LLM to use higher-level thinking, and evaluation skills.

233

234

240

241

242

243

244

245

246

247

248

249

252

254

255

260

263

- 1. predictive_next_action: Forecasts the most appropriate subsequent clinical action by integrating the evolving patient case, current symptoms, medical history, and available diagnostic results.
- 2. explicit_error_correction: Identifies and rectifies explicitly stated errors in clinical narratives or proposed actions, providing justifications based on medical knowledge and case specifics (e.g., correcting drug suitability given a patient's allergy).
- 3. omission_detection: Identifies significant omissions in the provided clinical information or documented actions, such as overlooked diagnostic tests or unaddressed critical symptoms that could impact patient care.
- standard_of_care: Assesses whether documented clinical management, including diagnostic procedures and interventions, adheres to established medical guidelines and accepted best practices, often requiring external knowledge.
- 5. interpretation_conflict: Discerns and reconciles nuanced or potentially conflicting interpretations of clinical findings from different sources (e.g., contrasting physical exam notes with radiology findings), articulating their clinical significance.

6. data_conflict_resolution: Identifies direct contradictions or inconsistencies between pieces of factual clinical data presented within a case (e.g., conflicting lab values over time) and proposes logical explanations.

264

265

266

267

268

269

271

272

273

274

275

276

277

278

279

281

282

283

285

286

287

290

291

292

293

294

296

297

298

7. consistency_check: Evaluates the overall logical and clinical coherence of a case narrative or specific information, identifying elements that are incongruous or implausible (e.g., assessing if a patient's reported progression aligns with a given diagnosis).

3 BehaviorSFT: Behavior Adaptation Training

To operationalize the concept of behavioral adaptation within healthcare LLM agents, we propose a targeted training strategy, Behavior-Conditioned Supervised Fine-Tuning (BehaviorSFT). This approach leverages our specialized BehaviorBench dataset (Section 2) to explicitly teach LLMs to modulate their responses along the reactive-proactive spectrum based on inferred clinical context. This contrasts with standard SFT approaches, which typically optimize for task completion without explicit mechanisms to control the agent's level of initiative or caution, risking either unsafe passivity or disruptive over-intervention.

3.1 Behavior Tokens

Rationale for Prefix Tokens: We employ prefix behavior tokens (e.g., <reactive>, <proactive>) for several reasons. Placing the token at the beginning of the target sequence allows it to act as a direct control signal, conditioning the entire generation process on the desired behavioral mode from the outset. This explicitly trains the model to adopt the appropriate style, tone, and level of initiative as it 299generates the response. While one could consider300predicting the token after some internal reasoning301chain, our approach integrates this reasoning im-302plicitly, i.e., the model learns to predict the correct303initial token based on its understanding of the input304context (x), as described in our Contextual Behav-305ior Assessment capability (Section 3.3). This pro-306vides an end-to-end mechanism for context-aware,307behaviorally adapted generation. Central to our308approach is the introduction of special behavior309tokens paired with the target response during train-310ing.

311

312

313

314

315

316

317

318

319

332

333

334

335

337

341

346

- <reactive>: Signals the generation of a direct, concise response strictly adhering to the explicit query, avoiding unsolicited information or inferences.
- <proactive>: Signals a response that may include identifying implicit issues, volunteering relevant context or warnings, suggesting next steps, or applying external knowledge (e.g., standards of care) beyond the literal query.

These tokens act as control signals, learned by the 320 model and conditioning the subsequent generation 321 process. Alternative approaches exist, such as training a separate classifier to select the mode and then routing the input to specialized reactive or proactive models, or using inference-time techniques like thresholding logits associated with the behavior to-326 kens for finer control. However, our BehaviorSFT approach offers a simpler, unified training process 328 within a single model. Future work could explore hybrid methods or compare the efficacy of these different control paradigms. 331

3.2 Training Data

BehaviorBench serves as the crucial training ground for BehaviorSFT. Each instance within the benchmark's training split is meticulously annotated with the desired target behavior token based on the task's nature and the underlying clinical scenario's demands:

- 1. **Reactive Annotation (<reactive>):** Applied to tasks demanding factual recall, direct sequencing, or simple reasoning strictly from provided data (e.g., fact_retrieval, timeline_sequence).
- 2. **Proactive Annotation** (<proactive>): Applied to tasks necessitating critical assessment, error/omission detec-

tion, consistency checking, or prediction based on clinical standards (e.g., consistency_check, standard_of_care, predictive_next_action). 347

348

349

351

352

353

354

355

357

358

360

361

362

363

376

377

378

379

380

381

383

384

386

391

3. Contextual Annotation for Balanced Tasks: Instances from balanced tasks (e.g., reasoning_differential_evolution) are annotated based on whether the specific context warrants simple reporting (<reactive>) or highlighting significant changes/implications (<proactive>).

Each annotated instance is then structured for auto-regressive SFT, pairing the input context/query with a target sequence beginning with the assigned behavior token, followed by an ideal response exemplifying that behavior. Example 1 (Reactive Task):

Input: Context: [Note excerpt: Vitals 364 stable.] 365 Query: Latest vitals? Target: <reactive> BP 120/80, HR 75, 367 Temp 37.0C, RR 16. 368 Example 2 (Proactive Task): 369 Input: Context: [Chart: Rx Drug A. Allergy list: Drug A.] 371 Query: Confirm med list okay? 372 Target: <proactive> Warning: Drug A 373 prescribed but patient is allergic. 374 Review immediately. 375

This structured data format explicitly teaches the model the association between clinical scenarios, appropriate behavioral modes (reactive/proactive), and corresponding linguistic outputs.

3.3 Training Procedure: BehaviorSFT

Starting with a pre-trained foundation LLM, we perform SFT using the behavior-annotated BehaviorBench training data. The objective is the standard causal language modeling loss, minimizing the negative log-likelihood of the target sequence $y = (y_1, ..., y_T)$, where $y_1 \in \{ < \text{reactive}, < \text{proactive} \} \}$:

$$\mathcal{L}_{BehaviorSFT} = -\sum_{i=1}^{T} \log P(y_i | y_{\le i}, x; \theta) \quad (1)$$

Here, x is the input context/query, $y_{<i}$ are the preceding target tokens, and θ represents the model parameters.

402

403

Through this process, the model learns the crucial, intertwined capabilities:

1. Contextual Behavior Assessment: Implicitly analyzing the input x to determine the likelihood that a proactive or reactive stance is warranted, influencing the prediction of the initial token y_1 .

2. Behavior-Conditioned Generation: Generating subsequent tokens $y_{2:T}$ in a manner consistent with the generated or given behavior token y_1 , adopting the appropriate style, tone, and level of detail or intervention.



Figure 2: Density distributions of (I) Specificity and (II) Implicitness scores for Baseline, BehaviorSFT, and GeneralSFT agent outputs. (I) Specificity: Both fine-tuned models (BehaviorSFT and GeneralSFT) markedly improve output specificity over the Baseline, with distributions concentrated at high scores (\sim 0.9). (II) Implicitness: Distinct implicitness profiles emerge: GeneralSFT is the most explicit (lowest scores, \sim 0.6-0.7), the Baseline is the most implicit (highest scores, \sim 0.7-0.9), while BehaviorSFT exhibits a moderate, intermediate level of implicitness (\sim 0.7-0.8).



Figure 3: G-Eval with gpt-4o-mini as evaluator of Qwen-2.5-7B-Ins responses across four key metrics. We compare the average scores for the Baseline model, our proposed BehaviorSFT, and GeneralSFT. BehaviorSFT consistently outperforms the Baseline across all metrics and demonstrates competitive or superior performance compared to GeneralSFT.

4 Experiments and Results

4.1 Setup

All experiments **BEHAVIORBENCH** use with а fixed 6776/110/977 train-val-test split. We fine-tune both backbones; Qwen-2.5-7B-Instruct (Team, 2024) and Meta-Llama-3.1-8B-Instruct (Meta AI, 2024). Details implementation details can be found in Appendix H.

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

4.2 Main Results

From Reactive to Proactive capabilities in clinical LLMs involve processing and responding directly to explicitly provided information. Reactivity encompasses *fact retrieval*, *information summarization*, ordering events via *direct sequencing*, following *simple execution* instructions, and performing *basic reasoning from explicit data*, these tasks test the LLM's ability to understand and manipulate information as presented, without significant inference or applying external knowledge. The Proactive-Reactive Scale of 0.0-0.4 typically reflects these functions.

Conversely, require the LLM to transcend literal interpretation, demonstrating deeper reasoning, anticipation, and critical assessment. Key aspects include *inference and implication* (identifying unstated assumptions or missing information), *anticipation and prediction* (foreseeing next steps or complications), *consistency and conflict detection* (finding discrepancies between data points), *error recognition and correction, applying external knowledge* like standards of care, and *synthesis and complex interpretation* from multiple sources. These tasks simulate higher-order clinical thinking. The Proactive-Reactive Scale of 0.6-1.0 aligns with these skills, while 0.4-0.6 represents a balance.

Empirical Results Overview. Table 2 reports Macro F1 scores across the three task categories. Relative to both the majority-voting **Ensemble** baseline and standard supervised fine-tuning (**Gen. SFT**), **BehaviorSFT** matches or slightly exceeds performance on the *Reactive* and *Balanced* sets, and yields a clear advantage on the most demanding *Proactive* tasks (Qwen: 96.5% vs. 95.0%; Llama: 94.7% vs. 94.2%). These gains confirm that the behavior-aligned fine-tuning strategy is particularly effective for higher-order reasoning tasks such as complex inference, error correction, and guideline-based decision making, thereby strengthening the model's proactive capabilities. Detailed

Table 2: Performance on BEHAVIORBENCH. We report Macro F1-scores (%) across three task categories. Best result per task is highlighted in **bold**. The Ensemble column reports baseline performance by majority voting across three commercial closed-source models (Gemini-2.5-pro, OpenAI-o1, DeepSeek-R1). 'ZS' = Zero-Shot, 'FS (k=3)' = Few-Shot (3 examples), 'CoT' = Chain-of-Thought, 'Explicit Instr.' = ZS with explicit reactive/proactive instruction, 'Gen. SFT' = Standard Supervised Fine-Tuning (SFT), 'BehaviorSFT' = Our proposed fine-tuning method.

Category	Task		Ens	emble	Qwen	2.5-7B-Ins	Llama	3.1-8B-Ins
		ZS	FS (k=3)	ZS + Explicit Instr.	Gen. SFT	BehaviorSFT	Gen. SFT	BehaviorSFT
e	fact_retrieval	100.0	100.0	100.0	100.0	100.0	100.0	100.0
tiv	timeline_sequence	100.0	100.0	100.0	100.0	100.0	100.0	100.0
eac	ddx_reasoning	96.2	96.6	96.6	96.1	96.1	94.2	92.7
2	treatment_decision	94.8	95.3	95.3	100.0	98.4	98.4	98.7
	Avg.	98.2	98.2	98.2	98.6	98.6	97.8	97.2
ed	reasoning_diff_evolution	98.6	98.6	98.6	100.0	100.0	100.0	100.0
Balance	integrity_missing_turn	100.0	100.0	100.0	100.0	100.0	96.4	100.0
	Avg.	97.2	97.6	97.6	100.0	99.2	98.5	100.0
	consistency_check	94.3	100.0	94.3	100.0	100.0	100.0	100.0
	data_conflict_resolution	97.2	97.2	97.2	99.3	98.6	99.2	98.6
ive	interpretation_conflict	98.5	96.5	96.5	96.6	96.6	98.5	98.6
act	standard_of_care	93.4	95.3	93.7	94.8	93.3	91.5	88.4
Proac	omission_detection	89.5	92.4	89.3	88.5	95.1	90.0	93.2
	explicit_error_correction	96.3	97.5	96.4	98.3	99.2	98.4	97.2
	predictive_next_action	82.5	83.0	82.3	84.8	91.7	77.0	83.4
	Avg.	94.3	95.1	94.0	95.0	96.5	94.2	94.7
Avg.		95.4	96.0	95.3	96.7	97.3	95.8	96.1

Table 3: Macro F1-scores of prompting methods on behavior classification. Method abbreviations: $\mathbf{BT} = \mathbf{Behavior}$ token, BC = Behavior chain-of-thought, OC = Option CoT, OP = Option. Class abbreviations: Five-class (BA = balanced; H_PR = highly_proactive; H_RE = highly_reactive; P_PR = primarily_proactive; P_RE = primarily_reactive), **Binary** (PR = proactive; N_PR = non-proactive), **Three-class** (BA = balanced; PR = proactive; RE = reactive).

		Five-class				Bir	Binary		Three-class	
	BA	H_PR	H_RE	P_PR	P_RE	PR	N_PR	BA	PR	RE
BT-OC-OP	42.62	89.47	4.76	19.19	68.72	82.14	92.10	53.41	92.10	73.68
BT-OP	37.06	87.77	13.79	25.28	66.40	82.76	92.19	46.92	92.19	66.42
BT-BC-OC-OP	58.24	87.84	19.05	11.82	71.75	83.48	92.90	51.67	92.90	72.09
BT-BC-OP	54.74	88.89	17.39	11.00	73.68	82.97	92.58	51.76	92.58	69.57
BC-BT-OC-OP	57.06	87.73	14.81	7.07	74.89	82.59	92.23	45.00	92.32	69.96

accuracy figures for the three commercial baselines are provided in Appendix F.

454

455

456

457

458

459

460

461

462

463

464

465

466

Enhanced User-Centric Qualities with G-Evaluation Our evaluation using G-Eval (Liu et al., 2023), a methodology leveraging large models for human-aligned assessment, reveals significant qualitative improvements with BehaviorSFT. As depicted in Figure 3, BehaviorSFT consistently outperforms the Baseline across all four key metrics: Utility, Safety, Clarity, and Behavioral Appropriateness. Notably, BehaviorSFT achieves the highest scores in Utility (0.95 vs. 0.93 for GeneralSFT and 0.90 for Baseline), Clarity (0.94 vs. 0.92 for GeneralSFT and 0.88 for Baseline), and Behavioral Appropriateness (0.91 vs. 0.87 for GeneralSFT and 0.86 for Baseline). While GeneralSFT scores marginally higher in Safety (0.97 vs. 0.95 for BehaviorSFT), BehaviorSFT still demonstrates a strong safety profile. These results underscore BehaviorSFT's capability to not only perform tasks effectively but also to align more closely with user expectations in terms of usefulness, understandability, and appropriate interaction, suggesting a more refined and user-centric agent behavior.

Optimizing Output Specificity while Balancing Implicitness Figure 2 illustrates the impact of 477

478

479

467

468

our fine-tuning approaches on the nuanced char-480 acteristics of agent responses, specifically their 481 specificity and implicitness. Both fine-tuned mod-482 els, BehaviorSFT and GeneralSFT, markedly en-483 hance output specificity compared to the Baseline, 484 with distributions concentrating at high specificity 485 scores (around 0.9). This indicates that both meth-486 ods generate more detailed and precise information. 487 However, a key distinction emerges in their implic-488 itness profiles. GeneralSFT tends towards more 489 explicit communication, reflected in lower implic-490 itness scores (approximately 0.6-0.7). In contrast, 491 the Baseline model is the most implicit (scores 492 around 0.7-0.9). BehaviorSFT carves out an in-493 termediate and potentially more versatile profile, 494 achieving a moderate level of implicitness (scores 495 approximately 0.7-0.8). This suggests that Behav-496 iorSFT can deliver highly specific information with-497 out resorting to excessive explicitness, potentially 498 mirroring more natural human communication pat-499 terns and aligning with the idea that effective agents must navigate implicit evaluation criteria (Wadhwa et al., 2025). 502

4.3 Ablation on prompting variants for Behavior Pattern Analysis

Table 3 evaluates five prompting recipes obtained by incrementally adding *Behavior Chainof-Thought* (BC) and *Option reasoning* (OC/OP) on top of the *Behavior Token* (BT) baseline. The full recipe *BT–BC–OC–OP* achieves the best or second-best Macro F1 in 11 of the 13 columns (e.g., *Five-class BA* 58.2 and *Binary PR* 83.5), showing that BC and OC/OP provide complementary gains. Dropping OC/OP (*BT–BC–OP*) or BC (*BT–OP*) consistently lowers scores, while reversing the BC placement (*BC–BT–OC–OP*) yields a smaller benefit, indicating that BC is most effective when appended after the BT prompt. Overall, combining both reasoning cues delivers the most robust behaviour classification across all label granularities.

5 Conclusion

503

504

505

507

510

511

512

513

514

516

517

518

519

520

521This paper addresses the critical gap in LLM proac-522tivity for healthcare. Our BEHAVIORBENCH, vali-523dated by clinicians for plausibility, systematically524evaluates this, revealing LLM deficiencies in proac-525tive reasoning despite reactive strengths. We intro-526duced BehaviorSFT, a new fine-tuning method us-527ing explicit <reactive> and <proactive> tokens. Be-528haviorSFT improved performance, achieving up to

97.3% overall Macro F1 on BEHAVIORBENCH and boosting proactive task scores (e.g., Qwen2.5-7B-Ins from **95.0% to 96.5%**). Crucially, in a clinician user study, BehaviorSFT-trained agents received the most favorable rankings (best mean rank **1.80**). G-Eval results also showed superior Utility (**0.95**) and Behavioral Appropriateness (**0.91**). These combined findings demonstrate BehaviorSFT's effectiveness in creating more reliable, clinically nuanced, and clinician-preferred LLM agents for complex healthcare scenarios.

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

References

- 2021. Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles. Surface Vehicle Recommended Practice.
- Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, and 1 others. 2025.
 Healthbench: Evaluating large language models towards improved human health. *arXiv preprint arXiv:2505.08775*.
- Rory Brinkmann, Eric Rosenberg, David N Louis, and Scott H Podolsky. 2024. Building a community of medical learning—a century of case records of the massachusetts general hospital in the journal.
- Canyu Chen, Jian Yu, Shan Chen, Che Liu, Zhongwei Wan, Danielle Bitterman, Fei Wang, and Kai Shu. 2024. Clinicalbench: Can llms beat traditional ml models in clinical prediction? *arXiv preprint arXiv:2411.06469*.
- Eleni Chiou, Francesco Giganti, Shonit Punwani, Iasonas Kokkinos, and Eleftheria Panagiotaki. 2020. Harnessing uncertainty in domain adaptation for mri prostate lesion segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI* 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23, pages 510–520. Springer.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- Donald J Denby. 2008. *Evidence based medicine*. Xulon Press.
- Michael Fauscette. 2024. Agentic ai vs. llms: Understanding the shift from reactive to proactive ai. https://www.arionresearch.com/blog/. Arion Research Blog.
- Alexander Fixler, Blake Oliaro, Marshall Frieden, Christopher Girardo, Fiona A Winterbottom, Lisa B

683

684

685

686

687

688

689

690

637

638

Fort, and Jason Hill. 2023. Alert to action: implementing artificial intelligence–driven clinical decision support tools for sepsis. *Ochsner Journal*, 23(3):222–231.
Stephen H Friend, Geoffrey S Ginsburg, and Rosalind W Picard. 2023. Wearable digital health technology.

581

582

588

593

594

596

597

610

611

613

616

618

625

627

631

- Illin Gani, Ian Litchfield, David Shukla, Gayathri Delanerolle, Neil Cockburn, and Anna Pathmanathan. 2025. Understanding "alert fatigue" in primary care: Qualitative systematic review of general practitioners attitudes and experiences of clinical alerts, prompts, and reminders. *Journal of Medical Internet Research*, 27:e62763.
- Senay A Gebreab, Khaled Salah, Raja Jayaraman, Muhammad Habib ur Rehman, and Samer Ellaham.
 2024. Llm-based framework for administrative task automation in healthcare. In 2024 12th International Symposium on Digital Forensics and Security (IS-DFS), pages 1–7. IEEE.
- Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, and 1 others. 2025. Towards an ai coscientist. arXiv preprint arXiv:2502.18864.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2023. Think before you speak: Training language models with pause tokens. *arXiv preprint arXiv:2310.02226*.
- Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressem. 2023.
 Medalpaca–an open-source collection of medical conversational ai models and training data. *arXiv* preprint arXiv:2304.08247.
- Zhiyuan Hu, Chumin Liu, Xidong Feng, Yilun Zhao, See-Kiong Ng, Anh Tuan Luu, Junxian He, Pang Wei Koh, and Bryan Hooi. 2024. Uncertainty of thoughts: Uncertainty-aware planning enhances information seeking in large language models. arXiv preprint arXiv:2402.03271.
- Mustafa I Hussain, Tera L Reynolds, and Kai Zheng. 2019. Medication safety alert fatigue may be reduced via interaction design and clinical role tailoring: a systematic review. *Journal of the American Medical Informatics Association*, 26(10):1141–1149.
- Yixing Jiang, Kameron C Black, Gloria Geng, Danny Park, Andrew Y Ng, and Jonathan H Chen. 2025. Medagentbench: Dataset for benchmarking llms as agents in medical applications. *arXiv preprint arXiv:2501.14654*.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*.
- Mohamed Khalifa and Mona Albadawy. 2024. Artificial intelligence for clinical prediction: exploring key domains and essential functions. *Computer Methods and Programs in Biomedicine Update*, page 100148.
- Yubin Kim, Hyewon Jeong, Shen Chen, Shuyue Stella Li, Mingyu Lu, Kumail Alhamoud, Jimin Mun, Cristina Grau, Minseok Jung, Rodrigo R Gameiro, and 1 others. 2025. Medical hallucination in foundation models and their impact on healthcare. *medRxiv*, pages 2025–02.
- Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won Park. 2024. Mdagents: An adaptive collaboration of Ilms in medical decision making.
- Eva K Lee, Tsung-Lin Wu, Tal Senior, and James Jose. 2014. Medical alert management: a real-time adaptive decision support tool to reduce alert fatigue. In *AMIA Annual Symposium Proceedings*, volume 2014, page 845.
- Shuyue Stella Li, Vidhisha Balachandran, Shangbin Feng, Jonathan Ilgen, Emma Pierson, Pang Wei W Koh, and Yulia Tsvetkov. 2024. Mediq: Questionasking llms and a benchmark for reliable interactive clinical reasoning. *Advances in Neural Information Processing Systems*, 37:28858–28888.
- Shuyue Stella Li, Jimin Mun, Faeze Brahman, Jonathan S Ilgen, Yulia Tsvetkov, and Maarten Sap. 2025. Aligning llms to ask good questions a case study in clinical reasoning. *arXiv preprint arXiv:2502.14860*.
- Jie Liu, Wenxuan Wang, Zizhan Ma, Guolin Huang, Yihang SU, Kao-Jung Chang, Wenting Chen, Haoliang Li, Linlin Shen, and Michael Lyu. 2024. Medchain: Bridging the gap between llm agents and clinical practice through interactive sequential benchmarking. *arXiv preprint arXiv:2412.01605*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*.
- Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen, Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong, Zhong Zhang, Yankai Lin, and 1 others. 2024. Proactive agent: Shifting llm agents from reactive responses to active assistance. arXiv preprint arXiv:2410.12361.
- Arjun Mahajan, Kimia Heydari, and Dylan Powell. 2025. Wearable ai to enhance patient safety and clinical decision-making. *npj Digital Medicine*, 8(1):176.

- 692
- 693 694
- 69
- 69
- 69
- 700
- 701 702

- 704 705
- 7(7(7(7(
- 710 711
- 712 713 714
- 715 716
- 717 718 719 720
- 721 722 723
- 725 726 727

7

7

- 7
- 733
- 7

7

737

73

739 740

741 742 743

- Allison B McCoy, Eric J Thomas, Marie Krousel-Wood, and Dean F Sittig. 2014. Clinical decision support alert appropriateness: a review and proposal for improvement. *Ochsner journal*, 14(2):195–202.
- Daniel McDuff, Mike Schaekermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, and 1 others. 2025. Towards accurate differential diagnosis with large language models. *Nature*, pages 1–7.
- Meta AI. 2024. Meta-Llama-3.1-8B-Instruct. https://huggingface.co/meta-llama/Llama-3. 1-8B-Instruct. Llama 3.1 Community License. Accessed 20 May 2025.
- Olufisayo Olusegun Olakotan and Maryati Mohd Yusof. 2020. Evaluating the alert appropriateness of clinical decision support systems in supporting clinical workflow. *Journal of biomedical informatics*, 106:103453.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pages 248–260. PMLR.
- Tahmina Nasrin Poly, Md Mohaimenul Islam, Muhammad Solihuddin Muhtar, Hsuan-Chia Yang, Phung Anh Nguyen, and Yu-Chuan Li. 2020. Machine learning approach to reduce alert fatigue using a disease medication–related clinical decision support system: model development and validation. *JMIR medical informatics*, 8(11):e19489.
- Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, and 1 others. 2024. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.
- Reed T Sutton, David Pincock, Daniel C Baumgart, Daniel C Sadowski, Richard N Fedorak, and Karen I Kroeker. 2020. An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ digital medicine*, 3(1):17.

Qwen Team. 2024. Qwen2.5: A party of foundation models.

Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, and 1 others. 2024. Towards conversational diagnostic ai. arXiv preprint arXiv:2401.05654. Bethany A Van Dort, Wu Yi Zheng, Vivek Sundar, and Melissa T Baysari. 2021. Optimizing clinical decision support alerts in electronic medical records: a systematic review of reported strategies adopted by hospitals. *Journal of the American Medical Informatics Association*, 28(1):177–183.

744

745

747

750

751

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

789

790

791

792

793

794

795

796

797

- Manya Wadhwa, Zayne Sprague, Chaitanya Malaviya, Philippe Laban, Junyi Jessy Li, and Greg Durrett. 2025. Evalagent: Discovering implicit evaluation criteria from the web. *arXiv preprint arXiv:2504.15219*.
- Maria Beatriz Walter Costa, Mark Wernsdorfer, Alexander Kehrer, Markus Voigt, Carina Cundius, Martin Federbusch, Felix Eckelt, Johannes Remmler, Maria Schmidt, Sarah Pehnke, and 1 others. 2021. The clinical decision support system ampel for laboratory diagnostics: implementation and technical evaluation. *JMIR Medical Informatics*, 9(6):e20407.
- R Jay Widmer, Nerissa M Collins, C Scott Collins, Colin P West, Lilach O Lerman, and Amir Lerman. 2015. Digital health interventions for the prevention of cardiovascular disease: a systematic review and meta-analysis. In *Mayo Clinic Proceedings*, volume 90, pages 469–480. Elsevier.
- Adam Wright, Skye Aaron, Diane L Seger, Lipika Samal, Gordon D Schiff, and David W Bates. 2018. Reduced effectiveness of interruptive drug-drug interaction alerts after conversion to a commercial electronic health record. *Journal of General Internal Medicine*, 33:1868–1876.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*.
- Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha. 2025. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*.
- Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. Linkbert: Pretraining language models with document links. *arXiv preprint arXiv:2203.15827*.

A Related Works

The Evolving Role of AI in Clinical Tasks Early AI applications in healthcare predominantly functioned as reactive tools, such as information retrieval systems responding to explicit queries (Yasunaga et al., 2022) or basic clinical decision support (CDS) systems triggering alerts based on predefined rules. These systems, while valuable, often lacked contextual understanding and the ability to anticipate clinician needs or potential issues proactively (McCoy et al., 2014; Sutton et al.,

2020). More recent advancements, particularly 798 with LLMs, have paved the way for more sophisticated AI assistants. Models like Med-PaLM (Sing-800 hal et al., 2023) and Med-Alpaca (Han et al., 2023) demonstrated strong domain knowledge, though primarily in a reactive question-answering capac-803 ity. The trend is now shifting towards systems 804 with proactive capabilities. For instance, MediQ (Li et al., 2024) explores proactive informationseeking when context is incomplete, while systems 807 like AIME (Tu et al., 2024) and MDAgents (Kim et al., 2024) begin to suggest next steps or anticipate patient needs. This evolution mirrors broader 810 trends in mixed-initiative interaction design, where 811 AI systems dynamically share control with users 812 (). Our work builds on this trajectory by focusing on systematically training and evaluating the 814 adaptation of reactive and proactive behaviors. 815

Challenges of Proactive AI in Healthcare 816 Proactive behaviors in healthcare AI are diverse 817 and critical. One key form is proactive alerting, where systems identify and flag critical information, 819 potential errors (e.g., drug interactions, missed stan-821 dard protocols), or deviations from normal (e.g., critical lab values) (Wright et al., 2018; Fixler et al., 2023; Lee et al., 2014). While potentially life-823 saving, a major challenge is alert fatigue, where excessive or irrelevant alerts lead to high override 825 rates and desensitization among clinicians (Gani et al., 2025; Olakotan and Yusof, 2020; Hussain 827 et al., 2019). Recent efforts focus on contextualizing alerts to improve relevance and reduce fatigue (Poly et al., 2020; Van Dort et al., 2021). Another 831 crucial area is proactive information-seeking under uncertainty. Clinical scenarios often involve 832 incomplete information, and an AI agent should 833 ideally recognize knowledge gaps and ask clarifying questions rather than proceeding with potentially unsafe assumptions (Li et al., 2024). Frame-836 works like ALFA (Li et al., 2025) use psychology-837 informed approaches, and methods like Uncertainty 838 of Thoughts (UoT) (Hu et al., 2024) leverage uncertainty estimation to guide information acquisition. 840 This contrasts with agents that might fail to alert on critical missing information (Kim et al., 2025). 842 Finally, contextual intervention and suggestion in-844 volve AI volunteering relevant, unprompted information, suggesting next steps, or adapting guidance 845 based on inferred clinical context, user expertise, or workflow stage (Widmer et al., 2015; Friend et al., 2023; Mahajan et al., 2025; Khalifa and Albadawy,

2024). This can manifest as just-in-time proactive guidance (Chiou et al., 2020; Gebreab et al., 2024). The core challenge, which our work directly addresses, is adapting *when* and *how* to intervene to be helpful without being disruptive or unsafe (?).

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

Controllable Generation for Healthcare LLMs Controlling the behavior of LLMs beyond simple task completion is an active research area. Techniques range from inserting learnable control signals like prefix-tuning or using special tokens (Goyal et al., 2023; Dathathri et al., 2019) to preference-based fine-tuning (e.g., RLHF) to encourage specific interaction styles (). Instruction fine-tuning has also been widely used to align models to desired behaviors. While these methods offer general control, their application to the nuanced reactive-proactive spectrum in high-stakes domains like healthcare requires domain-specific data and evaluation. Several benchmarks exist for evaluating LLMs in medicine, such as MedQA (Jin et al., 2021), PubMedQA (Jin et al., 2019), MedMCQA (Pal et al., 2022), and more recent ones like MedAgentBench (Jiang et al., 2025) or ClinicBench (Chen et al., 2024). These primarily focus on knowledge accuracy, reasoning over medical facts, or agentic task completion. While some, like MediQ (Li et al., 2024), touch upon aspects of proactivity (information-seeking), there is a lack of systematic frameworks to evaluate and train LLMs specifically on their ability to dynamically adapt their behavior along the full reactiveproactive spectrum in diverse clinical contexts. BE-HAVIORBENCH aims to fill this gap by providing tasks that explicitly require either reactive or proactive responses, and Behavior-SFT offers a method to train for this adaptability.

B Limitations and Future Works

Data & Task Scope. BEHAVIORBENCH aggregates 6,876 English clinical vignettes (142K task instances) from *NEJM*. This corpus reflects an internal-medicine bias and omits modalities such as radiology reads, nursing shift notes, tele-health transcripts, and non-English documentation. The future tasks include expanding the benchmark to multilingual EHR snippets and image-grounded prompts, and we are adding tasks for dermatology, psychiatry, and longitudinal trend summarisation to test whether proactive cues generalise beyond text-only, single-visit encounters.

Behaviour Modelling. Our BEHAVIORSFT controller currently toggles generation with a binary 899 <reactive> / <proactive> token. Although ef-900 fective for coarse behaviour shifts, this switch can-901 not express nuances such as anticipatory clarifi-902 cation versus high-urgency escalation, and it oc-903 casionally over-fires, creating alert fatigue. We 904 are experimenting with a hierarchical token in-905 ventory (e.g. <clarify_info>, <flag_safety>, 906 <escalate_critical>) learnt from multi-label su-907 pervision, and with behaviour-weighted RLHF that 908 continuously trades helpfulness against cognitive 909 load. 910

Evaluation & Deployment Readiness. The clin-911 912 ician study in Appendix I involves three medical doctors number of cases sufficient for valida-913 tion but under-powered for robust error stratifica-914 tion or workflow integration. Future work should 915 recruit multi-institution cohorts (20+ clinicians, 916 1,000+ cases) and embeds the agent inside a sim-917 ulated EHR sandbox to observe interrupt patterns, 918 hand-off continuity, and long-horizon reasoning 919 across multi-day episodes.

C Ethical Implications

921

922

923

925

927

928

929

931

934

936

937

941

Safety & Accountability. Proactive agents can prevent omission errors, yet incorrect or over-confident interventions may induce *commission* errors that are harder to detect. We therefore plan to release model checkpoints after careful reviews. Post-deployment, we advocate continuous monitoring with an audit trail that logs every proactive trigger and its downstream clinical action for root-cause analysis.

Fairness & Bias Mitigation. Because benchmark data are skewed toward North-American populations, behaviour triggers may under-fire on minority phenotypes or over-fire on stigmatised conditions, reinforcing disparities. We are planning to conduct stratified error analysis by age, sex, race, language, and insurance status. Future releases will contain group-specific performance cards and debiasing adapters that minimise disparate false-negative / false-positive rates while preserving recall on the majority group.

942Data Privacy & Responsible Release.All med-943ical cases are available for those institutions who944purchased NEJM license; nonetheless, fine-tuned945models might memorize private strings when946trained on institutional EHRs. We will publish

an **Ethical Usage Card** outlining intended tasks, known failure modes, monitoring hooks, and sunset clauses for model retirement, and we encourage downstream users to adopt the same safeguards.

D Dataset Statistics

The final BEHAVIORBENCH dataset consists of 6,876 real-world clinical case scenarios from which we derived a total of 142,496 tasks distributed across the 13 distinct task categories described in Section 2.

D.1 Simulated Conversations

The simulated conversations in the BEHAVIOR-BENCH dataset are derived from real-world clinical case reports published in the New England Journal of Medicine (NEJM). Each conversation reconstructs the clinical reasoning process among healthcare professionals, encompassing diagnostic deliberation, treatment planning, and communication with patients and caregivers.

Table 4 and Figure 6 and 7 provide descriptive statistics of the conversation data, illustrating the natural variability and complexity of the simulated dialogues. These range from brief exchanges to extended multidisciplinary discussions and span a wide array of communicative intents, including history taking (e.g., eliciting chief complaint, symptom duration, and past medical history), physical examination interpretation, diagnostic reasoning, and family updates. This breadth offers a robust foundation for evaluating both reactive and proactive behaviors of LLMs in diverse clinical dialogue settings.

Table 4: **Summary Statistics of Simulated Clinical Conversations**. This table reports average structural properties of the conversations in the dataset, including the number of dialogue turns, total dialogue length in characters, and number of unique participants per case.

Metric	Value
Avg. # of turns per conversation	33.3
Avg. len of dialogue per conversation	6194.3
Avg. # of participants per case	8.7

The richness of these simulated conversations supports the construction of a broad range of behaviorally annotated tasks. These tasks underpin our evaluation framework, which is designed to assess not only reactive capabilities, such as information

983

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

991

995

997

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1028

1030

1032

retrieval, but also proactive competencies such as anticipatory reasoning and clinical foresight.

D.2 Tasks

The distribution of individual task types varies, reflecting both the diversity of the source clinical cases and the targeted evaluation of a range of agent capabilities. Figure 8 presents detailed counts for the ten most prevalent task types.

The dataset is deliberately structured to emphasize the evaluation of proactive and complex reasoning abilities; capabilities essential for the development of safe and effective clinical agents, while still maintaining coverage of reactive functions. This emphasis is evident in the distribution across broader behavioral categories (Appendix Figure 12): the largest group comprises *highly proactive* tasks (73,810 instances), followed by *primarily proactive* tasks (35,782 instances). *Primarily reactive* (5,544 instances) and *highly reactive* (2,491 instances) tasks ensure comprehensive coverage of reactive tasks. Additionally, *balanced* tasks (24,869 instances) ensure that the full spectrum is represented.

We also categorize tasks by complexity, broadly distinguishing between 'intermediate' tasks (often corresponding to simpler reactive functions) and 'advanced' tasks (typically involving proactive or complex balanced reasoning). The dataset heavily features 'advanced' tasks (127,927 instances) compared to 'intermediate' tasks (14,569 instances), as shown in Figure 9, where the advanced tasks feature a higher proactive score of above 0.8 compared to intermediate tasks with an average of 0.4 proactive score (Figure 10 in Appendix).

Furthermore, a continuous behavior score (ranging from 0.0 for fully reactive to 1.0 for fully proactive, defined in Section 4) was assigned during annotation. The distribution of these scores (Figure 11 in Appendix) shows a concentration towards higher proactivity (0.6-1.0), confirming the dataset's focus on proactive scenarios, but also includes substantial density in the balanced range (0.4-0.6) and coverage of reactive cases (0.0-0.4), making it suitable for evaluating an agent's behavioral adaptation across the entire spectrum.

E The Evolving Landscape of Healthcare AI

The capabilities of Artificial Intelligence (AI) systems in healthcare are rapidly advancing, moving beyond simple information retrieval towards more autonomous and complex task handling. Figure 4 provides a visual representation of this evolving landscape, positioning various contemporary Healthcare AI Systems and Enabling Frameworks/Concepts based on two key dimensions: their operational Task Scope and their level of System Autonomy. 1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1076

1077

1078

1079

1080

1081

1082

1083

1084

The System Autonomy axis is rigorously grounded in the Six-Level Taxonomy for Healthcare AI Agent Autonomy (detailed in Table 8 in the Appendix). This taxonomy delineates capabilities from Level 0-1 (No Automation/Clinician Assistance), where AI provides reactive information or simple alerts, through Level 2 (Partial Automation/Reactive Support), where AI executes specific clinician-commanded tasks.

A critical transition zone, often referred to as the "Behavioral Chasm," exists as systems aim to move from Level 2 to Level 3 (Conditional Automation/-Contextual Proactivity). At Level 3, AI systems begin to perform proactive tasks and make some decisions within a limited, well-defined clinical context or Operational Design Domain (ODD), such as suggesting differential diagnoses or recommending next steps based on the ongoing clinical situation. This shift demands robust behavioral adaptation capabilities to ensure that proactive interventions are safe, appropriate, and effective. Our work on BehaviorSFT and the BehaviorBench evaluation framework is specifically aimed at addressing the challenges of training and assessing these crucial Level 3 behaviors, which are vital for the development of reliable AI co-pilots and assistants. As illustrated in Figure 4, many contemporary applied systems such as MediQ (Li et al., 2024), AIME (Tu et al., 2024), and Med-Gemini (Saab et al., 2024) are operating at or pushing the boundaries of Level 3 capabilities.

The higher autonomy levels, L4 (High Automation/Proactive Decision Support) and L5 (Full Automation/Autonomous Operation), represent the current research frontier for AI in healthcare. Systems like AI Co-Scientist (Gottweis et al., 2025) and AI Scientist v2 (Yamada et al., 2025), while focused on scientific discovery, demonstrate capabilities that conceptually align with L4 by making significant decisions and taking proactive actions within their research ODDs with minimal human oversight for extended periods. Achieving this level of robust autonomy in dynamic, direct clinical care across broad domains remains a significant longterm aspiration for the field.

1097

Enabling frameworks such as AutoGen (Wu et al., 2023) and general concepts like the Proactive Agent (Lu et al., 2024) are instrumental in this progression. They provide the tools and paradigms to build more sophisticated and autonomous AI agents capable of navigating higher levels of task complexity and autonomy. The continued development in this field underscores the critical importance of ensuring that as AI systems become more autonomous, their behaviors are rigorously evaluated and remain aligned, safe, and beneficial within the complex and high-stakes domain of healthcare.

BehaviorSFT Prompt

You are a helpful medical assistant. **Medical Information:**

The patient's history of present illness includes treatment with salve, Alpine lamp, intravenous and intramuscular injections, and Fowler's solution.

Question:

Based on the information in the case summary, how did the patient's treatment for his skin condition evolve from the initial presentation of 'eczema' to the administration of Fowler's solution (arsenic)?

Options:

....

A: "Initially treated with topical steroids...

B: "Initially treated with herbal ...

Instruction:

According to the previous information, give me the behavior first (highly_reactive, primarily_reactive, balanced, highly_proactive, primarily_proactive), then the Rationale and answer in <answer></answer>, later is the detailed option.

F Baseline Performance

1099	Tables 5, 6, and 7 compare o1, Gemini-2.5 Pro,
1100	and DeepSeek-R1 under three prompting regimes—
1101	Zero-Shot (ZS), Few-Shot with three examples
1102	(FS), and ZS augmented by explicit reactive/proac-
1103	tive instructions. All models score near-ceiling on
1104	the Reactive and Balanced subsets, but diverge on
1105	the harder Proactive tasks, where DeepSeek-R1 at-
1106	tains the highest average accuracy (95%), edging
1107	out Gemini and o1 (both \approx 93%). Across mod-
1108	els, FS generally yields the most consistent gains;
1109	especially on items such as predictive next action,
1110	while explicit instructions benefit DeepSeek yet
1111	can slightly reduce performance for Gemini and
1112	o1. These results underscore that, although lower-
1113	level clinical reasoning is largely saturated, proac-
1114	tive reasoning remains the principal differentiator
1115	among state-of-the-art LLMs.



Figure 4: **The Landscape of Healthcare AI Systems and Enabling Frameworks.** Systems are positioned based on their primary Task Scope (Narrow, Medium, or Broad) and their demonstrated level of System Autonomy. The autonomy levels are derived from the Six-Level Taxonomy for Healthcare AI Agent Autonomy (detailed in Table 8), ranging from L0-L1 (Assistance & Reactive Info) through L3 (Conditional Automation/Contextual Proactivity) to L4-L5 (High/Full Automation). Current systems demonstrating L4-L5 capabilities are typically within research frontiers for tasks like scientific discovery rather than direct, broad clinical deployment. Model placement reflects their predominant operational capabilities as described in recent literature (2023-2025). The progression towards higher autonomy, particularly the transition from L2 (Reactive Support) to L3 (Contextual Proactivity), necessitates significant advancements in behavioral adaptation to ensure safe and effective operation in nuanced healthcare contexts. Enabling frameworks and general proactive concepts are also shown, indicating their potential to facilitate the development of more autonomous systems.

H Implementation Details

Our BehaviorSFT has been trained with one 1119 epoch using the adamw_torch optimizer ($\beta_1=0.9$, 1120 $\beta_2 = 0.95, \ \epsilon = 10^{-8}$). The peak learning rate is 1121 1×10^{-4} , decayed with a cosine schedule after a 5 % 1122 warm-up. Training runs in bfloat16 on 4×H200 1123 GPUs with an effective batch size of 64 (per-GPU 1124 batch 4, gradient accumulation 4); weight decay is 1125 0.01 and gradients are clipped to a max-norm of 1126 1127 1.0. For BEHAVIORSFT we add the special tokens <reactive> and <proactive> and attach LoRA 1128 adapters (rank 8, $\alpha = 32$) to all linear layers. The 1129 best checkpoint, selected by validation accuracy 1130 every 100 steps, is reported. 1131

I Clinician-in-the-Loop Evaluation Study

To rigorously evaluate our BehaviorSFT agent and1133validate the proposed dataset, we conducted a comprehensive user study involving board-certified1134prehensive user study involving board-certified1135medical professionals. This study was designed1136to assess the clinical utility of BEHAVIORBENCH1137and to compare the performance of LLM agents1138exhibiting distinct behavioral characteristics.1139

1132

1140

I.1 Participant Recruitment and Profile

We recruited three medical doctors and each physical1141cian underwent a standardized orientation session1142to familiarize them with the study objectives, anno-1143tation tasks, and the custom-developed user inter-1144faces.1145



Figure 5: **Performance comparison on BEHAVIORBENCH for Few-Shot (k=3); Gen. SFT, and our proposed BehaviorSFT.** Tasks are colored based on task category: Reactive, Balanced, and Proactive. The radar plot illustrates that our BehaviorSFT achieves best or second-best performance across all task categories. While all methods perform strongly on Reactive and Balanced tasks, the gains from BehaviorSFT are most pronounced in complex Proactive scenarios, highlighting its effectiveness in enhancing nuanced behavioral capabilities of agents beyond standard fine-tuning approaches.

Table 5: **Performance Evaluation on BEHAVIORBENCH.** Accuracy (%) across task categories. Best result per task in **bold**. Baseline LLM is o1. 'ZS' = Zero-Shot, 'FS (k=3)' = Few-Shot (3 examples), 'Explicit Instr.' = ZS with explicit reactive/proactive instruction.

Category	Task		Base	line
		ZS	FS (k=3)	ZS + Explicit Instr.
e	fact_retrieval	100.00	100.00	100.00
ţi	timeline_sequence	100.00	100.00	100.00
eac	ddx_reasoning	93.92	91.96	91.92
Ř	treatment_decision	91.88	93.78	91.88
	Average	96.45	96.43	95.95
ed	reasoning_diff_evolution	98.05	100.00	100.00
Balance	integrity_missing_turn	100.00	98.46	100.00
	Average	99.03	99.23	100.00
	consistency_check	95.23	95.24	90.12
	data_conflict_resolution	96.52	96.44	95.11
ive	interpretation_conflict	98.48	98.30	98.29
act	standard_of_care	91.47	91.79	94.87
<u>e</u>	omission_detection	81.87	82.00	81.61
4	explicit_error_correction	96.30	98.12	95.54
	predictive_next_action	78.03	82.88	78.30
	Average	93.31	92.11	90.55
Average		93.86	94.25	93.55

I.2 Study Design and Procedure 1146

The study was structured into two principal phases,1147each targeting specific evaluation objectives:1148

Table 6: **Performance Evaluation on BEHAVIORBENCH.** We report Accuracy (%) across different task categories. Best result per task is highlighted in **bold**. Baseline LLM used is Gemini-2.5 Pro. 'ZS' = Zero-Shot, 'FS (k=3)' = Few-Shot (3 examples), 'CoT' = Chain-of-Thought, 'Explicit Instr.' = ZS with explicit reactive/proactive instruction.

Category	Task		Base	line
		ZS	FS (k=3)	ZS + Explicit Instr.
e	fact_retrieval	100.00	100.00	100.00
ti	timeline_sequence	99.10	78.65	99.10
eac	ddx_reasoning	95.33	93.99	94.56
R	treatment_decision	94.77	93.88	94.29
	Average	97.30	91.63	96.99
pə	reasoning_diff_evolution	98.59	82.33	97.26
anc	integrity_missing_turn	98.46	98.05	96.56
Bal	Average	98.53	90.19	96.91
	consistency_check	94.29	96.34	94.29
	data_conflict_resolution	97.18	97.24	98.53
ive	interpretation_conflict	96.70	95.11	94.95
act	standard_of_care	95.32	96.80	92.11
<u>p</u>	omission_detection	81.57	90.10	79.12
L	explicit_error_correction	96.34	94.23	95.55
	predictive_next_action	77.88	81.55	73.25
	Average	91.33	93.05	89.69
Average		94.27	92.17	93.04

Table 7: **Performance Evaluation on BEHAVIORBENCH.** We report Accuracy (%) across different task categories. Best result per task is highlighted in **bold**. Baseline LLM used is DeepSeek-R1. 'ZS' = Zero-Shot, 'FS (k=3)' = Few-Shot (3 examples), 'CoT' = Chain-of-Thought, 'Explicit Instr.' = ZS with explicit reactive/proactive instruction.

Category	Task		Base	line
		ZS	FS (k=3)	ZS + Explicit Instr.
ě	fact_retrieval	100.00	100.00	100.00
ţi,	timeline_sequence	100.00	100.00	100.00
Sac	ddx_reasoning	93.16	91.16	94.25
R	treatment_decision	94.22	95.70	94.77
	Average	96.84	96.71	97.26
pa	reasoning_differential_evolution	98.59	98.59	98.59
nc	integrity_missing_turn_inference	100.00	100.00	100.00
Bal	Average	99.29	99.29	99.29
	consistency_check	94.29	94.29	100.00
	data_conflict_resolution	97.18	95.68	97.88
ive	interpretation_conflict	100.00	96.53	98.22
act	standard_of_care	93.52	95.32	94.67
Ĵ,	omission_detection	93.78	90.75	93.57
4	explicit_error_correction	97.50	97.52	98.26
	predictive_next_action	78.54	80.86	82.69
	Average	93.54	92.99	95.04
Average		95.49	94.96	96.10

Phase 1: Dataset Validation

1149

1150

1151

1152

1153

1154

1155

In this phase, clinicians were tasked with validating a randomly selected subset of tasks (N=30) from the BEHAVIORBENCH. The primary goal was to ascertain the clinical soundness and appropriateness of the dataset components. For each presented task, which included a clinical 'Task Context', a specific 'Question', and multiple-choice 'Options'1156(as illustrated in Figure 14), clinicians utilized a1157dedicated evaluation panel (Figure 13). Their eval-1158uation encompassed:1159

• Correctness of Ground Truth: Verifying 1160 the accuracy of the designated correct answer 1161



Figure 6: **Distribution of total dialogue length (in characters) per conversation.** This metric captures the overall verbosity of clinical discussions. Most conversations range between 3000 and 5000 characters in length, indicating substantial detail per case.



Figure 7: **Distribution of the number of dialogue turns per conversation.** Each conversation represents a real-world clinical case discussion, with turns corresponding to speaker exchanges. The majority of cases fall between 15 and 30 turns.



Figure 8: **Distribution of instances across specific task types in BEHAVIORBENCH.** Each bar represents the frequency of a task type, colored by its average behavior score (blue = reactive, red = proactive). This illustrates the diversity of evaluation scenarios, spanning a wide range of communicative functions and behavioral expectations.



Figure 9: **Distribution of instances by task complexity level in BEHAVIORBENCH.** Tasks are broadly categorized as either 'intermediate' or 'advanced' based on reasoning depth and contextual demands. The dataset skews toward advanced tasks, aligning with the goal of evaluating high-autonomy agent behavior.



Figure 10: Average proactive score by task complexity level in BEHAVIORBENCH. Tasks labeled as 'advanced' exhibit a significantly higher average proactive score (above 0.8) compared to 'intermediate' tasks (around 0.4), highlighting the alignment between task complexity and expected behavioral autonomy in clinical reasoning.

1164

1165

1166

1167 1168

1169

1170

1171

1172

among the provided options.

- Annotator Confidence: Rating their confidence in their selected answer on a three-point scale (Low, Moderate, High).
- Task Proactivity Level Assessment: Evaluating the inherent proactivity level of the question itself on a continuous scale ranging from 0.0 (Reactive) to 1.0 (Proactive). This aimed to capture the degree to which the question prompted an anticipatory or forward-looking response.
- Clinical Plausibility: Determining if the task (question and options combined) was clinically plausible and relevant within the given case context, with options "Yes," "No," or "Unsure."

To ensure comprehensive understanding, clinicians1178had access to the broader 'Case Context', including1179a 'Case Presentation Summary', the 'Full Conversation' transcript leading to the task, and an option1181to refer to the original medical case for in-depth1182review (Figure 15).1183



Figure 11: **Distribution of continuous behavior scores across all tasks in BEHAVIORBENCH.** The behavior score ranges from 0.0 (fully reactive) to 1.0 (fully proactive), with the distribution skewed toward higher scores, indicating a dataset emphasis on proactive clinical reasoning.



Figure 12: **Distribution of tasks across discrete behavior categories in BEHAVIORBENCH.** Tasks are grouped into five categories, ranging from 'highly reactive' to 'highly proactive' to support structured evaluation of agent behavior along the autonomy spectrum.

Phase 2: Comparative Agent Behavior Evaluation

1184

1185

This phase focused on evaluating the quality and 1186 safety of responses generated by three distinct LLM 1187 agent archetypes when presented with N=10 clin-1188 ical tasks from BEHAVIORBENCH. The agents 1189 included: (1) BehaviorSFT: An agent fine-tuned 1190 1191 using our proposed BEHAVIORBENCH approach. (2) General SFT: An agent subjected to general 1192 supervised fine-tuning without specific behavioral 1193 guidance. (3) ZS + Explicit Instr.: An agent op-1194 erating in a zero-shot setting, guided by explicit 1195

instructions on desired behavior.

For each scenario, clinicians were first presented with the 'Question Posed to AI' and the 'Task Options' (with the correct answer highlighted for their reference). Subsequently, the responses from the three LLM agents were displayed side-by-side (Figure 17). The identity and order of these agents (Agent A, B, C) were anonymized and randomized for each task to mitigate bias. Using the feedback panel shown in Figure 16, clinicians performed the following evaluations:

• Comparative Ranking: Ranking the three 1207

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1285

1286

1257

1208agent responses from best (1st) to worst (3rd)1209using a drag-and-drop mechanism.

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238 1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251 1252

1253

1254

1255

1256

- **Safety Assessment:** Identifying and describing any instances of clinically unsafe information, critical errors, or significant omissions in any of the agent responses.
- **Proactivity/Reactivity** Appropriateness: Rating the appropriateness of each agent's proactivity or reactivity level on a 5-point Likert scale (1: Very Inappropriate, 3: Neutral, 5: Very Appropriate).

I.3 Interface Design for Annotation Tasks

Custom-designed web-based interfaces were developed to ensure a standardized, intuitive, and efficient annotation experience for the participating clinicians. The interfaces were tailored to the specific requirements of each study phase (see Figure 13, 14, 15, 16 and 17).

I.4 Annotation Results

This section presents the quantitative and qualitative findings from the clinician-in-the-loop evaluation study. All reported inter-annotator agreement scores were calculated among the three participating physicians.

I.4.1 Phase 1: BEHAVIORBENCH Validation

Clinicians evaluated a total of 60 unique tasks from the BEHAVIORBENCH.

MCQ Accuracy and Task Plausibility The physician annotators demonstrated a high level of accuracy in answering the multiple-choice questions, achieving an overall correctness of 83.3%. This proficiency underscores their expert understanding of the clinical scenarios presented within the dataset.

The clinical plausibility of the tasks was a key validation metric. As shown in Figure **??**, a substantial majority of tasks (**80.0**%) were rated as clinically plausible ("Yes"). No tasks (0.0%) were rated as definitively "No" for plausibility, while 20.0% were marked as "Unsure," suggesting areas where task framing or context might warrant further refinement or clarification for some annotators.

Annotator Confidence Levels Annotator confidence in their selected MCQ answers was recorded on a three-point scale. The distribution, illustrated in Figure ??, reveals that physicians were predominantly "High" in their confidence (55.0%). "Moderate" confidence was reported for 36.67% of answers, while "Low" confidence was expressed for

only 8.33% of answers. This general trend towards higher confidence aligns with the observed accuracy.

Inter-Annotator Agreement for Dataset Validation To ensure the reliability of the dataset validation process, inter-annotator agreement was quantified using the Intraclass Correlation Coefficient (ICC3) for continuous ratings.

The task proactivity/reactivity slider ratings (0.0-1.0 scale) demonstrated *good* reliability with an ICC3 of **0.61**. This robust agreement scores indicate that the physicians interpreted and applied the validation criteria consistently.

I.4.2 Phase 2: Comparative Agent Behavior Evaluation Results

Physicians evaluated agent responses across N=24 unique clinical tasks. The anonymized agents evaluated were BehaviorSFT, General SFT, and ZS + Explicit Instr.

Agent Response Ranking and Proactivity/Reactivity Appropriateness The primary evaluation involved ranking the three agents. Agent A (BE-HAVIORBENCH) received the most favorable rankings, achieving the lowest (best) mean rank of 1.80 (Figure 18). In terms of the appropriateness of proactivity/reactivity, Agent C (ZS + Explicit Instr.) scored highest with a mean Likert score of 4.20 out of 5 (Figure 19). Agent B (General SFT) had a mean rank of 2.08 and a mean Likert score of 4.08.

Annotator's Evaluation:
1. Your Confidence in Your Selected Answer:
○ High ○ Moderate ○ Low
2. Your Assessment of Task Proactivity Level:
0.0 (Reactive) 1.0 (Proactive)0.50
3. Is this task (question and options) clinically plausible within the case context?
\bigcirc Yes \bigcirc No \bigcirc Unsure
Previous Task Save & Next Task
Export Annotations
Annotations are automatically saved to Firebase as you click "Save & Next Task".
Annotated tasks for this case: 1 / 3

Figure 13: Interface for **Dataset Task Validation: Annotator's Evaluation**. Medical doctors used this panel to provide their confidence in the selected answer for a given task, assess the task's inherent proactivity level on a continuous scale (0.0 Reactive to 1.0 Proactive), and confirm the clinical plausibility of the task (question and options) within the provided case context.

Show/Hide Full Conversation

Task 1 of 3

Task Context:

Key Physical Exam Findings (breathing description), Vitals from Turn 12 (Respiratory rate 45), Patient's condition deterioration notes, Eventual outcome (expired).

Question:

The patient's respiratory rate is documented as 45. Considering this finding alongside the description of 'rapid, shallow, labored breathing' and the eventual outcome, does the documentation adequately describe the severity of respiratory compromise needed to guide modern supportive care interventions (like oxygenation goals or ventilation decisions)? What key physiological parameter, routinely monitored today, is potentially overlooked?

Options:

- A: No, the documentation, while noting the high respiratory rate and character of breathing, is inadequate by modern standards for fully assessing the severity of respiratory compromise and guiding interventions. A crucial missing parameter is oxygen saturation (SpO2), which would provide an objective measure of gas exchange efficiency.
 Blood gas analysis (pO2, pCO2, pH) would also be essential for assessing hypoxemia, hypercapnia, and respiratory acidosis/alkalosis, which are critical for managing severe respiratory failure.
- B: No, the documentation is insufficient, and the most critical missing detail is a subjective rating of dyspnea severity by the patient, which would better capture their experience of respiratory compromise.
- C: Yes, the documented respiratory rate of 45 and description of labored breathing provide sufficient initial assessment to guide critical interventions like high-flow oxygen therapy, making immediate objective measures less urgent.
- D: Yes, the description 'rapid, shallow, labored breathing' strongly implies significant ventilatory impairment, and while SpO2 is useful, a more direct assessment of work of breathing through documenting accessory muscle use or retractions would be the *key* overlooked parameter.
- E: No, while inadequate, the most crucial overlooked parameter is likely capillary refill time, which offers a rapid assessment of overall tissue perfusion and the severity of shock contributing to the respiratory distress.
- F: Yes, the respiratory rate of 45 alone in the context of labored breathing is highly indicative of severe respiratory failure, adequately justifying immediate intubation and mechanical ventilation without requiring further specific gas exchange data for initial management decisions.

Figure 14: Interface for **Dataset Task Validation: Task Presentation**. This view provided clinicians with the 'Task Context' (relevant excerpts from the case), the specific 'Question' being posed for the BehaviorBench task, and the multiple-choice 'Options', one of which was the ground truth answer they were validating.

ase Prese	ntation Summary (from Task File)
Case ID: NEJM	1192310251891705. (No summary).
	sation (from Conversation File)
🤒 Patient (Tu	ım 1):
Doctor, l've be	en having pain and trouble breathing. It started a few days ago like a bad cold.
F Doctor A (1	Furn 2):
Tell me more	about when this started and how it has progressed.
🤒 Patient (Tu	rn 3):
It was last Thi	ursday afternoon, March 23rd. I got a bad cold, headache, and backache. Felt really weak. Stayed in bec
most of Friday	y and all of Saturday. On Saturday, March 25th, I took aspirin for the headache and backache, which That affermoon L started coughing dry cough at first. Felt feverish and yery weak. Then around four
yesterday afte	ernoon, I had a really bad chill, shaking all over.
Doctor A (1	Furn 4):
And the pain a	and breathing difficulty, when did that start?
🤒 Patient (Tu	rn 5):
Download Or	iginal Case PDF

Figure 15: Interface for **Dataset Task Validation: Case Context Provision**. To ensure comprehensive understanding, clinicians had access to the broader 'Case Context', including a 'Case Presentation Summary' (if available from the task file), the 'Full Conversation' transcript leading up to the point of the task, and an option to download the original case PDF for in-depth review.

Evaluation Questions for Task 1 (Case 10111 - Type: reasoning_differential_evolution)							
Annotator ID:							
Enter your uniqu	e ID						
Rank the three /	Agent responses from Best	to Worst by dragging	them:				
Drag to order: Bes	t > Moderate > Worst						
	≡ Agent A		≡ Agent B	≡ Agent C			
Did any Agent ir	nteraction present informati	ion you consider clini	cally unsafe or a critical error/omission?				
Did any Agent ir Agent A	nteraction present informati Agent B Agent C se specify which Agent (ion you consider clini (e.g., Agent A) and	cally unsafe or a critical error/omission? d the issue				
Did any Agent ir	Neraction present informati Agent B Agent C se specify which Agent (vriateness of proactivity/rea	ion you consider clini (e.g., Agent A) and ctivity for each displa	cally unsafe or a critical error/omission? d the issue ayed Agent:				
Did any Agent ir Agent A If yes, plea: Rate the approp Agent A:	Agent B Agent C agent B Agent C as specify which Agent I riateness of proactivity/rea 0 1 (Very Inapp.)	(e.g., Agent A) and ctivity for each displa 2 3 (Neutral)	cally unsafe or a critical error/omission? d the issue ayed Agent: ○ 4 ○ 5 (Very Approp.)				
Did any Agent ir Agent A If yes, pleas Rate the approp Agent A: Agent B:	Agent B Agent C Agent (Agent Agent Agent) Agent (Agent Agent) A (Very Inapp.) 1 (Very Inapp.) 1 (Very Inapp.)	(e.g., Agent A) and ctivity for each displa 2 3 (Neutral) 2 3 (Neutral)	cally unsafe or a critical error/omission? d the issue ayed Agent: 4 5 (Very Approp.) 4 5 (Very Approp.)				
Did any Agent ir Agent A If yes, pleas Rate the approp Agent A: Agent B: Agent C:	teraction present informati Agent B Agent C se specify which Agent (arriateness of proactivity/rea 1 (Very Inapp.) 1 (Very Inapp.) 1 (Very Inapp.)	(e.g., Agent A) and ctivity for each displa 2 3 (Neutral) 2 3 (Neutral) 2 3 (Neutral) 2 3 (Neutral)	cally unsafe or a critical error/omission? d the issue ayed Agent: 4 5 (Very Approp.) 4 5 (Very Approp.) 4 5 (Very Approp.) 4 5 (Very Approp.) 4 5 (Very Approp.)				

Figure 16: Interface for **Agent Behavior Evaluation: Clinician Feedback Panel**. After reviewing the task and agent responses (shown in Figure **??**), medical doctors used this panel to: (1) Rank the three anonymized agent responses (Agent A, B, C) from best to worst via drag-and-drop. (2) Identify and describe any clinically unsafe information or critical errors/omissions presented by any agent. (3) Rate the appropriateness of the proactivity/reactivity level for each agent's response on a 5-point Likert scale (from Very Inappropriate to Very Appropriate).



Figure 17: Interface for **Agent Behavior Evaluation: Task and Agent Response Display**. For each evaluation scenario, clinicians were presented with the 'Question Posed to AI' and the 'Task Options' (with the correct answer highlighted for reference). Below this, the distinct responses from three anonymized LLM agents (Agent A, B, C), including their rationales, were displayed side-by-side for comparative assessment.



Figure 18: (a) Over half (55.0%) of the responses were marked as 'High' confidence, while 'Moderate' confidence accounted for 36.7%. 'Low' confidence was the least frequent category, representing only 8.3% of responses. (b) The vast majority (80.0%) of responses affirmed the clinical plausibility ('Yes') of the generated MCQs. A smaller portion (20.0%) of responses were 'Unsure', and no responses found the MCQs implausible ('No').



Figure 19: (a) Mean appropriateness scores for agent proactivity/reactivity (5-point Likert scale, higher is better). (b) BehaviorSFT received the lowest (best) mean rank (1.80), suggesting it was most frequently ranked highest by evaluators. Gen. SFT had a mean rank of 2.08, while ZS w/ explicit instruction had the highest (worst) mean rank of 2.12 in a system where lower ranks are better.

Table 8:	Six-Level	Taxonomy fo	or Healthcare	AI Agent	Autonomy
				0	

Level	Name	AI Agent's Role / Capability	Human Clinician's Role
0	No Automation	The AI system provides no assistance or automation for any clinical task.	Performs all tasks and makes all decisions related to patient care. The AI system is not involved.
1	Clinician Assistance	The AI system may provide information, simple alerts based on predefined rules (e.g., drug interaction warnings, out-of-range lab value notifications), or basic data visualization. It does not perform any part of the dynamic clinical task itself.	Performs all dynamic decision-making and actions. Uses the AI as a passive information source or a simple alerting tool. Responsible for interpreting AI-provided information.
2	Partial Automation (Re- active Support)	The AI system can execute specific, well-defined reactive sub-tasks under direct human supervision based on explicit clinician queries or predefined triggers (e.g., retrieving specific patient history, summarizing recent lab results, performing image segmentation on request). It does not manage the overall clinical situation.	Actively monitors the AI's execution of sub-tasks, provides necessary inputs, and must intervene if the AI's output is incorrect or inappropriate. Responsible for the overall task and integrating AI's contribution.
3	Conditional Automation (Contextual Proactivity)	The AI system can perform certain proactive tasks and make some decisions within a limited, well-defined clinical context or Operational Design Domain (ODD) (e.g., suggesting differential diagnoses based on current symptoms, flagging potential omissions in a standard care plan, recommending next tests). It can handle some dynamic aspects of the task.	Monitors the AI and the clinical environment. Must be ready to take over control if the AI encounters a situation it cannot handle, if its suggestions are inappropriate, or if the situation goes outside the AI's ODD.
4	High Automation (Proactive Decision Support)	The AI system can make significant clinical decisions and take proactive actions in most situations within its designed ODD without human oversight for extended periods (e.g., autonomously adjusting medication dosage based on real-time patient data within set parameters, initiating standard protocols for common conditions, triaging patients based on urgency).	Primarily acts as a fallback, intervening only in complex, novel, or out-of-ODD scenarios. Relies on the AI for most routine decisions and actions within the ODD. May oversee multiple AI-managed cases.
5	Full Automation (Au- tonomous Operation)	The AI system can perform all clinical tasks and make all decisions that a human healthcare professional can, under all conditions within its defined scope of operation. It can adapt to novel situations and operate entirely autonomously, potentially even taking on roles currently performed by specialized clinicians.	May not be required for tasks within the AI's full operational scope. Human role shifts to high-level oversight, system management, or handling tasks entirely beyond the AI's designed capabilities or ethical boundaries.

ODD: Operational Design Domain - The specific conditions under which a given AI system or feature is designed to function.