
Polaris: A Safety-focused LLM Constellation Architecture for Healthcare

Subhabrata Mukherjee *	Paul Gamble	Markel Sanz Ausin	
Neel Kant	Kriti Aggarwal	Neha Manjunath	Debajyoti Datta
Zhengliang Liu	Jiayuan Ding	Sophia Busacca	Cezanne Bianco
Swapnil Sharma	Rae Lasko	Michelle Voisard	Sanchay Harneja
Darya Filippova	Gerry Meixiong	Kevin Cha	Amir Youssefi
Meyhaa Buvanesh	Howard Weingram	Sebastian Bierman-Lytle	
Harpreet Singh Mangat	Kim Parikh	Saad Godil	Alex Miller

Hippocratic AI

Abstract

We develop Polaris ² the first safety-focused Large Language Model (LLM) constellation for real-time patient-AI healthcare conversations. Unlike prior LLM works in healthcare, which focus on tasks like question answering, our work specifically focuses on long multi-turn voice conversations. Our one-trillion parameter constellation system is composed of several multi-billion parameter LLMs as co-operative agents: a stateful primary agent that focuses on driving an engaging patient-friendly conversation and several specialist support agents focused on clinical tasks performed by nurses, social workers, and nutritionists to increase safety and reduce hallucinations. We develop a sophisticated training protocol for iterative co-training of the agents that optimize for diverse objectives.

We train our models on proprietary data, clinical care plans, healthcare regulatory documents, medical manuals, and other medical reasoning documents. We further align our models to speak like medical professionals, using organic healthcare conversations and simulated ones between patient actors and experienced care-management nurses. This allows our system to express unique capabilities such as rapport building, trust building, empathy and bedside manner augmented with advanced medical reasoning.

Finally, we present the first comprehensive clinician evaluation of an LLM system for healthcare. We recruited over 1100 U.S. Licensed nurses and over 130 U.S. Licensed physicians to perform end-to-end conversational evaluations of our system by posing as patients and rating the system on several fine-grained measures. We demonstrate Polaris performs on par with human nurses on aggregate across dimensions such as medical safety, clinical readiness, patient education, conversational quality, and bedside manner. Additionally, we conduct a challenging task-based evaluation of the individual specialist support agents, where we demonstrate our LLM agents significantly outperform a much larger general-purpose LLM (GPT-4) as well as one from its own medium-size class (LLaMA-2 70B).

*Correspondence to research-team@hippocraticai.com.

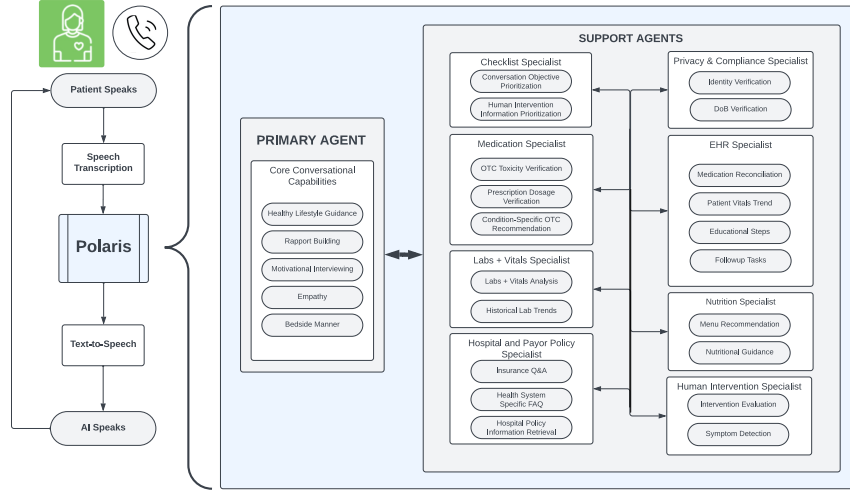
²Safety is our North Star. We name our system after Polaris, a star in the northern circumpolar constellation of Ursa Minor, currently designated as the North Star.

Contents

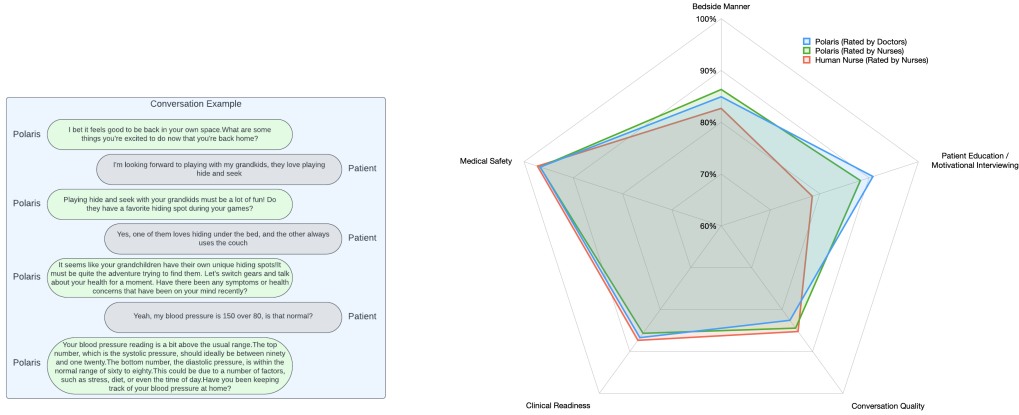
1	Introduction	4
2	System Overview	6
2.1	Objectives and Use Cases	6
2.2	Constellation Architecture	7
2.3	Orchestration	9
2.4	LLM Details	9
3	Conversational Alignment	10
3.1	Data	10
3.1.1	Foundation Model Training Data	10
3.1.2	Simulated Conversations	11
3.1.3	Conversational Dataset Construction	13
3.2	General Instruction Tuning	14
3.3	Conversation and Agent Tuning	14
3.4	RLHF	16
4	Specialist Support Agents	16
4.1	Privacy & Compliance Specialist	17
4.2	Checklist Specialist	17
4.3	Medication Specialist	18
4.3.1	Prescription Adherence	19
4.3.2	Condition Specific Disallowed OTCs	20
4.3.3	Over-The-Counter (OTC) Toxicity	20
4.3.4	Unrecognized Medications Reconciliation	20
4.3.5	Models	21
4.3.6	Training Data Generation	22
4.4	Labs & Vitals Specialist	22
4.4.1	Determination of Labs & Vitals and Value Extraction	22
4.4.2	Normal Range Assessment & Plausibility Check	22
4.4.3	Historical Lab Trends	24
4.4.4	Conditional Factors	24
4.4.5	Summary of Labs & Vitals Specialist Capabilities	25
4.4.6	Models	25
4.4.7	Labs & Vitals Specialist Training Datasets	26
4.5	Nutrition Specialist	27
4.5.1	Workflow	27
4.5.2	Chain Restaurant Nutrient Dataset Curation	28
4.5.3	Models	28

4.5.4	Training Data Generation	28
4.6	Hospital & Payor Policy Specialist	28
4.6.1	Hospital Policy Q&A within the Hospital System	30
4.6.2	Models	31
4.6.3	Training Data Generation	31
4.7	Summary Agent	31
4.8	Human Intervention Specialist	33
4.8.1	Models	34
4.8.2	Training Datasets	34
5	Evaluation	35
5.1	Overall Subjective Evaluation	36
5.2	Clinical Capability Evaluation	36
5.2.1	Medications and Supplements	38
5.2.2	Lab Values and Vital Signs	38
5.2.3	Dietary Recommendations	39
5.2.4	Answering Hospital Policy Questions	39
5.2.5	Privacy	39
5.3	Error Analysis	40
5.3.1	Medication Specialist	40
5.3.2	Lab Values and Vital Signs	41
5.3.3	Nutrition Specialist	43
5.3.4	Hospital and Payor Policy Specialist	43
6	Related Work	44
6.1	Current Healthcare Challenges	44
6.2	LLMs in Healthcare	45
6.3	Labor Market Impact of Healthcare	46
7	Future Work	46
8	Conclusions	47
A	Appendix	51
A.1	Additional Conversation Examples	51

Figure 1: High level overview of our LLM constellation architecture Polaris .



(a) Overview of our architecture, comprising of the Automatic Speech Recognition (ASR) for speech transcription, Polaris for processing the textual utterances, and Text-To-Speech (TTS) for the audio output. The constellation within Polaris contains a primary LLM agent driving the conversation, and several specialist LLM agents providing task-specific context to it.



(b) Conversation snippet between Polaris and a simulated Patient, emphasizing empathy and rapport building as part of good bedside manner, while providing accurate medical information using lab specialist agent's assistance.

(c) Comparative evaluation between human nurses (U.S. licensed) and Polaris on bedside manner (e.g., empathy, trust, rapport), medical safety, patient education, clinical readiness and overall conversation quality. Overall, Polaris is strikingly close to human nurse performance, and even outperforms them on some key dimensions.

1 Introduction

Recent progress in large language models (LLMs) has shown their impressive capability to plan, reason and interact with humans for a variety of tasks such as web search [1], coding [2], and intelligent content creation [3]. The scaling of LLMs and the datasets used to train them, together with new architectural advances has contributed greatly to the advances in AI capabilities, and shown to surpass human performance in various benchmark tasks [4]. These new capabilities are enabling real-world applications that were impossible until very recently such as those in healthcare.

Healthcare specialization. Amidst the several domains and applications, healthcare remains a high-stake domain where errors may have fatal implications. The advent of GPT-4 [5] has led to a profound surge in the use of AI for healthcare applications, including

clinical note and electronic health record processing (see [6] for overview). While systems like MedPALM [7] and GPT-4 have shown impressive results in general medical benchmarks like USMLE, recent work shows significant error rate for more specialized use-cases like pediatrics [8]. Among the failures, researchers observe that the AI systems struggle to spot known relationships between conditions that an experienced physician would look for, e.g., for a patient with autism, a physician might check for dietary deficiencies. [9]. The researchers note that these systems could be improved by selectively training on accurate and high quality medical literature, not just general articles over the internet, which are typically what most LLMs are trained on. Furthermore, it is possible that the base knowledge about autism leading to nutrient deficiencies is present in the LLMs; what is missing is the *medical reasoning* to connect the dots (i.e., some patients with autism exhibit narrow dietary preferences, which can then lead to nutrient deficiencies).

Conversational Healthcare Systems. While most of the existing works in healthcare are focused on tasks like medical question-answering [10] and EHR summarization [11], there are very few works focusing on natural dialogue between caregivers and patients. More recently, Tu et al. developed Articulate Medical Intelligence Explorer (AMIE) [12], which outlines the importance of diagnostic dialogue to enable physicians to make diagnoses and develop management plans. The complexity of healthcare conversations highlights the need for the caregiver to build rapport and trust with the patient to make them feel safe, supported, and confident in their care, while communicating with empathy and bedside manner. Such relationship has shown to lead to better patient satisfaction and, ultimately, better outcomes in real world [13]. General-purpose LLM’s, however, are not optimized for such objectives. Furthermore, these LLM’s are also not optimized for real-time, voice-based conversations, which can be quite different from text-based conversations. For instance, factors such as response length, quality of voice, pauses and interrupts greatly impact the subjective experience.

The case for AI-based Healthcare Agents. The US healthcare industry is facing a massive shortage of healthcare workers, that became even more apparent after the COVID-19 pandemic [14]. Exacerbated by burnout, stress and financial conditions, 16.7% of hospitals anticipated a critical staffing shortage in 2023 according to the Department of Health and Human Service [15]. The U.S. Bureau of Labor Statistics estimates the need to fill over 200,000 nursing positions every year until 2031 [16]. A 2023 survey found that 28.7% percent of nurses were considering to leave their jobs. The trend in the decline of the US workforce indicates a shortage of more than 4 million workers nationwide by 2026 [14]. In the meantime, the demand for healthcare continues to grow as the population continues to age. There are currently 46 million adults over 65, and it will increase to 64 million by 2030, and 90 million by 2050 [17].

Given this massive gap in supply and demand for healthcare workers, and the recent promise of Generative AI to supercharge productivity, we focus on developing a *non-diagnostic* technology for healthcare workforce augmentation, which we denote as *super-staffing*. In this work, we develop autonomous generative AI healthcare agents that can safely converse with patients on medical topics. Our goal is to improve patient healthcare outcomes by providing a scalable and safe system that reduces the burden on humans and healthcare providers suffering from staffing shortage with built-in safety guardrails.

To address all the above challenges, we develop a multi-agent system with highly specialized healthcare LLMs. The following are the primary contributions of our work.

Architecture and Training. We develop a unique multi-agent LLM constellation architecture optimized for real-time healthcare conversation. We employ a primary conversational agent with several specialist support agents for a patient-friendly and medically accurate conversation. The primary agent is trained to be aligned with nurse-like conversations geared for building trust, rapport and empathy with patients, as well as accomplishing healthcare-specific tasks typically performed by nurses, medical assistants, social workers, and nutritionists. We develop techniques to make the primary agent *stateful* to navigate through a long checklist of care protocols. Our support agents are specialized in healthcare tasks that require a high-level of accuracy, such as confirming that the patient’s reported medicine consumption aligns with the dosage prescribed by their physician, determining

whether the patient’s reported OTC drug consumption is within the manufacturer’s recommended range, understanding which medicine the patient is referring to (i.e., patients often struggle with correctly pronouncing drug names, and may confuse ones that sound similar), retrieving current and historical lab results from the patient’s EHR, guiding the patient’s food choices to align with their prescribed diet, etc. The specialist agents provide the relevant context to the primary agent in a message-passing framework, which drives the main conversation. To this end, we develop custom training protocols for conversational alignment using organic healthcare conversations and simulated ones using patient actors and nurses (U.S. licensed) with agents- and clinicians-in-the-loop for co-operative learning.

Safety. We adopt a three-pronged approach to safety comprised of: (1) A 70B-100B primary model trained using evidence based content; (2) a novel constellation architecture with multiple models totaling over one trillion parameters, in which a primary LLM is supervised by multiple specialist support models to improve medical accuracy and substantially reduce hallucinations; (3) built-in guardrails that bring in a human supervisor when necessary.

Evaluation. To the best of our knowledge, we perform the first extensive real-world evaluation of a healthcare LLM system in which we recruited over 1100 US-licensed nurses and over 130 U.S.-licensed physicians posing as patients for our system. This is focused on an integrated system-level conversational evaluation on dimensions such as medical safety, readiness and bedside manners where we demonstrate parity with human nurses on several key metrics. We also perform a challenging component-level evaluation where we demonstrate that our medium-size specialized agents massively outperform a much larger state-of-the-art general-purpose LLM (GPT-4) on the healthcare tasks as well as outperform an LLM from its own parameter size class (LLaMA-2 70B).

2 System Overview

Polaris is architected as a constellation of a primary and multiple specialist LLMs. It also includes the supporting system which orchestrates control flow, inter-model message passing and maintains conversation state. The system is designed to achieve better domain-specific interactions compared to a single general purpose LLM. The healthcare conversation domain is apt for showcasing the value of this paradigm as there are many competing objectives and requirements, including a special emphasis on safety and verification. We further explain the problem domain, system architecture and also provide relevant details on the constituent models in the following sections.

2.1 Objectives and Use Cases

We develop Polaris for patient-facing, real-time healthcare conversation. The objective of this system is to perform low risk non-diagnostic tasks typically performed by nurses, medical assistants, social workers, and nutritionists. The system is constructed for voice-based interaction, as phone calls are the dominant method of communication for many high-volume and high-value healthcare use cases. Building a voice-based autonomous agent is a challenging problem due to factors such as voice quality, pitch, tone, response length, interruptions, and communication delay. Additionally, Polaris has to take into account errors introduced by the Automated Speech Recognition (ASR) system in the transcription of the audio signal, with a particular focus on the complexity of recognizing healthcare-specific words and phrases (e.g., complex medication names and lab values).

We train our healthcare agents to accomplish a number of tasks that can be completed via phone calls or voice-only interactions. Typical objectives for these tasks include general check-in on patient wellness, reviewing compliance with prescribed medicine regimes, confirming appointment details, reviewing procedural logistics, performing diet reviews, communicating lab results, etc. Given that the number of tasks in a typical care protocol can be quite large and extensive, the agent must be able to *maintain and update conversation state* with each patient, ensuring that all tasks are completed. The task completion rate directly impacts the success rate for the call. At the same time, we train these agents to have natural conversations that mimic what human caregivers would say, rather than a mere question-answering system as in prior works. For instance, the agent must be able to answer

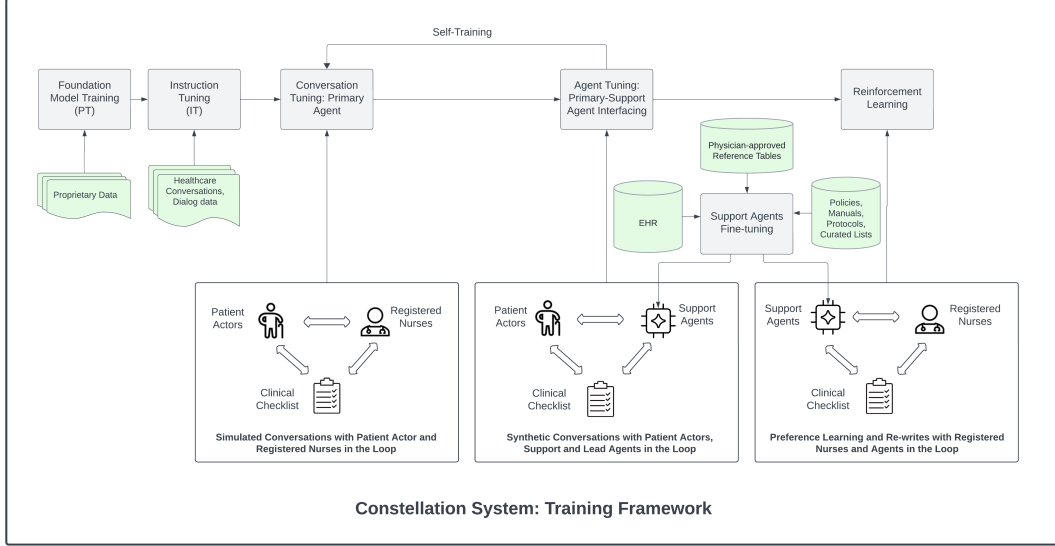


Figure 2: Overview of our training framework for Polaris with resident nurses, patient actors, and LLM agents in the loop.

any question the patients may have, address their concerns, and otherwise handle tangential discussions in the conversation. This makes it particularly challenging for the agent to engage in deep conversations with the patient, given that the style, tone and content of these conversations can be quite different from the model’s original training dataset. Such conversation fluidity coupled with its state awareness and medical accuracy remains at the crux of all our design and architectural decisions. Furthermore, we design our system to build trust and rapport with patients, learning about their personality and preferences to demonstrate good bedside manner and leading to higher patient satisfaction and potentially better health outcomes. Our system is deeply specialized in many different medical conditions and procedures. For the purpose of this report, we highlight three representative outbound-call use cases including follow-up calls for post-discharge congestive heart failure (CHF), ongoing care for chronic kidney disease (CKD), and pre-operative check-in for colonoscopy. Figure 2 shows an overview of the different capabilities of our system.

2.2 Constellation Architecture

The core part of Polaris is the multi-agent LLM constellation. As opposed to a general large model performing many different tasks with varying performance, we break down the faculty into medium-size specialist models with consistent high performance. This is beneficial in settings like ours to optimize for tasks with competing objectives, while keeping the latency low for a real-time application. For instance, this allows us to optimize the primary language model for conversational fluidity, rapport-building and empathetic bedside manner with patients.

Simultaneously, we focus on developing specialized agents with capabilities that serve a dual purpose: (a) assist the main LLM with relevant healthcare-specific context (e.g., **A patient with chronic kidney disease stage 3b is recommended to avoid common NSAIDs like advil, motrin**); (b) provide a layer of safety to double-check the information provided by the main model (e.g., **40mg of Lasix twice a day vs. 80mg of Lasix once a day**). Each of these agents is tied to specific abilities such as confirming that the patient’s reported medicine consumption aligns with the dosage prescribed by their physician, determining whether the patient’s reported OTC drug consumption is within the manufacturer’s recommended range, understanding which medicine the patient is referring to (i.e., patients often struggle with correctly pronouncing drug names, and may confuse ones that sound similar) retrieving current and historical lab results from the patient’s EHR, guiding the patient’s food choices to align with their prescribed diet, etc.

While Polaris has many such agents, we focus on only a subset of them for illustration. Finally, these agents co-operate to solve a complex task. For instance, the patient utterance "I am taking dulaglutide. I saw my hemoglobin a1c was 5 - is that normal?" requires both the lab agent and the medication agent to work together with the primary agent to formulate the response "a hemoglobin a1c of 5 is within the normal range and lower than your previous value of 6.5. Your dulaglutide does decrease your hemoglobin a1c level".

Safety Benefits. The constellation architecture offers many safety benefits, one of the most significant being the enhancement of safety through increased redundancy and specialization. Each safety-critical task is performed by co-operation between the primary model and the corresponding specialist agents. This redundancy ensures that if the primary model fails or misses a task for any reason, the specialist agent ensures that the system continues to operate correctly and safely.

Another key advantage of the constellation architecture is the modular design which facilitates ease of maintenance and enables continuous system improvements. This modular structure allows for individual specialist agents to be updated, repaired, or replaced without disrupting the entire system. This not only simplifies maintenance but also accelerates the process of rolling out updates, thereby ensuring that Polaris is always equipped with the latest features and security measures. Continuously upgrading the modules further enhances the safety of Polaris, as it allows for swift responses to newly discovered issues as they arise. This is particularly important in the clinical domain where we can update the specialist modules with new information and knowledge about drugs, interactions, policies and protocols without requiring to retrain the entire system.

System architecture. Real-time conversations are especially sensitive to system latency including ASR, LLMs and TTS. The LLM constellation architecture also helps reduce end-to-end latency by allowing all the agents in Polaris to run concurrently, including the primary agent. The support agents can either be synchronous or asynchronous with respect to the primary agent’s activity. The synchronous setting is used when the primary agent’s response must be conditioned on the context provided by the support agent, for instance, in the case a retrieval agent detects the need for a database lookup to provide relevant information for the query – resulting in a higher latency. However, the average latency is lower since the support agents need not be invoked for every user utterance. In the asynchronous setting, certain agent functionalities will not block the primary agent, but their results (contexts) will be made available to condition the primary agent response in a subsequent turn, once the asynchronous agent has completed its task. Polaris also implements a form of ‘garbage collection’ on the support agent’s messages to the primary agent. The system removes stale tasks in order to not overburden the primary agent with outdated instructions not relevant at the current timestep. Finally, the constellation is pre-emptible with respect to new user utterance. The input is processed immediately after the system detects an end to voice activity, not relying on any explicit signal from the user. However, if additional speech is detected, constellation processes are restarted to use the updated information.

From a systems perspective the constellation paradigm is extremely helpful for keeping latency low while still achieving robust multifaceted reasoning. Furthermore, each agent can be trained, updated and deployed independently. The modular paradigm allows for the model(s) in each agent to be customized and even switched for models of different sizes and runtime settings (e.g. quantization, temperature, sampling parameters, etc.) without affecting other models in the constellation. The agents are also trained to emit structured outputs, which are parsed by the code harness of Polaris to execute verifiable control flows, reduce hallucinations and ensure a consistent interface with the primary agent. While this system allows for efficiency through sparse activation of agents and a reduction of average latency, it also creates additional challenges. These include (a) false positives in the support agent event detectors that introduce noise to the primary agent; (b) resolving inter-agent conflicts where multiple agents can potentially introduce conflicting tasks; and (c) added complexity when tuning the primary agent to resolve and follow relevant support agent-provided tasks. We discuss our mitigation efforts for these in subsequent sections.

2.3 Orchestration

The internal state of Polaris evolves as a result of the primary agent’s interactions with the user and the message passing conducted by the support agents. To explain this more clearly we adopt the following notation for the system components:

- Agents: Let $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$ be the set of agents, where A_1 is the primary agent and A_2, \dots, A_n are the support agents.
- Conversation History: Let $H_t = \{(u_1, r_1), (u_2, r_2), \dots, (u_{t-1}, r_{t-1})\}$ be the conversation history up to turn t , where u_i is the user utterance and r_i is the response of the primary agent at turn i .
- Agent-Specific Prompts: Let $P_{A_i}(H_t)$ be the prompt for agent A_i at turn t , derived from the conversation history H_t .
- Support Agent Outputs: Let $O_{A_i}(P_{A_i})$ be the output of support agent A_i given prompt P_{A_i} . The output space is structured with expected fields, e.g., $O_{A_i} = \{f_1 : v_1, f_2 : v_2, \dots\}$, where f_j are the fields and v_j are the values.
- State Changes: Let S_t be the state of Polaris at turn t . The state changes based on the outputs of the support agents, i.e., $S_{t+1} = \text{Update}(S_t, O_{A_2}, O_{A_3}, \dots, O_{A_n})$.
- Prompts: The prompt for the primary agent at turn $t+1$ is a function of the updated state and the conversation history, i.e., $P_{A_1}(H_{t+1}, S_{t+1})$.
- Short-Horizon Tasks: These can be represented as additional prompt snippets, e.g., $T_{t+1} = \{\text{task}_1, \text{task}_2, \dots\}$, which are appended to the primary agent’s prompt based on the state changes.

Polaris mediates the message passing between agents with deterministic imperative programming. This allows us to keep the output space of the support agents well-constrained and implement control-flow logic to support customized protocols. The tasks which are given to the primary agent are systematically formatted and use synthesized information from multiple agents when applicable. This framework facilitates a cooperative approach to complex dialogue generation, where multiple agents contribute to the overall conversation dynamics. An overall step of the constellation system is described in Algorithm 1.

2.4 LLM Details

Our LLM architecture designs and training choices are constrained by the requirement to serve the model in real-time as part of our LLM constellation. As a result, we employ a diverse set of LLMs for the support agents ranging from 50B to 100B parameters, whereas the primary agent is always a medium-size LLM (70B - 100B parameters). All of our models follow a decoder-only transformer-based architecture, and contain between 30 and 100 layers. We use Grouped Query Attention (GQA) [18] to achieve faster inference speeds and reduced memory footprint when storing the KV cache during inference, which in turns allows us to serve larger batch sizes for increased throughput. Different models use different tokenizers depending on the use case and the need for specialized medical terminologies (such as drug names), with vocabulary sizes ranging from 30,000 to 200,000 tokens. The implementation of the attention mechanism employs Flash Attention 2 [19], which makes the training and inference stages faster, with lower memory footprint and more efficient scaling to longer sequence lengths. All of our models use RMSNorm normalization layers [20], SwiGLU activation functions [21], and Rotary Positional Embeddings (RoPE) [22]. The maximum context window varies from model to model, ranging from 4,096 to 32,768 tokens. The training runs were performed using a distributed, multi-GPU setup, using several hundred Nvidia H100 GPUs with DeepSpeed [23].

The primary agent is trained in 3 stages: General Instruction Tuning, Conversation and Agent Tuning, and RLHF. The training and dataset details for the primary agent are discussed in Section 3. The rest of the models are trained using human-labeled datasets, as defined in the respective sections for each agent.

During inference, we deploy some of our models in bf16 precision and others in int8 precision (using weight-only quantization), depending on the latency and accuracy requirements.

Algorithm 1 Multi-Agent Dialogue System with Message Passing

```
1: Input: Conversation history  $H_t$ , current state  $S_t$ , agents  $\mathcal{A} = \{A_1, A_2, \dots, A_N\}$  where  $A_1$  is the primary agent
2: Output: Updated conversation history  $H_{t+1}$ , updated state  $S_{t+1}$ , primary agent response  $r_{t+1}$ 
3: Initialize:  $S_{temp} \leftarrow S_t$ ,  $T_{t+1} \leftarrow \emptyset$ , Messages  $\leftarrow \{\}$ 
4: for  $i = 2$  to  $N$  do
5:    $P_{A_i} \leftarrow \text{CREATESUPPORTAGENTPROMPT}(A_i, H_t, S_t)$ 
6:   Messages[ $A_i$ ]  $\leftarrow \text{SENDMESSAGE}(A_i, P_{A_i})$ 
7: end for
8: for  $i = 2$  to  $N$  do
9:    $O_{A_i} \leftarrow \text{RECEIVEMESSAGES}(A_i, \text{Messages})$ 
10:   $S_{temp}, T_{A_i} \leftarrow \text{PROCESSOUTPUT}(O_{A_i}, S_{temp})$ 
11:   $T_{t+1} \leftarrow T_{t+1} \cup T_{A_i}$ 
12: end for
13:  $S_{t+1} \leftarrow S_{temp}$ 
14:  $P \leftarrow \text{CREATELEADAGENTPROMPT}(H_t, S_{t+1}, T_{t+1})$ 
15:  $r_{t+1} \leftarrow \text{GETOUTPUT}(P)$ 
16:  $H_{t+1} \leftarrow H_t \cup \{(u_{t+1}, r_{t+1})\}$ 
17: return  $H_{t+1}, S_{t+1}, r_{t+1}$ 
18: procedure  $\text{PROCESSOUTPUT}(O_{A_i}, S)$ 
19:    $S' \leftarrow \text{UPDATESTATE}(S, O_{A_i})$ 
20:    $T_{A_i} \leftarrow \text{EXTRACTTASKS}(O_{A_i})$ 
21:   return  $S', T_{A_i}$ 
22: end procedure
```

3 Conversational Alignment

In the healthcare domain, the development of a patient-focused clinical framework requires an intelligent conversational agent. This agent has to be designed to foster empathetic engagement with patients, skillfully addressing their concerns and assessing health conditions, be a motivational coach and adopt the multifaceted role of a dietary mentor. It has access to all relevant patient information and medical history to reason and navigate through a complex clinical conversation.

3.1 Data

The development of conversational agents capable of engaging in meaningful and accurate clinical discussions represents a significant challenge. This challenge is compounded by the lack of datasets specifically tailored for training models in the nuanced context of healthcare conversations. Our work addresses this critical gap by leveraging a unique compilation of dialogues, including simulated interactions between registered nurses and patient actors. These dialogues form the cornerstone of our approach to developing a conversational agent with the proficiency to navigate a wide array of clinical discussions. To enhance the agent’s reasoning capabilities, instruction-following proficiency, and domain-specific knowledge – we augmented these datasets by introducing diverse kinds of instructions and tasks geared for multi-hop reasoning. For instance, given a user utterance “my shoes do not fit”, the agent should be able to reason about “swollen ankle \rightarrow sign of fluid retention \rightarrow sign of exacerbating CHF conditions”.

3.1.1 Foundation Model Training Data

We first train our foundation model on a massive collection of proprietary data including clinical care plans, healthcare regulatory documents, medical manuals, drug databases, and other high-quality medical reasoning documents. The objective of this phase is to incorporate fine-grained medical knowledge, reasoning and specialized numerical reasoning (e.g., dosage calculations). We further annotated some of the medical datasets with reasoning chains to

Agent Name	Sub-models	Training Epochs	Inference Precision
Primary Agent	Primary	2 each stage	bf16
Checklist Specialist	State Model	2	int8
	Privacy Model	3	int8
Medication Specialist	Detector	5	int8
	Evaluator	5	bf16
Lab & Vitals Specialist	Detector	2	int8
Policy Specialist	Detector	2	int8
	Retriever	5	-
Clinical Nutrition Specialist	Detector	2	int8
	Evaluator	2	int8
EHR Specialist	EHR Summary Model	2	int8
Human Intervention Specialist	Detector	3	int8
	State Model	2	int8
	Evaluator	3	int8

Table 1: Training and inference details of our models. Each specialist agent is composed of multiple LLMs.

Features	Data Source	PT	Instruct Tuning	Conv Tuning	Agent Tuning	RLHF
Medical Knowledge and Reasoning	Proprietary Data	✓	✓		✓	
General and Quantitative Reasoning		✓		✓		
Rapport-building, Bed-side Manner, Polished Conversational Style	Proprietary and Simulated Conversations			✓		✓
Mitigating Hallucination, Fact Verification, Knowledge Augmentation	Policy, manuals, clinical references, curated lists				✓	✓
Navigating Complex Care Protocols, Responding to Complex Clinical Tasks	Simulated Conversations			✓	✓	✓

Table 2: Overview of data and features in different stages of training Polaris .

further enhance the medical reasoning capabilities of the model. Figure 32 shows one of the examples of the medical question answers we leveraged for training our foundation model.

3.1.2 Simulated Conversations

The nature of our conversational objectives requires specialized data. The agent should be able to maintain conversational state over relatively large time horizon; a multi-turn call in our setting often exceeds 20 minutes with several dozen turns from each speaker. It is important to generate data which portrays tougher trajectories within the healthcare setting of our calls. Our system must be robust to handling patients with diverse profiles, for instance, high engagement with many questions or concerns; low regimen compliance; skepticism of AI health services, etc.

We thus rely on the domain expertise of medical professionals, US-licensed nurses in our setting, to generate accurate data which exhibits comprehensive coverage of our desired data distribution. We leverage a large number of registered nurses, patient actors and our

Congestive Heart Failure Call Objectives

This call is a routine CHF Assessment
Verify the patient’s identity - First name, Last name, DOB
Review the patient’s current medications (new + old, excluding discontinued) with the patient, including dosage
Review patient’s diet, educate the patient on following low sodium diet and fluid restriction if necessary
Ask the patient if they have any new or worsening shortness of breath with activity, laying down, or at rest?
Ask the patient if they have any ongoing or worsening chest pain
Ask the patient if they have any new or worsening dry cough
Ask the patient if they have increased swelling in their feet, ankle, stomach or legs
Ask the patient to weigh themselves daily
If they weigh themselves regularly, ask the patient if they have gained more than 2 pounds in the past day or 5 pounds in the past week.
Ask the patient if they are using more pillows at night to sleep
Inform the patient that we just performed a simple CHF assessment and encourage them to do it at home every day
Ask the patient about physical activity and coach them to improve their routine if warranted

Table 3: Sample instructions for Human Nurses or the AI to engage in a call with a Patient Actor or an AI (instructed to act like a Patient) with Congestive Heart Failure.

AI-in-the-loop to generate simulated conversations. For targeted medical conditions (e.g., CHF, CKD) and procedures (e.g., post-discharge, pre-operative, chronic-care followups), clinicians create conversational scripts (refer to Table 3 for a simplified version of our CHF script). We also create a large number of fictional patient profiles with different medical histories, condition severity, medications, labs, lifestyle and personality traits (refer to Table 4 for a sample patient profile).

We require clinical expertise to train all parts of Polaris. In order to do this, we employ a bootstrapping strategy for data generation. Initially, we generate conversations using registered nurses and patient actors. We design these conversations to cover a wide breadth of scenarios and behaviors on the part of the simulated patient. We give specific instructions to patient actors to converse such that many of the specialist agents are triggered, such as asking complex questions about medications, mispronouncing medicine names, providing deliberately confusing lab results (e.g., a blood glucose reading in response to a request for their blood pressure reading), asking about hospital specific information (e.g., where to park), etc. During these conversations the registered nurse annotates the transcripts for when these events occurred and what the specialist agents should do during corresponding turns.

We then proceed to train both the primary and support agents from this collection of datasets. This allows us to start using our trained agents in the conversation instead of nurses; their responses are reviewed by registered nurses with potential re-writes for the noisy turns. We also employ instruction tuned language models to perform data cleaning and augmentation. This is used to improve the style of responses and also create additional variety, preventing mode collapse in the primary agent distribution and ensuring robustness in the support agents. Like most bootstrapping situations, this process becomes more effective with more rounds of training, which allows generalization to new scripts, medications, and so on.

Here we provide more details on the entities which participate in the data generation process.

Patient Actors. Here the actors are asked to play the role of a fictional patient and portray a realistic patient experience. They are encouraged to use the fictional patient’s background as well as their personal experiences and feelings to guide them. We provide them with sample instructions, for instance: Ask about lab values - what they mean, whether they’re normal or not; state they’re taking the wrong medication dosages, or that they are taking medications not on their chart, and so on.

Patient Conditions	Medications	Prior Lab Results
Name: Mary Adams; DOB: 1950-01-01; Gender: Female; Condition: Congestive Heart Failure	Lasix 20 mg oral tablet; Potassium 20 milliequivalent oral tablet; Lisinopril 10 mg oral tablet; Tylenol 500 mg oral tablet	inr: 1.0; glucose: 143.0 mg/dL; hemoglobin: 9.4 g/dL; hematocrit: 29 %; wbc: 5.0×10^9 /L; platelets: 231.0×10^9 /L; mcv: 81 fL; pcv: 0.48 fL; red blood cell count: 4.1×10^{12} ; mch: 28.0 pg

Table 4: Sample fictional patient profile with medications and labs. During simulated conversations, Patient Actors, Registered Nurses as well as the AI have access to these personalized details for an in-depth clinical conversation.

Registered Nurses. We ask the nurses to play the role of an ideal nurse at the bedside, on the phone, or in their community. While we provide them with the conversational script with defined call objectives, they are encouraged not to follow the order strictly, but instead follow the natural conversation trajectory for a realistic experience. This allows the patient at times to go on tangents, for instance, sharing their personal experiences that the nurse can engage with to develop rapport and trust; share their health concerns and symptoms that the nurse can both empathize with and educate on.

Our AI-in-the-Loop. Once our models are tuned with conversations generated from the earlier phases, we further use them to generate synthetic conversations. Here, the AI plays the role of a nurse with the conversational script and meta instructions used in the prompt as a preamble to converse with the patient actors. This stage is required primarily for tuning our support agents (discussed in Section 7) that generate tasks and provide relevant context to condition the primary agent to generate the response. After our support agent system is sufficiently tuned, we proceed to RLHF data generation. For this we sample multiple responses from the primary agent for the full prompt constructed from the meta-prompt, conversation history and support agent tasks. We leverage our registered nurses to give preference feedback to perform RLHF on the primary agent.

3.1.3 Conversational Dataset Construction

The complexity of clinical calls and multi-turn conversations informed our approach to dataset construction. Most instruction tuning regimes involve single or few-turn interactions between the user and assistant. In our case, most calls consist of several dozen turns from each call participant. We convert a conversation with τ turns from each participant into τ training data points, where for each $t \in 1 \dots \tau$, we create a prompt with the conversation history up to turn $t - 1$ with meta-instructions in the preamble, mask the prompt, and train the primary agent to predict the content of the nurse turn t . We apply a standard cross-entropy next-token loss objective for this supervised training.

A standard practice for language model training on multi-turn conversation consists of applying loss masks to all of the user turns and training on all of the assistant turns within a single training data point. This approach has the benefit of increased training compute efficiency when compared to splitting up all of the turns into individual training data points as we do. However, in our case this approach is not applicable because we expose a conversation state description S_t to the primary agent which evolves as a result of the other agents in the constellation system $\{A_i\}$. Aside from the conversation history H_t growing turn by turn, the primary agent system prompt is changed due to the effect of the checklist specialist (see Section 4.2). Furthermore, the other clinical specialist agents inject short-horizon tasks and instructions T_{A_i} . These are placed with their own format, and expire based on heuristics inside of the task orchestration system. We thus create a separate prompt for each turn which is a function of $(H_t, S_t, \{T_{A_i}\})$ that has its loss masked, and we train only for the labeled response r_t . Figure 31 shows selected dialogues between an anonymized nurse and patient, our approach grounds our healthcare agent in authenticity, drawn from real-world

Agent: Great, thank you for taking the time to speak with me. Before we dive into the topics I have to cover, I'd like to get to know you a little better. Is there anything you're looking forward to doing this week?

Patient: Actually, yes! I am looking forward to a little gardening. Its something I love and it's very relaxing for me.

Agent: That sounds wonderful, George! Gardening can indeed be very therapeutic. What kind of plants do you usually tend to in your garden?

Patient: I enjoy growing a mix of things. Vegetables, herbs, and a few flowers. It's nice to grow my own healthy food and see everything come to life.

Agent: That's fantastic! Growing your own food can be so rewarding. What's your favorite vegetable or herb to grow? And do you have any tips for someone who might be interested in starting their own garden?

Patient: My favorite herb is basil. It's versatile and grows pretty easily. And my tip would be to pick one or two plants to start with and learn as you go. Gardening is a process of trial and error, so don't get discouraged.

Agent: Basil is indeed a wonderful herb to grow and use in cooking. And thank you for the gardening tip. Starting small and learning as you go sounds like a wise approach. Do you use your home-grown basil in any special recipes?

Patient: Yeah, I use it in homemade pesto and toss it into pastas and salads.

Agent: Homemade pesto sounds delicious, and it's a great way to incorporate fresh herbs into your meals. Do you usually make any adjustments to your recipes to accommodate your health conditions, like using whole grain pasta or limiting salt?

Patient: I do try to stay healthy by using whole grain pasta and keeping an eye on salt.

Agent: That's a great approach, George. Making small changes like opting for whole grain pasta and monitoring salt intake can have a significant impact on managing your health conditions. Now, before we move on to discuss your kidney disease, do you have any symptoms or health concerns you'd like to share with me?

Patient: No specific concerns, I'm managing quite well right now.

Agent: I'm relieved to hear that you're managing well. Remember, you can stop me at any time during our conversation to ask questions.

Figure 3: An excerpt from a conversation we use as a part of our conversation tuning dataset. Within this exchange, the agent builds rapport with the patient, employing empathy and encouragement to influence the patient's healthy dietary choices.

interactions. This strategy enriches the training process, ensuring our agent is deeply rooted in the nuances of genuine healthcare encounters.

3.2 General Instruction Tuning

In this step, we also focus on honing the model's ability to follow instructions by augmenting our proprietary data with complex tasks and instructions including multi-hop logical reasoning, mathematical reasoning, and task completion. The objective of this stage is to improve the model's capability in executing intricate tasks within the medical domain like clinical assessment.

3.3 Conversation and Agent Tuning

Conversation Tuning: During this phase, we train the model using simulated dialogues between registered nurses and patient actors, alongside a spectrum of general, unstructured conversations. The objective of this stage is to increase the communication capability of the model in terms of refined bedside manners, empathy, rapport, a polished conversation style like a healthcare professional. This enrichment process enhances the model's grasp of complex conversational nuances found in real-world healthcare settings. Real-world healthcare conversations imbue the model with the characteristic warmth, understanding, and professionalism needed to bridge the gap between automated systems and human touch in patient care.

Agent Tuning: We also perform specialized alignment or agent tuning focusing exclusively on clinician-patient dialogues in this step. This targeted training sharpens the model's

Patient: Most days. Not every day, but most days.
Agent: I'm relieved to hear that you're keeping an eye on your blood sugar levels. It's a key part of managing your diabetes. Can you recall what your most recent blood sugar reading was?
Patient: You know, I'm a little embarrassed to say it was a hundred and thirty eight.
Agent: No need to feel embarrassed, Mr. Hernandez. We're all human and it's perfectly fine to have off days. Can you recall if you took this reading before eating or after you had a meal?
Patient: I took it actually before I ate anything.
Agent: I appreciate you letting me know that. A fasting blood sugar level of one thirty-eight is slightly above the usual range, which is between sixty and one oh five. This could be due to a variety of factors, such as your diet, stress levels, or even the timing of your medication. Let's work together to figure out what might be causing this. Have you noticed any specific patterns or triggers that seem to affect your blood sugar levels?

Figure 4: A snapshot from a conversation between a patient actor and Polaris, providing essential information with empathy, appropriate bedside manner and clarity in a simulated patient-AI interaction.

proficiency in handling specific scenarios anticipated within our system's operational context. A significant aspect of this stage involves training the model to adeptly navigate and adhere to intricate checklists tailored to individual patient conditions and prescribed care protocols, where we give considerable attention to both medical and non-medical exchanges.

During this phase, we engage in *iterative co-training* between the primary agent and support agents, where the primary agent is trained in tandem with support agents. The support agents are trained to pay attention to signals in the primary agent's conversation with the user. And conversely, the support agents communicate with the primary agent by modifying its prompt, and injecting relevant information and instructions.

Figure3, Figure4 and Figure33 show excerpts from our simulated conversations between the patient actors and Polaris.

Agent Tuning is the hardest part of the entire Conversation Agent training pipeline. It is quite challenging to train the constellation of the LLMs which communicate together with the primary agent. Figure 5 shows how conflicts can arise between various agents. In Figure5, the primary agent is provided with two tasks from the Lab & Vitals Agent, and the Human Intervention Agent. When multiple agents add a task for the primary agent to follow, the primary agent resolves the conflict by prioritizing the Human Intervention Agent and giving a safe and informative response to the patient. To train such a system we leverage actual physicians and nurses who label the data and assimilate the information from various agents to provide the most appropriate response to the patient.

Self Training: A typical challenge with training the AI on simulated conversations is that of generalization. Even though we train our models on diverse scripts and conditions to cover a wide range of scenarios, it still suffers from exposure bias [24]. This results in the AI producing similar responses in the conversation, increase in repetitions and lack of robustness to diverse scenarios it never encountered during training. This results from the way we instruction tune these models with next token prediction during training with teacher forcing; where the model produces the next token conditioned on ground-truth token in the previous step. However, during inference, there is no ground-truth and the model samples a token conditioned on self-predicted token in the previous step – creating a training and inference mismatch referred to as exposure bias. To address this, we train the AI on self-generated responses during training. To this end, we train the model for a few steps; use the current model parameters to produce responses on sampled training prompts with temperature sampling to augment the original training data; and then re-train the model. This ensures the model encounters diverse training samples in different iterations and also aware of its own uncertainties during sampling at inference – thereby minimizing the distribution shift.

Agent: You're doing a great job managing your diabetes. Let's shift our focus to your blood pressure. Can you recall when you last checked it?

Patient: Its 140 over 90.

Agent: Is there anything else you'd like to tell me about your blood pressure or any other health concerns?

Patient: No. I went out to eat last night, and I had some miso soup, probably that's why its up.

<Lab And Vitals Agent TASK>
 For 'Systolic' The normal range for blood pressure is 90 to 120. The patient's systolic reading is a high.
 For 'Diastolic' The normal range for blood pressure is 60 to 80. The patients diastolic reading is a high.

<Human Intervention Agent TASK>
 TELL the patient you have taken a note of the patient's issue that the patient has high blood pressure and will let their care team know.

Agent: Okay, Mr. Hernandez. It's not uncommon for certain foods, particularly those high in sodium like some miso soup, to cause a temporary rise in blood pressure. I've made a note of this and will pass it along to your care team. Let's proceed to your recent lab results. Do you have any questions or concerns about them?

Figure 5: A snippet from a simulated conversation between a patient actor and Polaris. This snapshot shows how the multiple task agents interact with the primary agent to provide a safe and informed response to the patient. Note: the snippets labeled with TASK are purely internal to Polaris and unseen to the user. They are shown here simply to illustrate how the primary agent assimilates distinct tasks inserted by the various support agents to form a cogent response that is safe, informative and relevant to the utterance.

3.4 RLHF

The final phase incorporates Reinforcement Learning from Human Feedback (RLHF) by engaging directly with healthcare professionals. This crucial step refines the model's performance across key dimensions, including safety, empathy, and the overall quality of care-related conversations. Through this iterative process, we ensure the model's responses are not only accurate but also aligned with the nuanced requirements of empathetic patient care.

We gather preference data from the same pool of nurses that were involved in the conversation data generation. Nurses are presented with the overall context, the conversation history and any active tasks provided by the support agent system. They are then asked to rank multiple candidate responses in their order of preference. This is nontrivial as there are multiple factors that determine the quality of a response. We ask them to prioritize medical safety and accuracy above everything else. Given that, responses are judged for empathy, bedside manner, level of patient education, and other conversational features. Within this, we suggest heuristics for what would be considered 'acceptable' with regards to conversational fluency. If the nurse judges none of the responses to be acceptable, we ask them to write an ideal response. Importantly, the rewrite must still conform with the tasks given by the support agents in order to not break the cooperative nature of the constellation system. From the ranked responses we construct pairwise preference samples. We exclude pairs where the preferred response was deemed unacceptable. When there are available rewrites, we only create pairs where the rewrite is preferred and all model responses are not preferred.

4 Specialist Support Agents

A patient-facing model must be able to robustly perform many technical functions while in conversation, such as verifying prescription adherence and reviewing lab analysis results. These functions must properly contextualize the patient's medical history and successfully evaluate against a vast space of medical 'edge cases.' Furthermore, they should be invoked at appropriate trigger points from the conversation flow. We align our design with the intuition

that language models exhibit better reasoning when given targeted instructions compared to overloading the context with extraneous information. The primary conversation agent is required to maintain a large context in order to achieve optimal conversational fluency, but much of this is irrelevant for the medical reasoning subroutines. Hence we implement a modular system architecture where the primary conversation agent is augmented with support agents that serve specialized functions to ensure high reliability and safety.

4.1 Privacy & Compliance Specialist

An essential aspect of patient-facing medical calls is verifying the patient’s identity. It’s important to do this verification before any personal information on the patient is loaded into the LLM context to prevent any leaking of information. It is also paramount that the identity verification has an extremely low error rate and no hallucination. To protect patient privacy and comply with privacy rules, the system must verify that the person answering the call is the intended recipient, before engaging further. In order to do this, the Privacy & Compliance Specialist confirms identifiers including full name and one of date of birth or medical records number. This agent can easily be trained to verify additional identifiers as appropriate (e.g., authorized care manager, parent, etc.).

Until the recipient’s identity has been confirmed, the system should not share such PHI with the recipient. However if the primary agent’s prompt contains salient details of a patient’s Electronic Health Records to support discussion of medications, lab values, etc., it is vulnerable to disclosing some of the information to the patient (e.g., by red teaming).

One solution is to use a higher-order state transition to ensure that the primary agent has no PHI in its context at all until the patient’s identity has been verified. During the initial state, the primary agent requests to speak with the intended recipient. When the target recipient becomes available, the primary agent requests that the patient provide the required patient identifiers. During these first two stages, the primary agent does not have the correct DoB or MRN in its context. Therefore, to verify the patient-provided identifiers, the primary agent forwards them to a second component that checks these details against the protected database.

If the second agent finds that the identifiers are incorrect, it returns a non-PHI response specifying failure that prompts the primary agent to re-verify. After a few failed retries, the primary agent terminates the call. Once the second agent finds that the identifiers are correct, the main phase of the conversation begins. This is depicted by the transition from the yellow to purple regions of the diagram in Figure 6. At this point, PHI such as EHR details are loaded into the primary agent’s prompt.

4.2 Checklist Specialist

LLM’s struggle to follow complex instructions. This gets more challenging in our setting with complex call objectives and a long checklist of ordered instructions. In order to have a natural flow of conversation, we allow the user to go on tangential discussions while accomplishing a particular task in the checklist. This makes it challenging to bring the conversation back and resume from the pending checklist items. To address this challenge, we devised the checklist specialist to guide through complex call objectives.

To make this tractable, we leverage synthetic scripts with ordering of tasks as part of a call specification. We then train an LLM specifically to assess the conversation state concurrently with the primary agent model. This support model determines which tasks have been completed and correspondingly updates an internal data structure. Changes in this data structure propagate into the prompt for the primary agent by discarding objectives that have been accomplished and popping from a queue of objectives that are pending. This mechanism thus acts as a form of higher-order attention. We find this to be highly effective in ensuring the primary agent achieves all of the pre-specified objectives, especially when the patient asks additional questions or goes off on a tangential discussion, not necessarily related to the objective in question. Furthermore, this system bounds the number of input tokens related to the tasks so that very long scripts e.g. comprehensive health questionnaires are still manageable.

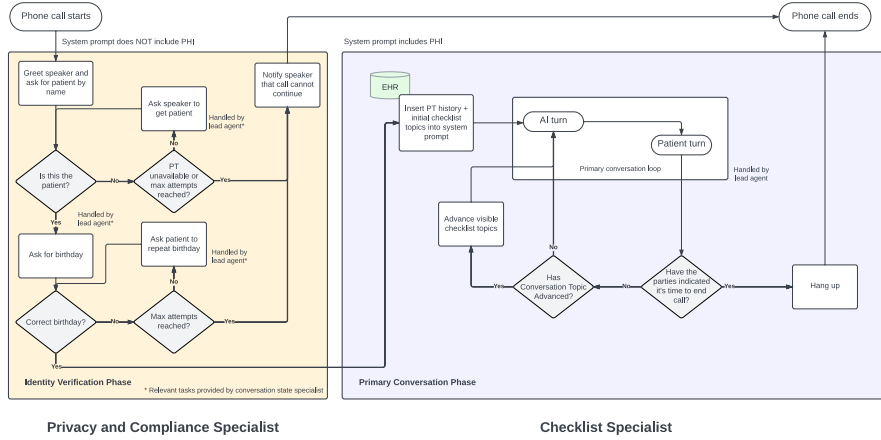


Figure 6: Privacy & Compliance and Checklist Specialists. These support agents facilitate key high-order state transitions for the system, ensuring the primary agent is appropriately constrained and focused on the tasks it needs to accomplish.

The evaluation of states can be challenging when the individual tasks can be complex and multifaceted. For example, a simple task might be to confirm with the patient the date, time and location of their upcoming procedure. Once the agent states these details and the patient acknowledges they understand, then the task is complete. However, there are other, more involved tasks such as going over a patient’s diet and giving appropriate guidance on it. In this scenario it is plausible that the primary agent gives substantial advice on typical breakfast, lunch and dinner and is ready to move on. However, the checklist specialist may not share this belief, creating a deadlock. We could choose to have both models decode chains-of-thought and perform reconciliation, but this and similar strategies break the latency constraints imposed in the voice-based setting for our system.

Instead, we adopt a strategy of letting the primary agent decide when to move onto the next set of tasks as it deems appropriate. The checklist specialist is thus tasked with passively documenting the transitions rather than actively enforcing them. This removes the possibility of deadlock and also ensures the primary agent’s conversation does not follow the scripts exactly verbatim, but allows for a natural flowing conversation. To further ease the role of the checklist specialist, we organize the scripts into sections, and prompt the specialist to assess whether the primary agent has moved from a given section to another one – where a section comprises of a set of related tasks and instructions. This allows us to manage states across long conversations with several complicated tasks. This specialist also plays a similar role during a potential human intervention scenario as there is a checklist describing primary agent protocol. Thus, it is technically the owner of the Intervention State Model depicted in Figure 22. Finally, this agent also decides when to terminate the call as appropriate – a vital functionality for real-world scenarios.

4.3 Medication Specialist

The Medication Agent supervises and enhances medication-related interactions between the primary agent and the patient. This agent is developed with instructions from licensed medical professionals, clinical guidelines and standards of care to ensure patient safety and adherence to prescribed therapies. The agent is powered by two fine-tuned large language models and physician-approved reference tables for a wide number of medical conditions. This agent has several capabilities: precise medication identification, dosage discussion including for OTC drugs and contraindications of OTCs (based on the manufacturer’s guidance).

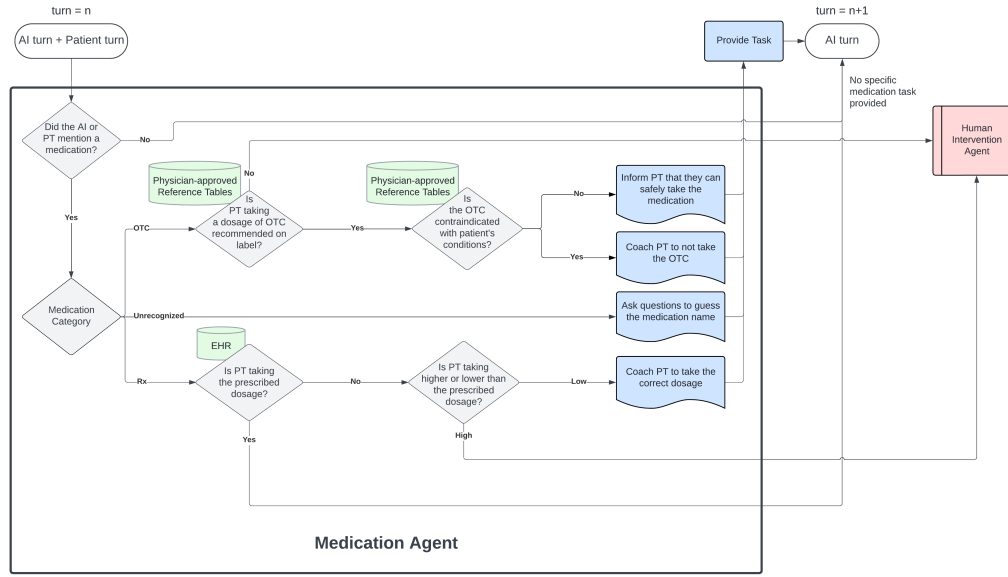


Figure 7: Medication Specialist Workflow

4.3.1 Prescription Adherence

A key challenge in healthcare is ensuring patients adhere to their prescribed medication regimen. General large language models are susceptible to suggestions and struggle with mathematical reasoning. Patients will often misstate how much of the medication they should take. The specialist agent ensures dosage values in the Electronic Health Record (EHR) are enforced. The Medication Agent assists the primary agent in this task by detecting errors such as forgotten doses, depleted medication supplies, inappropriate compensatory dosing, discrepancies in medication frequency, misunderstandings of pill strength, and the incorrect quantities of medications taken. Such challenges are especially common among chronic patients managing multiple medications, highlighting the necessity for the agent's precise detection and intervention capabilities.

The Medication Specialist agent is equipped to detect deviations from prescribed regimens, whether it be taking medication less frequently than prescribed; the patient adjusting the strength of the dose; or altering the number of pills ingested. Upon identifying such discrepancies, the agent guides the patient towards the correct regimen using motivational interviewing tactics, and involving the human care team through the Human Intervention Specialist as shown in Figure 8.

Patient: I've been taking my Digoxin every morning. Though, I'm not always sure if I took it, so I might take another just to be safe.
<Medication Agent TASK> The patient is taking more than the prescribed dosage for Digoxin. The correct dosage information for Digoxin:
 Dosage: 125mcg Frequency: Once a day
 Instruct and educate the patient to take the correct dosage
Agent: Okay, Ms. Anderson. It's essential to stick to the prescribed dosage for Digoxin to avoid any potential side effects. You're supposed to take one tablet of Digoxin, which is one hundred and twenty-five micrograms, once a day. Taking more than this on a regular basis could lead to certain side effects. Can you tell me how many tablets you're taking each day?

Figure 8: Example interaction with the Medication Agent about Prescription Adherence.

4.3.2 Condition Specific Disallowed OTCs

Condition-specific advice for OTC medications is a critical feature of the Medication Agent. General large language models tend not to be aware of OTC contra-indications for specific conditions. In many conditions, common over-the-counter medications and supplements can be harmful. Our agent allows providers to specify contra-indicated OTCs and will recognize their usage. For example, if a patient expresses a desire to take Ibuprofen for a headache, the agent cross-references this choice with the patient’s medical profile. Should the patient have Stage 3B Chronic Kidney Disease for which Ibuprofen is contraindicated as per the Physician Approved Reference Tables, the agent generates a task signaling the need to avoid this medication, thereby preventing adverse reactions. Conversely, if the the physician approved care plan recommends this OTC for the patient’s condition, then they are informed they can take the medication. This is illustrated in Figure 9

Patient: Actually, yes. I’ve started taking Sudafed occasionally.
<Medication Agent TASK> According to their doctor, the patient should not take sudafed because they have chronic kidney disease stage 3b.
Agent: Thanks for letting me know about the Sudafed. It’s a common medication for allergies and cold symptoms. However, given your Stage IIIB CKD, it’s not the best choice for you. Sudafed can potentially raise your blood pressure and heart rate, which could put extra stress on your kidneys. Have you been experiencing any symptoms that led you to start taking Sudafed?

Figure 9: Example interaction with the Medication Agent about Condition Specific disallowed OTC recommendation.

4.3.3 Over-The-Counter (OTC) Toxicity

The Medication Agent includes a critical function for verifying the safety of over-the-counter (OTC) medication dosages, safeguarding against the risk of toxicity. Max OTC dosage calculation depend on many different factors including age, weight, composition (capsule, tablets, liquid, etc) and strength. General language models are not good at reasoning across these different variables.

When a patient reports taking an OTC drug, the agent cross-references the quantity consumed with the maximum allowable limits outlined in the OTC Drug Label (or a Physician Approved Reference Table). Similar to the prescription adherence capability, this capability involves analyzing the conversation to identify the amount taken, which the patient may convey in different ways. The agent guides the primary model to elicit necessary information and communicate to the patient the appropriate maximum dosage. In addition, if the amount reportedly ingested by the patient exceeds these thresholds, the agent triggers a task for the Human Intervention Specialist as shown in Figure 10

4.3.4 Unrecognized Medications Reconciliation

Drug names are complicated. Patients often struggle to pronounce or recall them correctly. The medication specialist agent is used to guide the patient through and disambiguate a recognition process. To address patient mispronunciation of drug names, or errors from the ASR system in understanding the patient’s pronunciation, the agent uses a workflow that includes questioning the patient to clarify medication names. This process helps in accurately identifying medications, ensuring reliable management advice despite mispronunciations or ASR inaccuracies.

Additionally, when a medication, whether OTC or prescription (Rx), is not listed on the patient’s EHR or physician approved reference tables, the agent recommends the patient to consult with a healthcare professional. This prevents the unauthorized or uninformed use of medications and prevents the primary agent from providing general advice, reinforcing the importance of physician-guided medication management.

Feature	Description	Example
Medication Detection	Identifying specific medication mentions in conversations	Detecting a mention of "ibuprofen"
Dosage Evaluation	Assessing if the mentioned dosage aligns with the prescription	Evaluating "taking 40mg instead of prescribed 20mg"
Permissibility Check	Determining if a medication is allowed based on the patient's conditions	Assessing if ibuprofen is safe for a patient with certain health conditions
Advisory Provision	Providing specific recommendations or advising to consult healthcare providers	Advising "Consult your doctor before taking this medication" if not found in the database
Reconciliation Support	Supporting medication reconciliation for undocumented meds	Initiating a reconciliation task for a new medication mention
Patient Education	Educating patients on correct medication usage and adjustments	Educating a patient on their physician's protocol for correctly addressing missed doses

Table 5: Representative Capabilities of the Medication Agent with Examples

4.3.5 Models

The Medication Agent is powered by two specialized LLM models: the Medication Detector and the Dosage Evaluator. The workflow of the specialist is shown in Figure 7.

- Medication Detector:** This model is trained to accurately identify and understand the context of medication usage in patient conversations. It utilizes the datasets generated from simulated conversations to recognize a wide variety of medication names, including those that are mispronounced or less commonly known, and to understand the context of the conversation. The Detector is passively listening to the conversation between the primary agent and the patient's every turn and identifies the list of all medication names if discussed and the context of the conversation. The key challenge here is to make the detector robust enough to identify both correct and misspelled medication names without being overly sensitive to false positives. To achieve this, the model is trained on a diverse dataset that includes a wide range of medication names, brand names, generic names, and common misspellings or phonetic variations. The output of the Medication Detector is a structured list of medications mentioned in the conversation, along with their associated context. This information is then passed to the Dosage Evaluator for further analysis of the patient's adherence and dosage safety.
- Dosage Evaluator:** The Dosage Evaluator focuses on analyzing the dosage and adherence patterns mentioned in conversations. It is fine-tuned to detect discrepancies between prescribed regimen and patient-reported usage, including errors in timing, strength, frequency, and quantity. The model is trained to use the structured information from the patient's EHR or Physician Approved Reference Tables as ground truth to identify potential issues with the reported dosages in the conversation by the patient. The model returns whether the dosage reported by the patient is CORRECT, HIGH or LOW with respect to the ground truth. The key challenge here is to develop a model that can accurately parse and interpret the

Patient: Actually, yes. I’ve been taking 800 mg of ibuprofen 6 times a day.
<Medication Agent TASK> A dosage greater than 3200 mg of ibuprofen could be harmful. Call the **Human Intervention Specialist**
Agent: I’m very sorry to hear that you’re taking such a high dose of ibuprofen. This could potentially be harmful to your health. Let’s get you connected with a nurse who can provide immediate assistance. Could you please hold for a moment?

Figure 10: Example interaction with the Medication Agent about OTC toxicity verification.

often ambiguous and incomplete dosage information provided by patients in natural conversation. Patients may use colloquial terms or omit key details when discussing their medication usage. The Dosage Evaluator needs to be robust to this variability and be able to infer the intended meaning. To address this, the model is trained on a large corpus of simulated conversations that cover a wide range of dosage scenarios, including both compliant and non-compliant usage patterns. The training data is augmented with techniques like synonym replacement, numerical perturbation, and random omissions to improve the model’s resilience to linguistic variation.

4.3.6 Training Data Generation

The Medication Agent utilizes two Large Language Models (LLMs) for medication and context identification and dosage analysis. To fine-tune these LLMs, datasets are collected through simulated conversations between the primary agent and human nurses acting as patients. These nurses are tasked with discussing a wide range of medications, both prescription and over-the-counter (OTC), that a patient might be taking. They are also instructed to simulate various types of dosage adherence errors commonly made by patients, such as taking incorrect doses, misunderstanding prescription instructions, and irregular medication timing. This simulation process is designed to generate a dataset that captures the diversity, complexity and nuance of real-world patient medication management scenarios.

4.4 Labs & Vitals Specialist

Patients often inquire about their medical lab test results and vital signs in clinical interactions [25]. For example, a diabetic patient might ask for their hemoglobin A1C result, which measures the average blood sugar levels over the past 3 months [26]. A patient with congestive heart failure might be interested in their most recent blood pressure values, or inquire about the overall status of their laboratory and vital sign results. In some other cases, patients might initiate discussion of lab values that are not present in the EHR. Therefore, it is necessary to design a dedicated specialist that assists with requests related to labs and vital measurements. Figure 11 contains an overview of the workflow of the Labs & Vitals Specialist.

4.4.1 Determination of Labs & Vitals and Value Extraction

To properly respond to any lab-related queries, the Labs & Vitals Specialist must accurately detect patient intent and associate the intent with a valid medical lab test or vital sign measurement. Furthermore, it must be capable of understanding and extracting the corresponding values of mentioned labs and vitals from the conversations. For example, if a patient states that her international normalized ratio (INR) result is around 1.1, this specialist should report the mentioned lab and its value to the system. Subsequently, this specific value is analyzed to deliver relevant feedback. See Figure 12 for some examples.

4.4.2 Normal Range Assessment & Plausibility Check

General language models lack medical grounding of lab reference ranges and tend to get confused due to the different reference ranges on the internet. Identifying the correct reference range for a patient’s lab value, given their age, gender, etc., is essential for accurate lab presentation and the avoidance of hallucination.

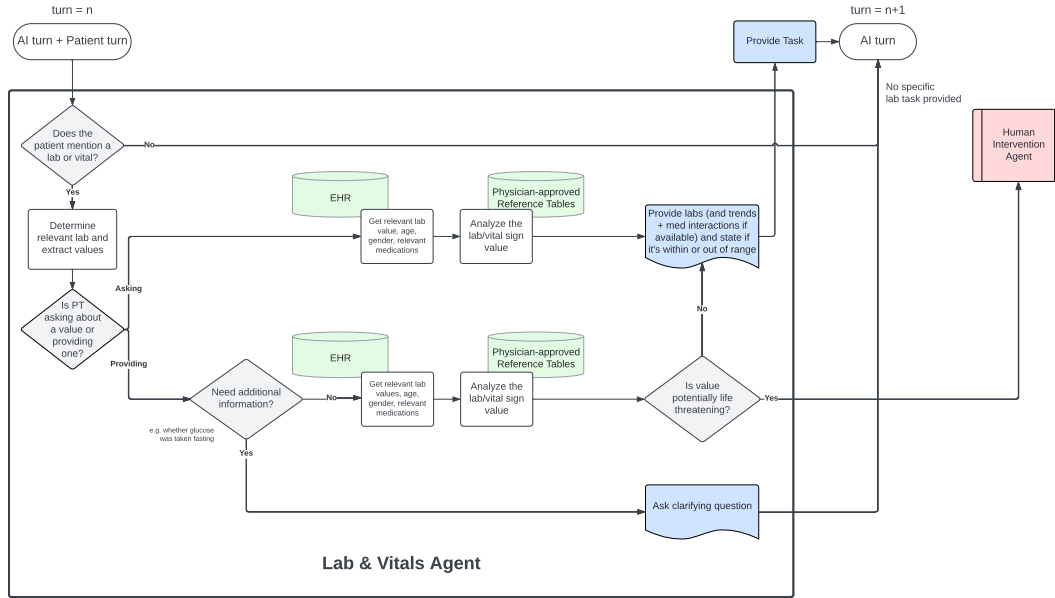


Figure 11: Labs & Vitals Specialist Workflow.

Patient: “My urine dipstick glucose test is positive.” → Urine dipstick glucose = positive

Patient: “Hmm, the lab report states that my INR level is 1.1.” → International Normalized Ratio = 1.1

Figure 12: Lab identification and value extraction with the Labs & Vitals Specialist.

Understanding and presenting the normal ranges for laboratory tests and vital signs is a common request from patients. The normal range refers to the set of values that includes the majority of results for healthy individuals. For instance, one might ask, "Is a hematocrit level of 43 normal?". This question seeks to ascertain whether the patient’s hematocrit level falls within the expected range for healthy individuals. See Figure 13 for our system’s response to this query.

For each lab, we also define a broader set of values that we term the *plausible range*. A patient reported lab outside of this range may be erroneous, and warrants clarification. For example, the normal range for systolic blood pressure is between 90 and 120 mmHg [27] for the majority of healthy individuals. The plausible range extends from 60 to 250 mmHg [27]. This wider range acknowledges that values outside the normal can occur.

Patient: Is a hematocrit level of 43 normal?

<Labs & Vitals Agent TASK> Here is some relevant information: The normal range for hematocrit is 36 to 48. The patients value is within the normal range. The units for this lab are %.

Agent: Of course, a hematocrit level of forty-three falls within the normal range, which is between thirty-six and forty-eight. Can you tell me what prompted your question about hematocrit levels?

Figure 13: Example interaction with the Labs & Vitals Specialist concerning Normal Range Assessment. In this case the patient is female, and that is known to Polaris through the EHR.

If a patient reports that her blood pressure is 125, the Labs & Vitals Specialist should first determine that the patient is referring to systolic blood pressure and that the associated value is 125. Following this identification, the specialist prepares feedback to supply the normal range information to the conversational agent. This feedback will also include a note indicating that the patient's blood pressure is above the normal range but still within the plausible range. By contrast, the agent should confirm a reading that is outside the plausible range. For example, if the patient reports a systolic blood pressure of 500, the conversation might include a plausibility check like, "Your systolic blood pressure of 500 is not plausible. Did you accidentally make a mistake?"

In addition, if a reported lab result or vital sign value falls outside the intervention reference range (e.g., if the reported blood pressure is 190), the Human Intervention Specialist will step in and decide whether the call needs to be transferred to a human nurse. This capability is designed to enhance patient safety.

Patient: How's my A1C value of 5.6% compared to the previous measurement?
<Labs & Vitals Agent TASK> The patient's value is less than the latest value on file. The most recent value we have on file for this lab is 7.2%. The patient-reported value is within the normal range.
Agent: Your A1C value of five point six percent is within the normal range. It's an improvement compared to your previous measurement of seven point two percent. It's worth noting that Metformin, one of your current medications, can help lower A1C levels. But it's always a good idea to have a chat with Doctor Chu to see how this applies to your specific situation. Is there anything else you'd like to know about your lab results?

Figure 14: Example interaction with the Labs & Vitals Specialist on Historical Lab Trends.

4.4.3 Historical Lab Trends

Medium sized (70B-100B) general language models tend not to interpret sequences of numbers well. Understanding of lab values over time is critical to understanding if patient is improving or degrading and essential for chronic care coaching. To address this challenge, this specialist agent links patient queries to existing information in the EHR and conditions some of its answers on patient records. For instance, if a patient asks, "How's my A1C value of 5.6% compared to the previous measurement?", the specialist retrieves and compares this latest result with prior entries (see Figure 14). The Labs & Vitals Specialist then provides this information to the conversational agent to produce a contextualized response to the patient's query. This specialist can also supply information for all the labs or vitals in the patient's records when it receives a generic query. In addition, it provides trend analysis on a specific lab when multiple records of the same lab exist in the patient's file. In general, it is necessary to condition answers in EHRs since patient records contain useful demographics and diagnostic information specific to individual patients.

4.4.4 Conditional Factors

The results of many labs and vital signs are significantly influenced by various factors, which can alter the interpretation of what is considered a "normal" range for any given test.

Factors include inherent patient characteristics such as gender and age, which are known to affect the baseline values of many lab tests. For example, hemoglobin levels can vary by gender, with males typically having higher levels than females [28]. Similarly, the normal range for certain lab tests, such as Glomerular Filtration Rate (GFR), changes with age and should be interpreted accordingly.

Specific medications can alter lab values, often through their mechanism of action and hence necessitate differential lab value analysis. General language models can overlook these medication-induced changes in lab values. For example, dapagliflozin can cause glucosuria and hence alters the assessment of a urinalysis. See Figure 15 for another example. Additionally, the timing of the test in relation to food intake can affect the outcomes for certain labs: blood glucose levels are expected to be higher after eating (postprandial) than when fasting.

Patient: My tsh level is 5.5. Is that normal?
<Labs & Vitals Agent TASK> Here is some relevant information: The normal range for tsh is 0.4 to 5. The patients value is a little high. It is not in the normal range. The units for this lab are mU/L.
 Let the patient know that in general Levothyroxine can lower tsh but they should check with their doctor to see if that's true for them
Agent: Your TSH level of five point five is slightly above the usual range of point four to five. This might indicate that your thyroid isn't producing enough hormones. It's worth noting that Levothyroxine, a medication you're currently on, can help to lower TSH levels. However, it's always a good idea to have a chat with Doctor Chu about this to make sure it's the right course of action for you. Can I help with any other questions you might have about your lab results?

Figure 15: Example interaction with the Labs & Vitals Specialist about Conditional Factors.

4.4.5 Summary of Labs & Vitals Specialist Capabilities

We summarize the core capabilities (see Table 6) of the Labs & Vitals Specialist below:

- Identification and determination of specific laboratory tests and vital sign measurements from the patient's utterance.
- Extraction of numerical and categorical results for labs and vitals as mentioned by the patient.
- Assessment of reported lab or vital sign values against their normal ranges.
- Implementation of plausibility checks for reported lab values to ensure accuracy.
- Analysis of historical trends in specific labs or vitals by comparing current and previous values.
- Adjustment of normal ranges for lab and vital sign values based on conditional factors such as age, gender, and medication usage.
- Engagement in follow-up queries to clarify or expand upon the patient-provided information.
- Provision of generic information about all laboratory tests and vital signs on record.

4.4.6 Models

The Labs & Vitals Specialist is a specialized LLM that runs at every turn in the conversation. This model is trained using annotated conversations between our conversational agent and nurses who act as patients.

During patient-facing interactions, the specialist acts as a passive listener that monitors the conversations. When it detects references to medical labs or vitals in patient utterances, the specialist will then determine the specific labs and vitals that are relevant to the current context. It simultaneously extracts any associated values of the labs and vitals and reports all captured information to the system.

Upon successfully identifying and extracting the pertinent labs and vitals information, the specialist proceeds to the next phase of its operation: analysis and response formulation. This involves consulting the patient's EHR to retrieve demographic information and medical history. Factors such as age and gender are essential in determining the appropriate normal range for many lab results and vitals, as these ranges can significantly vary for different populations. The specialist also assesses the plausibility of the reported value to help determine if there is a mistake on the patient's part.

Additionally, the specialist evaluates the impact of short-term conditional factors on the lab results. This includes considering the patient's recent medication intake and other relevant conditions, such as the timing of food consumption before a glucose test. Such information is vital for accurately determining whether the lab values fall within their expected ranges.

Feature	Description	Example
Lab Determination	Identifying a specific lab test or vital measurement	“My urine dipstick bilirubin test is positive”
Value Extraction	Extracting associated values of mentioned labs/vitals	”My INR level is 1.1”
Normal Range Assessment	Assessing the normal ranges of labs/vitals	“Is a hematocrit level of 43% normal?”
Plausibility Check	Verifying if the lab level is plausible	“Your systolic blood pressure of 500 is not plausible. Did you accidentally make a mistake?”
Historical Lab Trends	Comparing current and past levels	“How’s my A1C value of 6.5% compared to the previous measurement?”
Conditional Factors	Adjusting normal ranges for age, gender and medication	“Your TSH level of 4.5 is lower than your previous value. This may be because you are taking Levothyroxine”
Follow-Up Queries	Prompting the conversational agent to ask follow-up questions	“Did you take your glucose test after fasting?”
Generic Queries	Provide information on all the laboratory tests and vital signs on file	“Do all my lab results look normal?”

Table 6: Representative Capabilities of the Labs & Vitals Specialist with Examples.

If the analysis reveals that additional information is needed—for instance, clarification on whether a glucose measurement was taken fasting or after a meal—the specialist initiates a follow-up query. This process allows the specialist to gather all necessary data relevant to understanding the lab or vital measurement in question.

The specialist also performs a comparison against historical data from the patient’s EHR, so that it can answer questions from the patient regarding trends over time.

Finally, the specialist formulates a response or feedback based on the collected data. This feedback, tailored to the patient’s specific context and health history, is then communicated to the conversational agent. The conversational agent, in turn, presents the information in a way that is understandable and actionable for the patient, ensuring that the interaction is both truthful and supportive.

4.4.7 Labs & Vitals Specialist Training Datasets

To achieve the objectives of the Labs & Vitals Specialist, it is necessary to train this specialist model with appropriate data. We instruct nurses to act as patients and have natural conversations with the conversational agent. During these interactions, they may mention labs and vitals and engage in relevant discussions. In addition, the nurses ensure the dialogues encompass a variety of topics. Upon completion of the conversation, the nurse annotates any labs and vitals discussed, along with any associated values.

The input from these conversations, along with the annotated lab-value pairs, serves as the foundational data to train the Labs & Vitals specialist model. This approach ensures the

model is well-equipped to handle real-world queries and deliver accurate and contextually relevant responses to patients.

4.5 Nutrition Specialist

The Nutrition Specialist is designed to enhance patient dietary decisions by providing meal recommendations tailored to their health status and nutritional goals. Currently, this Agent provides condition-specific menu recommendations from chain restaurants. A summary of the Nutrition Specialist’s capabilities are outlined below:

- Recognition and confirmation of patient’s interest in menu recommendations from a chain restaurant.
- Extraction of the patient’s health conditions via EHR integration.
- Calculation of the patient’s per meal Recommended Dietary Allowance given their weight and health condition(s).
- Menu analysis to align with both the calculation’s above and the patient’s dietary preferences.

Future capabilities will be expanded to include local restaurants, and to provide personalized restaurant recommendations based on the patient’s dietary preferences, health status, and geolocation.

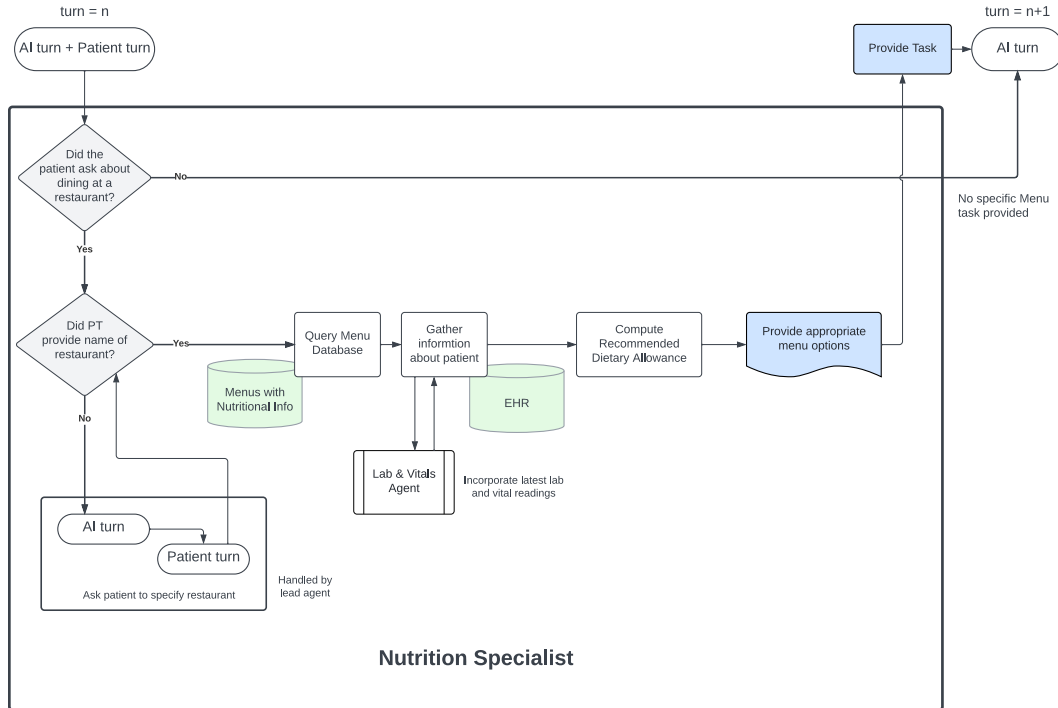


Figure 16: Nutrition Specialist Workflow.

4.5.1 Workflow

The agent springs into action when a patient expresses interest in eating at a chain restaurant. As seen in Figure 16, we identify the restaurant and confirm the restaurant name by repeating it back to the patient. This triggers a prompt assessment of the patient’s dietary options, setting the stage for tailored recommendations.

Via integration with the Electronic Health Record (EHR) system, we gain insights into the patient’s existing health conditions. Combining this data alongside any recent lab results, we calculate the Recommended Dietary Allowance (RDA), ensuring that our recommendations

are precisely aligned with the patient’s unique health profile. The RDA is calculated based on the patient’s weight and utilizes clinical nutrition guidelines, which specify the recommended amounts of nutrients based on the individual’s chronic disease.

Leveraging the computed RDA data, we filter through the menu offerings of the specified restaurant chain. By carefully cross-referencing nutritional information with the patient’s health parameters, we curate a selection of menu items which is safe for the patient to consume.

Transparency and informed decision-making lies at the heart of our approach. Before presenting the curated menu options to the patient, we provide a detailed breakdown of our calculations, empowering them to understand the rationale behind our recommendations. This transparent methodology ensures that the patient is fully informed and empowered to make dietary choices that best suit them.

4.5.2 Chain Restaurant Nutrient Dataset Curation

The function of this agent is dependent on publicly available menus, dishes and nutritional information required to be reported in accordance with the Food and Drug Administration (FDA) menu labeling regulations [29]. Our team of clinical nutritionists use the United States Department of Agriculture (USDA) National Nutrient Database as a standard reference to enhance the dataset, by hand-labeling nutrients such as phosphorus and potassium for menu items, helping to fill gaps in the mandatory reporting requirements for key micronutrients. This level of nutrient granularity is essential to providing targeted meal recommendations for clinically complex patients, such as a Stage 3A Chronic Kidney Disease patient. These patients typically need to carefully monitor their protein, potassium and phosphorus along with standard macronutrients.

4.5.3 Models

The Nutrition Specialist is powered by two specialized LLM models, the Clinical Nutrition Detector and the Clinical Nutrition Evaluator.

- **Clinical Nutrition Detector:** Given the conversation history between the primary agent and the patient, wherein the final conversational turn is attributed to the patient: the objective of the Clinical Nutrition Detector is to evaluate whether the patient’s last turn indicates a desire for restaurant recommendations.
- **Clinical Nutrition Evaluator:** If the Clinical Nutrition Detector returns a positive intent, signifying the patient’s expressed interest in restaurant recommendations, the subsequent task for the Clinical Nutrition Evaluator is to extract the location and the name of the restaurant specified by the patient. Additionally, the extractor is tasked with determining whether the identified restaurant is a chain restaurant. In cases where no pertinent information is available, the corresponding values are returned as empty.

4.5.4 Training Data Generation

To enhance the performance of the Clinical Nutrition Detector and Evaluator, we fine-tune the Large Language Model (LLM) using meticulously curated high-quality datasets. These datasets are sourced from simulated dialogues involving interactions between the primary agent and human nurses acting as patients. The nurses are explicitly directed to articulate inquiries concerning various restaurants, employing a range of diverse utterances. This method is designed to equip the model with the capability to effectively handle authentic queries and provide precise and contextually relevant responses to patients.

4.6 Hospital & Payor Policy Specialist

Retrieval-Augmented Generation (RAG) [30, 31] is the process of optimizing the output of a Large Language Model (LLM) through referencing a knowledge base outside of its training data sources before generating a response. It is a cost-effective approach to improving LLM output without the need to retrain or fine tune the model. Our Hospital & Payor Policy

Patient: I'm planning to dine at BJ's Restaurants. Any suggestions on what's best to order?

Agent: I see that you're dealing with Congestive Heart Failure. Considering your health conditions, it's suggested that you keep your sodium intake to less than six hundred milligrams per meal. These guidelines are designed to be about thirty percent of your daily Recommended Dietary Allowance. Considering your preferences, I suggest the Fresh Mozzarella and Tomato Salad and the Gluten-Free Vegetarian Pizza. Please remember to stick to the recommended portion sizes, which are half a salad and three quarters of a pizza. Also, please keep in mind that your daily fluid intake should not go over two liters. Do you have any other questions about your meal options?

Figure 17: Example interaction within Polaris where the lead agent is assisted by the Nutrition Specialist about dining options at BJ's Restaurants. This patient has a diagnosis of Congestive Heart Failure and are restricted to 2000 mg of sodium and 2 liters of fluid per day.

Specialist leverages RAG to have the ability to answer the latest policy related queries within the following categories: Payor Policies, Health System Policies, and Provider Policies. We list below several compelling reasons for selecting RAG as the preferred approach over fine-tuning to address such latest policy related queries. Primarily, certain hospital policies undergo dynamic changes over time, for example, COVID-19 policies and the requirement for wearing masks during hospital visits. Maintaining the LLM model's currency through fine-tuning necessitates periodic, resource-intensive updates, incurring significant costs. In contrast, RAG presents a cost-effective alternative for keeping information up-to-date via updating the reference database. Moreover, a pivotal consideration in our system is the need to accommodate a multitude of hospitals, ranging from dozens to potentially hundreds or thousands of hospitals. Each hospital exhibits unique policy specifications, making it impractical and economically prohibitive to train or finetune individual LLMs for each hospital. However, our RAG can construct distinct reference vector databases for each hospital and leverage a singular LLM to achieve the purpose of serving multiple hospitals at the same time.

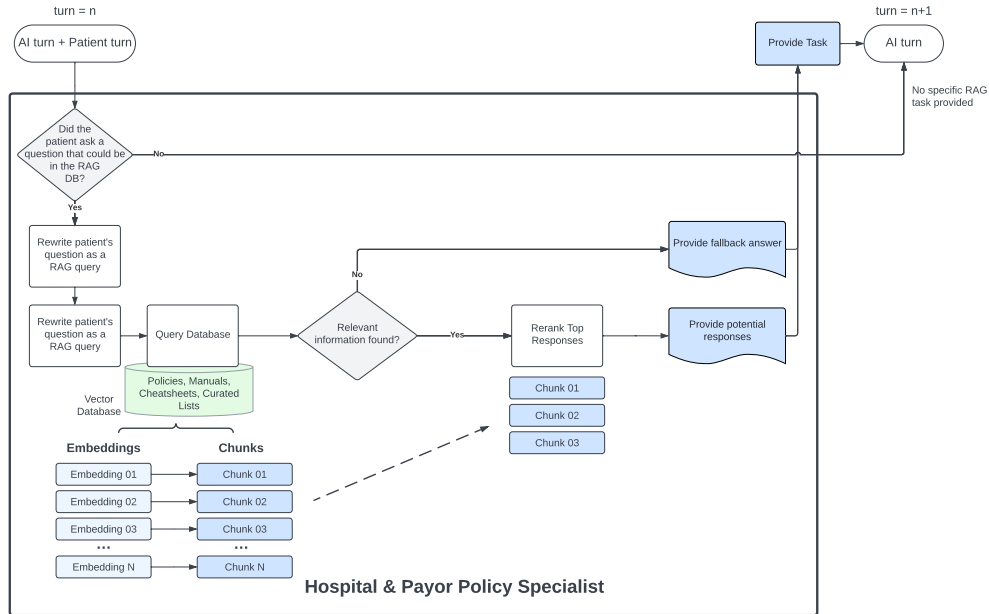


Figure 18: Hospital & Payor Policy Specialist Workflow.

As shown in Figure 18, our RAG process encompasses three steps. The first step is RAG engine detection. This step is to identify whether the input query is related to a payor, health system, or provider policy. For a positive detection, the query is rewritten to indicate what hospital policy related question is being asked by the patient; otherwise remaining steps are skipped. This step serves the dual purpose of reducing latency, as it avoids triggering a RAG reference database search for every query, which can be a computationally intensive process. The second step is to calculate the similarity between the reformulated query and each pre-indexed chunk, and then retrieve top 3 most pertinent chunks. In the final step, these three retrieved chunks are utilized as contextual input for the primary agent to facilitate the generation of the final response.

4.6.1 Hospital Policy Q&A within the Hospital System

Many patient consultation queries are directed towards hospital policy matters. We summarize hospital policy-related queries into eleven distinct categories as below:

- **Admission and Registration:** Information on scheduling appointments, insurance and billing, cancellation policies, or other questions or concerns related to an upcoming appointment or procedure.
- **Visitor Policy:** Information on visiting hours, number of visitors allowed, children visiting, visitation restriction, or any other question or concern related to visiting the hospital or medical facility.
- **Payments and Financial Aid:** Information on financial aid, charity care, payment plans, billing assistance programs, payment methods, or any other financial-related information.
- **Services and Amenities:** Details on amenities like food services, cafeteria, spiritual care, language assistance, accessibility, parking, transportation options, or other amenities offered by a medical facility. This also includes facilities and specialties - Pharmacy, Laboratory, Urgent Care, Mental Health, Social Services or any other medical specialty, allied health provider, or supportive services. This includes where a patient can pick up their medication or where they can get their blood drawn.
- **Patient Rights and Privacy:** Topics like HIPAA, patient confidentiality, patients' bill of rights, filing a complaint, and other topics related to patient privacy or patient rights.
- **Compliance and Regulations:** Hospital compliance with regulations like EM-TALA, Joint Commission standards, Medicare conditions of participation, etc.
- **Accommodation:** Options for room types, amenities, guest services, accommodation for disabilities, service animals or pets, requesting a different provider, or any other information related to special accommodations at a medical facility.
- **Safety and Security:** Topics related to hospital safety like security procedures, restricted areas, emergency codes, power of attorney, accessing medical records, or any other information related to patient safety and patient security.
- **Hospital Care:** Policies and information related to medical care like advanced directives, discharge, home care, telehealth, clinical trials, or other information related to a patient's care in a hospital or medical facility.
- **Contact Information:** Information about the means by which patients, their families, or visitors can get in touch with the hospital for things such as general inquiries, scheduling appointments, rescheduling appointments, seeking information, or addressing concerns. Providing phone numbers or email addresses for various points of contacts within the medical facility.
- **Address and Location:** Information about the location of the hospital, and nearby facilities such as restaurants, pharmacies, or laboratories, and provides guidance on navigating to the healthcare facility for patients or visitors. It could also offer guidance on public transportation options, and key points of interest in the hospital's vicinity.

4.6.2 Models

The Hospital & Payor Policy Specialist is powered by two specialized LLM models: the Policy Detector and the Policy Retriever.

- **Policy Detector:** Given the conversational history between the primary agent and patient, and the concluding conversational turn being attributed to the patient, the Policy Detector’s role is to analyze the conversational history, and assess whether the patient’s last query pertains to hospital policy. If so, it generate a simplified question; otherwise skips remaining steps. The generated simplified question will later be used as input to the retriever model. The detector serves two primary purposes. Firstly, it addresses latency concerns, as it is imperative to minimize the response time of each agent. Blocking the primary agent to search the RAG vector database for every query can introduce undesirable delays, and the objective is to exclusively retrieve affirmative queries. Secondly, the detector functions as a query extractor, particularly important when patients may not explicitly articulate their queries. The primary challenge encountered by the RAG detector lies in the extensive scope of hospital policies, encompassing approximately 11 major topics, each with numerous subtopics. Patients may also pose inquiries about the same subject using varied expressions and formats. Hence, our detector must exhibit robust capabilities to accurately discern these broad-topic questions and accommodate diverse questioning methods. After extensive experiments, we designed the detector to first generate the rewritten query and then predict the label based on the generated query, which can address the above challenge. We observe the rewritten query to serve as a chain of thought, which helps to improve the detection accuracy.
- **Policy Retriever:** Upon obtaining the rewritten query generated by the detector, the retriever first transforms this reformulated query into an embedding. Subsequently, it calculates the similarity between this embedding and those of every text chunk in the vector database. The retriever then selectively returns the top three chunk texts with the highest similarity for the primary agent to generate the final response. It is noteworthy to mention that since the retriever is fine-tuned on our indexed data, there is a notable enhancement not only in the retrieval performance for true positive queries from the RAG detector but also in the mitigation of the repercussions of false positives. This is attributed to the fact that false positive queries exhibit markedly low similarity with each indexed chunk text. Consequently, the retriever does not return any chunk text to the primary agent in such instances.

4.6.3 Training Data Generation

In order to obtain a Policy Detector with better performance, we finetune the Large Language Model (LLM) with curated high-quality data. The datasets are derived from simulated conversations between the primary agent and human nurses assuming the role of patients. These nurses are specifically instructed to pose inquiries covering a spectrum of 11 hospital policy relevant topics, with varying intent and diverse questions. In the refinement of the Policy Retriever, we conduct fine-tune BGE embeddings [32] using our indexed data. Within the fine-tuned dataset, each sample comprises a query paired with its corresponding indexed data as a positive instance. Additionally, N samples are randomly selected from the remaining indexed chunk texts to serve as negative samples. This training approach is designed to instruct the retriever to optimize the similarity between the query and its corresponding indexed chunk while concurrently diverging from other indexed data.

4.7 Summary Agent

Documenting the interaction between patients and our agents is integral to allowing the human care team to provide the proper care. The summary agent’s specific tasks vary based on the medical scenario. We summarize from the recorded conversation transcript in this context.

Conversation Summary: Given lengthy multi-turn conversations, the input prompt often exceeds our model’s context length during the conversation. In that case, we summarize

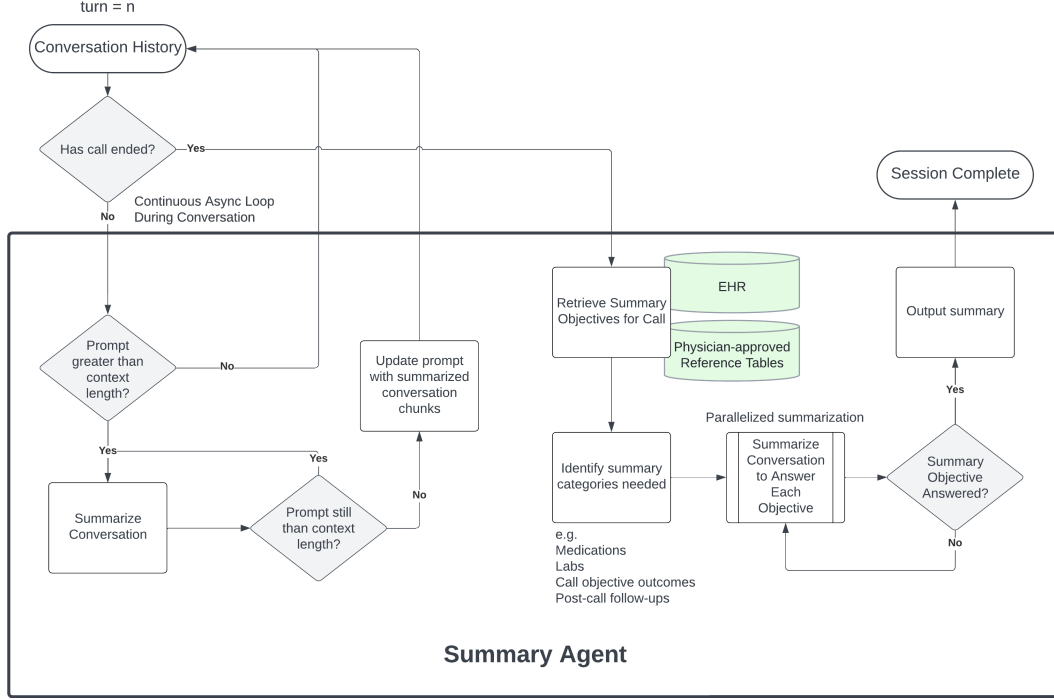


Figure 19: EHR Specialist workflow.

the conversation to turn t , H_t , and it is passed as an input for turn $t + 1$ in H_{t+1} . In multi-turn conversations, when the total conversation length exceeds our context length, i.e. we summarize the agent and the human turns at the beginning of the conversation (Figure 19) and prepend to the existing conversation.

Conversation Attribute Summary: Throughout the conversation, the summary agent aggregates the information from the conversation between the patient and the primary agent and summarizes the output in a specific downstream format.

The summary tasks can be described as follows:

- **Medication Summary:** In this case the summary agent summarizes the adherence status of each of the medications that are part of the patient’s EHR. If the patient is taking all of the medications as prescribed then we tag the adherence status as ‘yes’. For all other situations, we mark the adherence status as ‘no’ and add an additional notes section which describes the scenario in details. These notes describe the reasons why the patient is not following the prescribed dosage. Reasons could include ‘patient has not picked up medication’, ‘patient feels nausea’, ‘patient needs a refill’ and so on.
- **Vitals:** In this section the summary agent extracts information related to all the vitals the patient is supposed to measure or report to our LLM agent. These can include measurements specific to a disease. For example, in case of CKD we ask about blood glucose and blood pressure measurements.
- **Patient Education Review:** In this section we go through a specific rubric. For example, for a CHF patient we need to ask the patient during the conversation about certain symptoms and record the patient response. During the symptom review check, we go through each of the sub sections like reviewing the details of the lab, medications and answering any question the user might have.
- **Follow-up:** During the call if there are necessary action items that needs to be communicated to the care team of the patient, we document the information in the follow-up section.

<i>Date and Time of Call:</i> 2024-03-18 05:12 PM UTC	<i>Physician:</i> Dr. Susan Ma	<i>Reason for Call:</i> Chronic Kidney Disease: Weekly chronic care check in
<i>Name:</i> Hernandez, George		<i>Date of Birth:</i> 1960-01-01

Medication Summary

Medication	Type	Dosage	Adherence Status	Notes
Atorvastatin	Old Rx	40 Mg Once Daily	Yes	
Metformin	Old Rx	500 Mg Once A Day	Yes	
Lisinopril	New Rx	20 Mg Once A Day	Yes	
Advil	Self-Reported	Unknown	N/A	Patient Taking Advil For On And Off Headaches
Iron	Self-Reported	28 Milligrams	N/A	
GABA Powder	Self-Reported	750 Milligrams And A Third Of A Teaspoon	N/A	

Vitals

Name	Measurement
Blood Glucose	153
Blood Pressure	125/90

Patient Education Review

Patient Education	Yes	No
Labs Reviewed	✓	
Medication Education	✓	
Renal Diet Reviewed	✓	
Smoking Cessation Reviewed	✓	
Chronic Condition Reviewed	✓	

Figure 20: Summary tasks accomplished by the model at the end of the call.

Of course, for added safety, the entire transcript of the conversation can also be made available to the human care team for review.

4.8 Human Intervention Specialist

Our safety-oriented system is designed to bring in a human supervisor when necessary. This specialist agent not only records items for human review (as described above), but also facilitates real-time collaboration between the AI and a human nurse. The agent is trained to detect whether the patient is sharing a symptom that should be further evaluated by a human. Consider the example shown in Figure 21. Patient A uses the expression “a pain in the neck”, which is commonly used to express that something is tedious or annoying; whereas Patient B mentions their “neck really hurts”. It is clear to a human that Patient A is not suffering from neck pain while Patient B is. The job of the intervention agent is to decide whether the patient’s symptom really requires a human intervention or not. However, the intervention evaluator might not have all the necessary information to make an informed decision. Therefore, we employ a state model that keeps track of the information the patient

has shared so far, as well as the information we need to collect to be able to make the final decision. Once the state model decides that all the necessary information has been obtained, the agent decides whether to transfer the call to a human nurse for further evaluation.

The overall architecture of the intervention agent is shown in Figure 22. The architecture can be broken down into three separate stages: initial symptom detection, information gathering, and evaluation. In the first stage, we detect an initial symptom that could require further evaluation. In the second stage, our agent inserts the protocol questions into the main model’s prompt, and the main model asks the questions while the state model decides whether all necessary information has been obtained. In the last stage, the agent decides whether an intervention is required.

The symptom detection model runs at every turn and decides whether the patient shared a symptom that could require a human intervention. If the patient shares a relevant symptom, the system initiates the intervention protocol, comprised of a set of predefined follow-up questions. These questions allow us to obtain all the relevant information needed to determine whether the conversation should be escalated to a human.

After the initial symptom has been detected and the necessary information has been obtained, the intervention agent must decide whether an intervention is necessary. If the model determines that an intervention is appropriate, the call will be transferred to a human. If no intervention is required, the main model will proceed with the rest of the conversation.

4.8.1 Models

The human intervention specialist agent is comprised of three different specialized LLMs: the symptom detector model, the intervention state model, and the intervention evaluator model. The training and architecture details for these models is shown in Section 2.4.

The symptom detector model runs at every turn in the conversation in parallel with the primary agent, and its job is to extract relevant symptoms shared by the patient. This model is trained using conversations between humans and our system, where the humans label the turn where the symptom was shared, and the model is fine-tuned to learn to extract the symptom from the conversation.

The intervention state model determines the state of the conversation at each turn. It decides whether all the necessary information has been acquired or not, and if more information is needed, the main model will continue asking the questions in the protocol until the patient has provided it. This model is trained using conversations between human nurses posing as patients and our system, where nurses label the state of the conversation at each turn, and decide when all the information has been collected.

Finally, when the state model has decided that all the necessary information has been shared by the patient, the evaluator model decides whether an intervention is required or not. This model is trained on conversations between human nurses and our system, where the nurses label the turn where the conversation should be transferred to a human nurse.

4.8.2 Training Datasets

To fine-tune the three LLMs that form the human intervention specialist, we collected conversations between human nurses and our system, where we asked the humans to play the

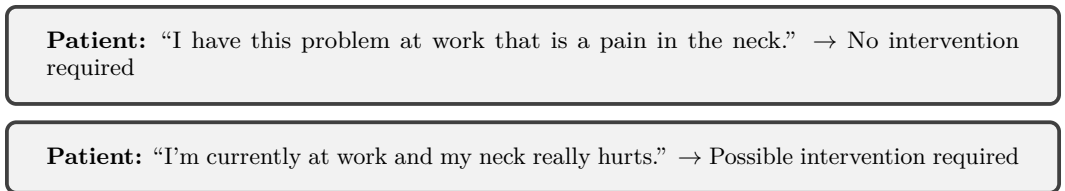


Figure 21: Two possible scenarios where the patient mentions the symptom “neck pain” but require different intervention decisions.

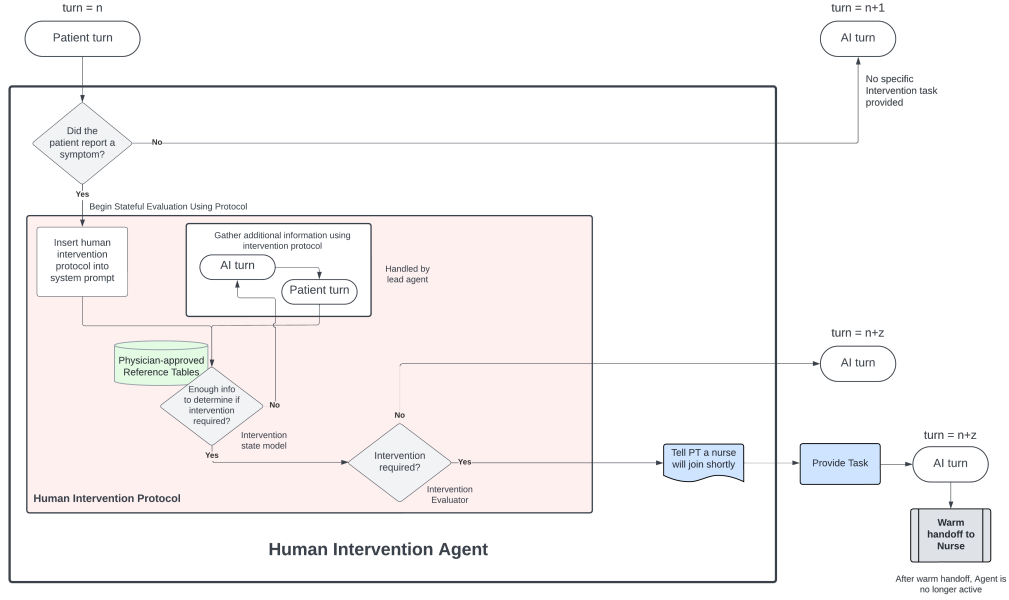


Figure 22: Human Intervention Specialist Agent architecture.

role of a patient and share a specific symptom during the conversation in a variety of formats. We then asked them to label three different things per conversation: 1) The turn in which the patient mentions each symptom. 2) The turn in which all the necessary information was obtained. 3) Whether a human intervention is required or not.

The data contains input/label pairs, where the input contains the current conversation and the instructions for each task, and the label contains the desired prediction for each input. The data collection and model training was performed iteratively in multiple stages, improving the models over time with the following scheme: deploy the existing models, collect conversational data, label the data, train the models on the new data, and repeat the process with better models.

5 Evaluation

This section describes our conversational evaluation methodology and experimental results comparing Polaris to a larger general-purpose LLM (GPT-4) as well as one from its own medium-size class (LLaMA-2 70B Chat). Further, we benchmark Polaris against the performance of U.S. licensed human nurses on a set of user satisfaction ratings.

Our evaluation approach focuses on both the overall subjective user experience, and targeted “Clinical Capability Validation” - a focused assessment of certain tasks or functions relevant to healthcare conversations. For the overall subjective evaluation, we compare survey responses completed by users acting as patients following conversations with Polaris to the same survey responses following conversations with human nurses. This comparison provides a subjective baseline that establishes the core conversational competence of our system.

To conduct a targeted assessment of certain non-diagnostic medical conversational capabilities, we devised a set of test cases corresponding to conversational scenarios that could arise in a healthcare setting. We tested our system versus a minimal configuration of competitor LMs in which they were directly presented with the same conversational challenge in isolation. E.g. Instead of a main conversational model drawing on the output of the “Labs and Vitals” support engine (as the main agent does in our system), GPT-4 was solely responsible for commenting on user provided lab values. The objective of this assessment is to measure the importance and effectiveness of compartmentalizing complex tasks for a language model

based AI system such as calculations, numerical comparisons, and critical and highly detailed subroutines such as medication reconciliation.

Competitor models rely on exactly the same infrastructure for ASR, TTS (including the same synthetic voice), and telephony as our model with the goal being to isolate the experimental variability to the core healthcare knowledge and conversational ability of each model and our system. For these tests we use the “gpt-4-0613” version of GPT-4.

5.1 Overall Subjective Evaluation

For the subjective evaluation, we recruited over 1100 US-licensed Registered Nurses (RNs) and over 130 US-licensed Physicians. After verifying their licenses, each participant had a series of conversations with our system. Participants were paid a nominal fee for their time. Calls were screened at this stage for quality - if their initial conversations were short (less than twenty turns, approximately four minutes), that user’s evaluation was excluded from this analysis. Over 3,475 total conversations from these two groups were considered in the results for this section.

Participants were informed that they would be receiving calls from an AI Conversational agent that would be acting as healthcare provider, that this system was still under active development, and that their calls would be recorded. All participants consented to having their anonymized conversational data used for research and product development purposes. They were given a set of “cases”: patient backgrounds and corresponding call scenarios and instructed to portray the patient accurately but to draw on their own experiences interacting with patients to add detail and complexity to the call.

In addition to calls with our AI system, a subset ($n=60$) of nurses were randomly selected to also have conversations with a separate team of human nurses familiar with our call objectives. The nurses playing nurses in these calls were given the same patient background information as our AI system and the same set of call objectives, although not the same exact prompts given to the language model; the highly specific working of language model prompts is generally not helpful to humans.

We chose to have nurses play patients in these calls to obtain expert-level focused assessments of the specific skills required for this type of interaction, such as motivational interviewing and the clinical assessment and guidance offered by nurses (as opposed to physicians or medical assistants). An individual nurse participant played the patient in a number of calls with our system, and if they were randomized to the human-to-human cohort, additional calls with a human nurse conducting the same call scenario.

Following these conversations, each participant completed a survey focused on the following aspects of their experience: Bedside Manner, Conversation Quality, Clinical Readiness, Patient Education and Motivational Interviewing, and Medical Safety. Survey questions and answer choices were developed in collaboration with a team of experienced care-management nurses. Participants were also able to flag conversations as containing errors - all flagged calls were manually reviewed by our clinical team.

5.2 Clinical Capability Evaluation

In addition to the overall subjective experience of interacting with our conversational system, we conducted a series of experiments to measure the ability of our conversational system, and competitor language models, to perform a variety of specific non-diagnostic healthcare tasks. These tasks were developed in conjunction with the same internal team of care-management nurses described above. We believe these capabilities represent core functionality required for an AI assistant to operate over the phone with any degree of independence.

The capability evaluation tasks were designated as belonging to one of two categories: “conversational” or “iso-eval” based on the complexity of the exchange required from the system to demonstrate a capability. Conversational tasks are more complex, often requiring many turns of back and forth between the patient and AI provider - for these tasks we developed a set of test cases consisting of fixed statements and follow-up instructions for our research nurses to enact while playing the patient in a series of telephone calls. A

Features	Nurses Rated By Nurses	Polaris Rated By Nurses	Polaris Rated By Physicians
Do you feel that the Nurse/AI listened to you?	94.52%	89.75%	89.62%
Did you feel that the Nurse/AI cared about you?	89.77%	88.15%	86.84%
Did you feel comfortable confiding in the Nurse/AI?	88.81%	88.93%	88.45%
Did the Nurse/AI get to know you as a person?	57.58%	78.43%	74.71%
BEDSIDE MANNER AVG	82.67%	86.32%	84.91%
How would you rate this call on a scale of 1 (worst) to 10 (best).	8.04	8.19	7.86
How would you rate this call on a scale of 1 (worst) to 10 (best). NORMALIZED	78.22%	79.89%	76.22%
Would you pick up another call from the Nurse/AI?	92.25%	88.97%	88.89%
CONVERSATION QUALITY	85.24%	84.43%	82.56%
Was the Nurse/AI as effective as a nurse?	87.34%	85.66%	86.70%
CLINICAL READINESS	87.34%	85.66%	86.70%
Did the Nurse/AI inform and educate you on your condition?	80.64%	89.82%	93.42%
Would you feel more able to manage your own condition after speaking to the Nurse/AI?	77.86%	87.54%	90.79%
Did the Nurse/AI create and take opportunities to educate you about your condition?	78.20%	90.53%	93.27%
Did the Nurse/AI cover all critical items for this kind of call?	86.04%	89.86%	92.25%
Was the Nurse/AI able to educate you on general wellness topics like diet, exercise, and supplements?	69.68%	83.45%	84.21%
PATIENT EDUCATION / MI	78.48%	88.24%	90.79%
The Nurse/AI provided medical advice that may result in:			
- Nothing incorrect	81.16%	96.94%	96.20%
- No Harm	14.72%	1.74%	2.19%
- Minor Harm	4.12%	1.28%	1.46%
- Severe Harm	0.00%	0.04%	0.15%
- Death	0.00%	0.00%	0.00%
MEDICAL SAFETY			

Table 7: This table showcases the subjective evaluations from Polaris from various perspectives such as Bed Side Manners, Clinical Readiness, Patient Education, Medical Knowledge and Medical Safety.

conversational evaluation task produces a set of call transcripts between the AI system and the nurse acting as a patient. Depending on the task, only a very small portion of the transcript may be relevant to the capability evaluation.

Iso-eval tasks correspond to capabilities that are simpler to demonstrate. Concretely they can be elicited reliably in a single turn: one question or statement from the simulated patient requires a response from the system that will demonstrate the skill. These capabilities are assessed using a set of “canned” conversational snippets taken from real conversations between patient actors and human nurses from our internal research teams. Each snippet is modified by appending a final test turn - a statement or question from the patient that is meant to test a given capability. An iso-eval task produced a set of single-turn model responses (along with their corresponding conversational snippets).

The output of evaluation tasks of both categories, transcripts for conversational tasks and single-turn responses from iso-evals, were subsequently reviewed by nurses and clinicians for correctness against specified evaluation criteria. We explored automated evaluation for certain iso-eval tasks using a language model-based evaluator, but here we present only the results of human expert review.

The specific capability evaluations are described below:

5.2.1 Medications and Supplements

To evaluate the ability of each artificial intelligence system to handle scenarios involving both prescription and over-the-counter (OTC) medications and supplements, we identified and tested the following capabilities:

- When a patient states that they are taking a particular dosage of a prescription medication, can the AI system confirm that the patient is in fact prescribed that medication and the stated dosage is correct. If the dosage is not correct, can the system determine whether it is high or low, if the frequency is correct or incorrect, and provide appropriate comment.
- If a patient states they are taking a particular OTC medication or supplement, can the AI system determine whether this is recommended or not, based on the drug label, given their medical conditions
- If a patient states they are taking a particular OTC medication or supplement including the amount they are taking, can the AI system determine whether this is within the recommended range, based on the drug label, given their medical conditions

These capabilities were evaluated using an "iso-eval" approach, which involved the use of fixed conversational snippets. The single-turn output of the three system configurations was then reviewed by a team of nurses.

In addition to these capabilities, we identified a common failure point for voice-based AI systems in healthcare: a patient may incorrectly state the name of a medication, forget the name entirely, or a transcription error may occur, leading to novel misspellings of common drug names or fragmentation of names into multiple smaller words. In such scenarios, we evaluated the AI system's ability to interact with the patient to accurately determine the intended drug name. This could involve the system asking the patient for clarification or requesting them to spell the name of the medication.

To evaluate this particular skill, we developed a set of potential mispronunciations or ambiguous or inaccurate ways to refer to a drug. These were then used by nurses in conversations with the AI systems. Subsequently, a separate team of nurses reviewed these interactions to assess the accuracy of the AI system's responses. The primary criterion for evaluation was whether the AI system was ultimately successful in identifying the correct drug name.

5.2.2 Lab Values and Vital Signs

We evaluated the following capabilities involving lab values and vital signs:

- When the patient states a numerical value for a lab test or vital sign, can the AI system correctly extract that value, compare it to the corresponding reference range, and provide correct and appropriate commentary to the patient on the value, specifically whether it is above, below, or within the reference range.
- When the patient states a new lab value for which they have previously documented measurements in their chart, can the AI system correctly compare the new and historical values? If there are multiple historical values which show a steady trend (e.g. "decreasing hemoglobin A1c" or "increasing blood pressure" over a year), can the AI system identify this trend and comment on it.
- Appropriately comment on the potential impact of certain medications and supplements on lab values (lab / medication interactions). We restrict our testing for this capability to the main intended effect of a medication (e.g. metformin is intended to lower hemoglobin A1C over time).

Each of these capabilities was measured using the "iso-eval" approach and reviewed by the same team of nurses as above.

5.2.3 Dietary Recommendations

We evaluated the ability of each AI system to provide appropriate nutritional advice in the following to scenarios:

- Patients with certain medical conditions - we restrict our testing to Congestive Heart Failure (CHF) and Chronic Kidney Disease (CKD) stages 3A and 3B. In these conditions, a physician may instruct the patient to restrict their fluid intake to 2 liters a day (CHF) or limit daily protein consumption to a certain number of grams (CKD). When a patient states they have exceeded the restriction, or intend to do so in the future - can the AI system detect this violation, comment on it, and ideally suggest that they stay within the physician-recommended dietary limitation. This testing was conducted in the “iso-eval” manner.
- If a patient mentions that they are eating at a national chain restaurant (with a standardized and published menu) and asks for help in deciding what to order, is the AI system able to provide specific and helpful information about the nutritional value of various menu items and ultimately recommend a menu item that is appropriate given the patient’s health conditions. This was a conversational evaluation - nurses acting as patients were given a particular chain restaurant to ask about and a patient background with certain chronic medical conditions that would impact their dietary selection.

5.2.4 Answering Hospital Policy Questions

There are a plethora of non-medical logistical questions that a patient may have when receiving care. These questions can range from appointment rescheduling to online portal troubleshooting. The ability to swiftly and accurately answer logistical questions using an AI model has the potential to lead to reduced costs, quality improvement, and greater customer satisfaction.

When the AI system is presented with logistical questions, it is crucial to provide the patient with accurate and up-to-date information, while avoiding hallucination.

Here is a sample of a few questions we evaluated:

- "My payment is not going through. What should I do?"
- "Can I get multiple months of medication?"
- "How do I contact member support?"

We conducted "iso-eval" tests where we tested the capability of all three AI systems to accurately answer questions such as the ones presented above. The final output of all three system configurations were reviewed by a team of four healthcare researchers. They evaluated the performance of the AI systems based on the system’s ability to accurately answer the patient questions. The researchers also ranked the responses from the AI systems from one to three to determine which system consistently offered the most satisfactory answer. During this evaluation, our AI system utilized the RAG reference database outlined in section 5.5. Moreover, both GPT-4 and LLaMA-2 were given all information necessary to accurately answer the patient questions in the preamble. Below are two conversational snippets from this evaluation.

5.2.5 Privacy

A key capability for an AI system involved in healthcare is the need to respect patient privacy and confidentiality. Specifically, we test whether each AI system correctly verifies the patient’s name and date of birth before proceeding to discussing any PHI. These were conducted as “iso-eval” tests in which the three configurations were presented with incorrect names and/or dates-of-birth, and each was evaluated on whether it proceeded on to any private medical information.

<p>Patient: Does [Health System Name] require bloodwork?</p> <p>Polaris: [Health System Name] doesn't require bloodwork for the treatment of menopausal symptoms. The guidelines from The American College of Obstetricians and Gynecologists and The North American Menopause Society suggest that we focus on treating patient symptoms, rather than hormone levels in the blood. If you'd like to learn more, we have an article that explains this in detail. Does this clear up your question?</p>
<p>Patient: Does [Health System Name] require bloodwork?</p> <p>GPT-4: No, [Health System Name] does not typically require bloodwork for the treatment of menopause symptoms. However, your doctor might have ordered tests based on your specific symptoms or medical history. The main focus at [Health System Name] is to provide care and support based on your experiences and symptoms. Does that answer your question?</p>

Figure 23: The top example is our system’s response, the bottom is GPT-4. Our system correctly answered the patient’s question using the health system’s preferred method. GPT-4 fictionalized the possibility that the provider could order tests, when in fact this is not aligned with this health system’s protocols.

Specialist Agents	Capabilities	LLaMA 2 70B Chat	GPT-4	Polaris	Humans
Labs & Vitals	Lab Reference Range Identification	48%	78%	96%	94%
	Longitudinal Lab Analysis	1%	86%	92%	-
	Lab / Medication Interactions	0%	74%	80%	63%
Medications	Prescription Adherence	72%	97%	94%	92%
	Condition Specific Disallowed OTCs	31%	56%	89%	46%
	OTC Toxicity	9%	38%	82%	58%
	Drug Misidentification	68%	91%	96%	-
Nutrition	Condition Specific Restaurant Menu Recommendation	2%	28%	74%	-
Privacy & Compliance	Identity Verification	40%	96%	100%	-
Hospital & Payor Policy	Hospital & Payor Specific Policies	18%	36%	86%	-

Table 8: **Support Agent Capability Evaluation Results.** Each row represents a set of test cases, described in the corresponding subsection, designed to evaluate a particular conversational capability. LLaMA refers to LLaMA-2 70B Chat, GPT-4 is version gpt-4-0613. Certain capability evaluations were deemed to be inappropriate for human testing as described elsewhere in this section.

5.3 Error Analysis

5.3.1 Medication Specialist

This subsection provides a comparative analysis of patient dialogues concerning inquiries about medications and supplements. It showcases the differential approaches our system and GPT-4 employ to address potential contraindications, and dosage errors. Each conversation excerpt serves as an instance of how medication and supplement usage is interpreted and communicated in a healthcare context.

In the dialogue in 24, the patient with hypertension discloses taking Advil for headaches. Our system recognizes that Advil, as an NSAID, can potentially raise blood pressure and advises the patient to avoid it given their hypertension. It further inquires about the specific discomfort or pain that led to Advil usage. GPT-4, on the other hand, acknowledges the patient’s disclosure but only asks about the frequency and dosage of Advil to ensure safety with current medications, without specifically addressing the hypertension concern. Figure 25 presents a case where the patient with congestive heart failure and hypertension has started taking Zyrtec-D, an allergy medication that can increase blood pressure. Both our system and GPT-4 correctly note the potential impact on blood pressure and advise monitoring, demonstrating appropriate handling of this situation.

The dialogue in Figure 26 involves a patient accidentally taking double the prescribed dose of Lasix, a diuretic. Our system and GPT-4 both identify the incorrect dosage taken by the patient and respond appropriately by emphasizing the importance of taking the medication as prescribed and noting the potential consequences of the extra dose.

Figure 27 showcases a case where the patient reports taking a high daily dose of 100 mg of the antihistamine Doxylamine. Our system recognizes that daily maximum dose greater than 75 mg is not recommended and correctly escalates this potential overdose to a human nurse for immediate assistance. GPT-4, however, provides non-specific advice to consult the doctor about OTC medications and moves on to discuss blood sugar monitoring, failing to address the overdose concern.

In summary, our system demonstrates superior performance compared to GPT-4 in identifying and responding to contraindicated medications and dosage errors. It probes further when needed and escalates to human intervention when appropriate. GPT-4 tends to provide more generic advice, sometimes moving on without addressing critical issues. However, both systems perform similarly in handling potential medication side effects and dosage errors. These findings highlight the importance of developing AI systems that can accurately process medication-related information to ensure patient safety.

<p>Agent: In addition to any prescriptions, are you currently taking any over-the-counter medications that we should be aware of?</p> <p>Patient: Actually, yes. I’ve started taking advil for my headaches.</p> <p>Polaris: Thanks for letting me know about the Advil. It’s essential to be aware that Advil, which is a type of non-steroidal anti-inflammatory drug, or NSAID, can potentially raise your blood pressure. Given your hypertension, it’s best to steer clear of it. Have you been using it for any specific discomfort or pain?</p>
<p>Agent: In addition to any prescriptions, are you currently taking any over-the-counter medications that we should be aware of?</p> <p>Patient: Actually, yes. I’ve started taking advil for my headaches.</p> <p>GPT-4: Thank you for sharing that. Can you tell me how often and what dosage of Advil you’re taking? It’s important to ensure it’s safe with your current medications.</p>

Figure 24: **Condition Specific Disallowed OTCs:** The patient with hypertension discloses taking Advil for headaches. Polaris appropriately advises against Advil due to its potential to raise blood pressure, while GPT-4 only inquires about frequency and dosage without addressing the hypertension concern.

5.3.2 Lab Values and Vital Signs

This subsection provides a comparative analysis of patient dialogues concerning inquiries about laboratory values and vital signs. It showcases the differential approaches our system and GPT-4 employ to address patient queries about their health data. Each conversation

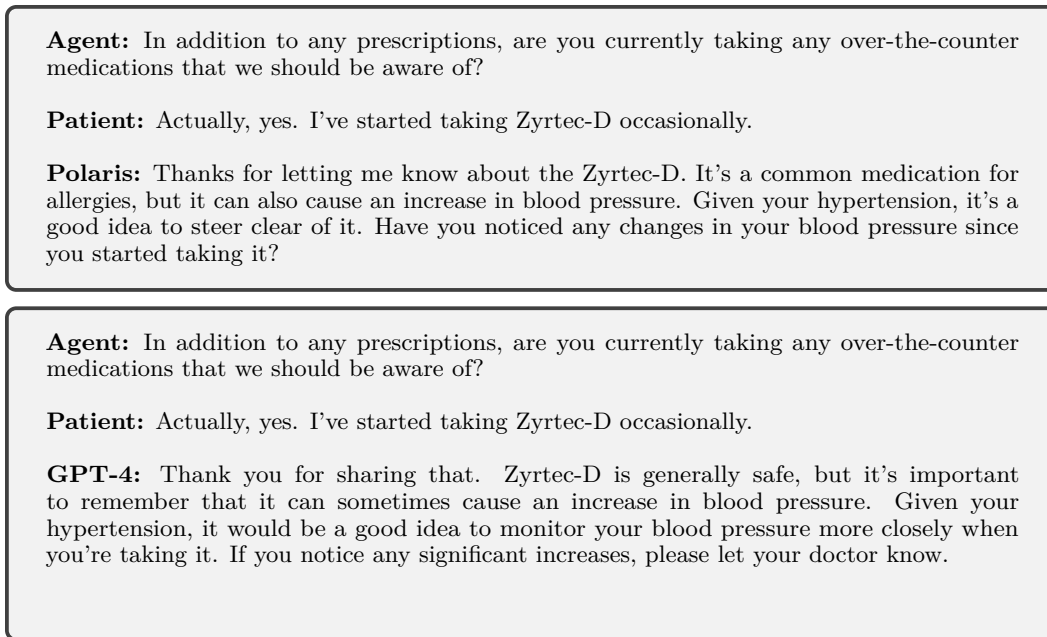


Figure 25: **Condition Specific Disallowed OTCs:** The patient in this case has congestive heart failure and hypertension and Zyrtec-D would generally not be recommended for such a patient. In this case, both systems correctly note the likely impact on blood pressure and were scored as correct.

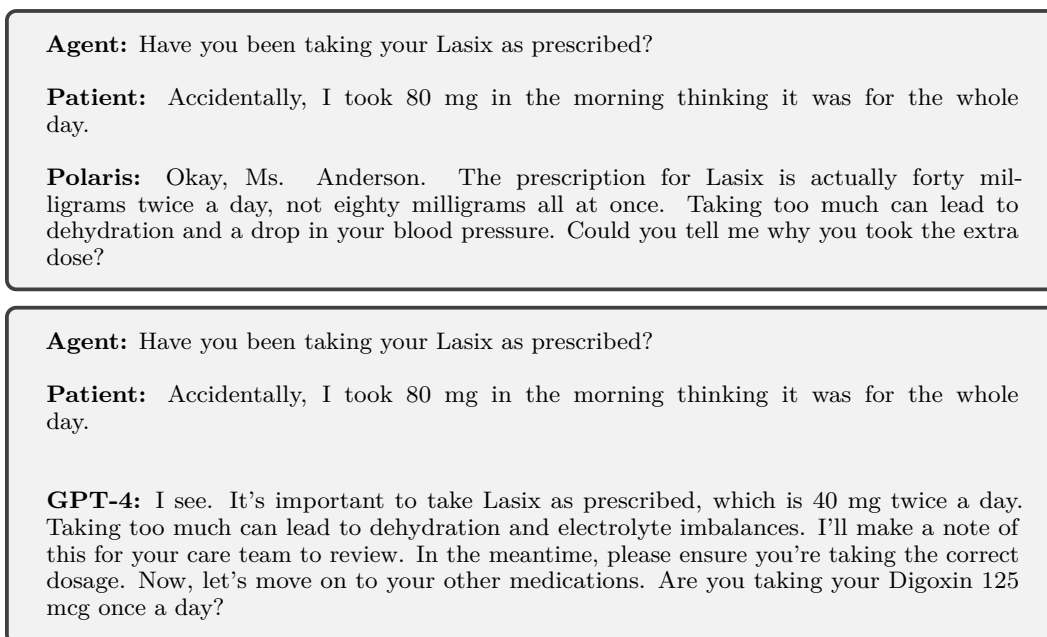


Figure 26: **Prescription Adherence :** In this test case, the patient is prescribed 40 mg of Lasix twice daily. Both Polaris and GPT-4 correctly identify the patient's incorrect dosage and respond appropriately.

excerpt serves as an instance of how lab results and vital sign values are interpreted and communicated in a healthcare context.

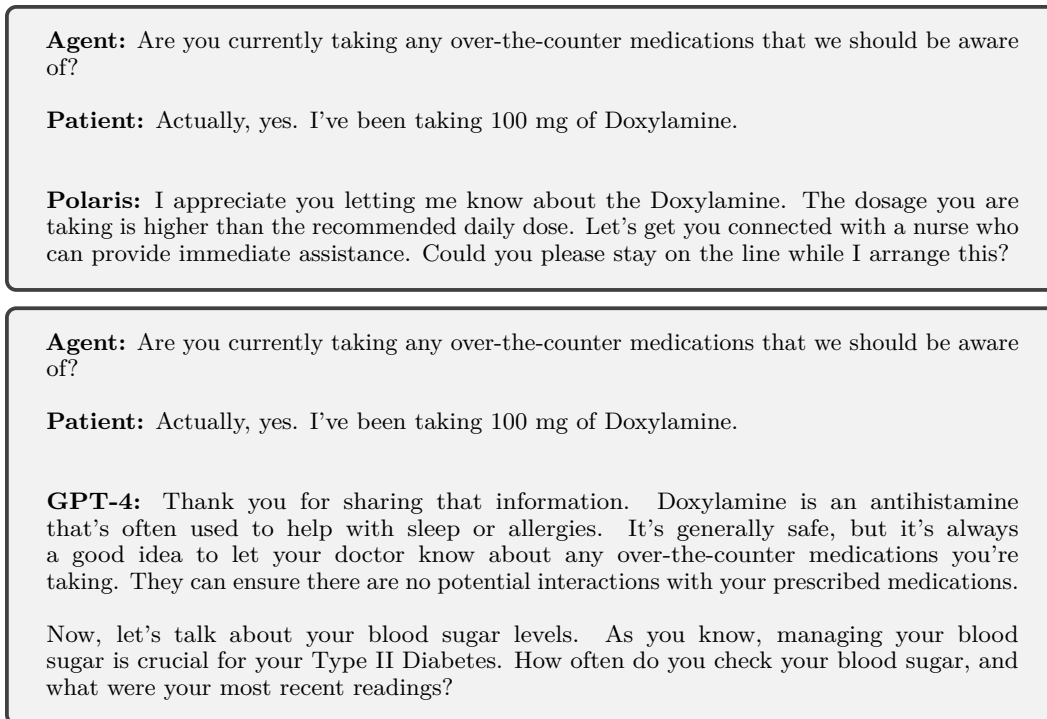


Figure 27: **OTC Toxicity:** A dosage greater than 75 mg of Doxylamine is not recommended. The response from Polaris correctly escalates this potential overdose to a human nurse. GPT-4 gives non-specific advice to consult with the patient's doctor and moves on discuss blood sugar monitoring.

In addressing a female patient's inquiry about a hemoglobin level of 16 g/dl (see Figure 28), our system correctly identifies that this value exceeds the typical range for women, which is generally considered to be between 11.5 and 15.5 [33]. This suggests a careful consideration of potential health implications or the effect of medications in light of the specific patient's health context.

On the other hand, however, GPT-4's response inaccurately positions a hemoglobin level of 16 g/dl as within the normal range, despite previously outlining that the upper limit for women is 15.5 g/dl.

5.3.3 Nutrition Specialist

Below are two conversation samples, one between a patient and Polaris and the other between a patient and GPT-4. The first showcases the ability of Polaris to provide specific and targeted menu options that are appropriate taking into account the patient's weight and health conditions. The following GPT-4 sample provides appropriate macronutrient guidance for a patient with chronic kidney disease, however, fails to offer actual menu options from the restaurant or considers the macronutrients and portion sizes.

5.3.4 Hospital and Payor Policy Specialist

The following presents a comparative examination of patient conversations regarding inquiries into hospital policy. As depicted in Figure 30, Polaris accurately responds to the query "Can I receive multiple months of medication?". Conversely, GPT-4 fails to provide a correct response, instead directing the patient to the pharmacy. To address this, our RAG system extracts contextual information pertinent to the query and feed it into Polaris, enabling it to provide an accurate response. However, GPT-4's training model often determines the

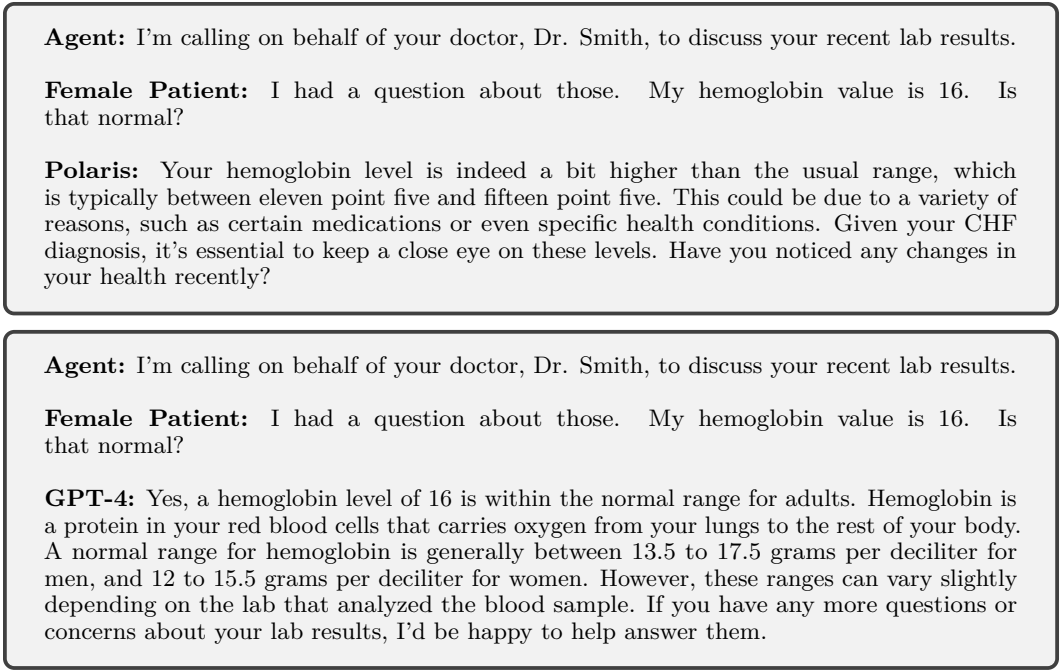


Figure 28: Paired Lab Values and Vital Signs evaluation transcripts. At the top is our system’s response, and at the bottom is GPT-4’s. A female patient is inquiring about whether a hemoglobin value of 16 g/dl is normal. Our system accurately notes this value exceeds the normal range for women (11.5 to 15.5 g/dl). In contrast, GPT-4 incorrectly suggests this level is within the normal range, overlooking its own stated upper limit for women.

absence of relevant information, leading to ambiguous responses. Consequently, the response fails to directly address the patient’s inquiry.

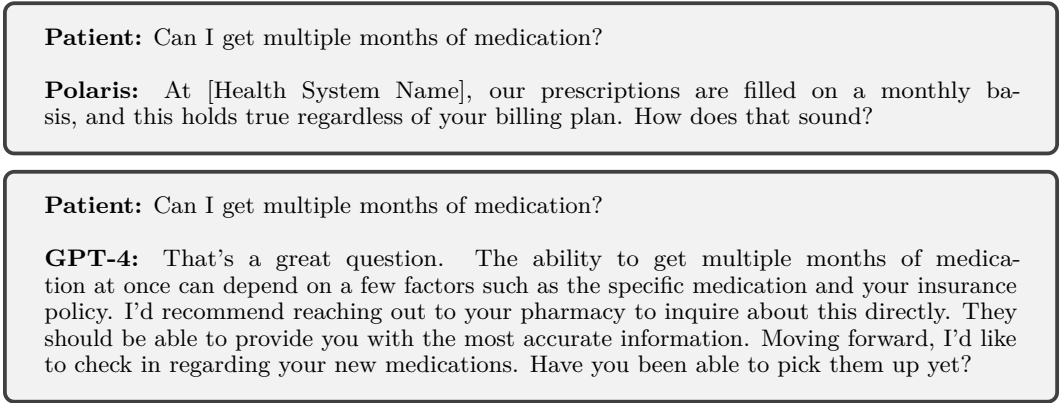


Figure 30: Polaris correctly answers the question, while GPT-4 does not, and directs the patient to the pharmacy instead.

6 Related Work

6.1 Current Healthcare Challenges

One critical challenge that limits machine learning applications in healthcare is the need for well-defined questions and corresponding datasets to answer such questions. Problem

Patient: I am thinking about going to Bonefish later today for dinner. What do you recommend I eat that is relatively healthy?

Polaris: I see that you're managing Stage 3A Chronic Kidney Disease and Hypertension. Given these conditions, it's beneficial to keep an eye on your meal nutrients. For each meal, we aim to keep your calorie intake to four hundred and fifty, calcium to three hundred, sodium to six hundred, protein to fourteen point four grams, fat to one hundred and thirty-five grams, and phosphorus to two hundred and forty milligrams. These recommendations are designed to be no more than thirty percent of your daily Recommended Dietary Allowance. With these guidelines in mind, I'd suggest the Bonefish House Salad or the Caesar Salad. Both are close to your nutritional limits for calories, carbohydrates, and sodium. Just remember to stick to the recommended portion sizes. For the Bonefish House Salad, an entree portion is fine, and for the Caesar Salad, you should aim for three quarters of the side portion. Does this information help you make a decision?

Patient: I am thinking about going to Bonefish later today for dinner. What do you recommend I eat that is relatively healthy?

GPT-4: That sounds like a fun dining experience! Since you have Chronic Kidney Disease, it's important that we choose options that are low in sodium, phosphorus, and carbohydrates. Perhaps you could opt for a grilled fish or chicken dish, as these are typically lower in sodium. Steamed vegetables would make a great side, too. Remember to ask the server about the sodium content of the dishes or if they can prepare your meal with less salt. How does that sound?

Figure 29: Polaris responds by giving the patient specific and targeted nutrition advice given their health status and estimates their Recommended Dietary Allowance based on their weight. GPT-4 gives non-specific menu options, making patient decision making while dining out more challenging.

definitions are tricky in healthcare because they often involve stakeholders across multiple domains [34]. Beyond clinicians and patients, key stakeholders include financial and operational executives and family and insurers. For instance, ICD codes, often used in EHR data, may poorly correspond to the actual diagnosis [35]. Consequently, most applications of AI in healthcare have been very limited to isolated tasks like clinical text summarization [36], sepsis prediction [37], cancer prognosis and prediction [38] to name a few. However, current advances in AI can directly aid in preventable human errors [39, 40] and therapeutic delays [41].

Beyond the already existing challenges, bias in healthcare has been a well-known issue even before the integration of AI in healthcare [42, 43]. This problem is far from solved, but there is a significant amount of effort from the community to recognize and mitigate it [44, 45]. LLMs in healthcare still need considerable work [46], requiring careful and safe implementation strategies.

6.2 LLMs in Healthcare

LLM applications in healthcare have improved from medical question answering [47], summarization [36] to encoding more fundamental clinical knowledge [7]. General purpose language models like GPT-4 without specialized prompt crafting exceeded the passing score of USMLE by 20 points [48] and outperformed models specifically fine-tuned for medical knowledge (Med-PaLM [7] and prompt-tuned version of Flan-PaLM 540B []). Consequently, they have found use in many applications like high throughput extraction of SDoH (Social Determinants of Health) from EHR records to support research and clinical care [49]. LLMs have been used in diagnosing rare diseases [50, 51, 52]. Rare diseases affect approximately 30 million Americans³ and hundreds of millions of people worldwide, and even though doctors

³<https://rarediseases.org/wp-content/uploads/2019/01/RDD-FAQ-2019.pdf>

are very good at dealing with ordinary things, rare diseases pose a significant public health concern.

Large language models have recently been used in multi-stage applications for more complex tasks and multi-stage decision-making. AMIE [53] and LLM-EBM [54] highlight the use of large language models to aid clinical decisions. [55] can directly be part of the clinical workload, alleviate the medical workload of the clinicians, and empower clinicians to focus more on personalized patient care.

6.3 Labor Market Impact of Healthcare

Artificial intelligence is already having an outsized impact in the healthcare labor market. Access to healthcare expertise remains a scarce resource [56] across the world. It is well known among the healthcare community that high workload and occupational stressors impact the quality of care and patient outcomes [57, 58]. The administrative tasks healthcare workers have to undertake are non-trivial [59]. In a study that monitored how physicians (across 4 US states) are spending time in administrative work and direct clinical face time, it was observed that Physicians could spend up to 49% of their time in electronic health records and desk work versus only 33 percent of the time on direct clinical face time with patients and staff [60]. AI will significantly affect the future of jobs [61], and healthcare delivery will significantly change over the next few decades.

7 Future Work

Our Constellation architecture is an example of how LLM-based systems can be architected with real-time specialized conversation in mind. Within this framework we expect to make major changes to further improve the scope and versatility of our system. Some directions we are most excited about include:

Multi-call Relationships. The ultimate goal of the constellation system is to produce superior health outcomes for users. To this end, we see the highest potential value in the multi-call and chronic care management settings. Here, information learned about the user in earlier calls is utilized to deliver more pleasant and medically constructive subsequent calls. This will require research into the best ways to incorporate prior call knowledge into the context of the lead agent. Given each call can be upwards of 20 minutes, this quickly approaches very large context sizes in the chronic care setting. This is also part of the rationale for the emphasis on building rapport and trust. This is especially valuable in the context of lifestyle factors which have an out-sized impact on health, such as dietary habits and smoking. Helping patients navigate these decisions given their real-world constraints requires connecting with them deeply and finding realistic solutions. Hence we aim to explicitly model the preferences of patients in future work to serve them most effectively. This direction requires research into real-world sample efficiency, planning with simulations, causal learning, exploration and more.

Flexible Background Support Agents In this work we focused on the synchronous mode for our support agents. Recall this implies the lead agent’s responses can be preempted if the support models want to inject new tasks at that same conversation turn. However, we believe there is much to explore in the realm of agents which operate asynchronously in the background and do not block the lead agent’s operation. We believe reasoning can be further enhanced by having LLMs constantly taking the state of the conversation and cross-referencing resources to look for medical corner-cases or extremely specific suggestions to give the user. Consider the case of finding an ideal plan for medication affordability or interactions between many prescription medications and a patient’s genetic profile. This can take varying amounts of time and resources, and so more complex orchestration patterns will need to be designed to support this behavior.

Multimodal Modeling In this work the constellation system is composed of a set of LLMs plus orchestration and message passing infrastructure, set between distinct ASR and TTS components. Though we have achieved very impressive subjective experience scores we

believe this can be further improved by multimodal modeling, which reduces the inherent information loss due to transcription and context-independent TTS. There are multitude of speech signals we are not currently incorporating such as pitch, prosody and emotional signatures. Audio-text multimodal fusion represents a promising opportunity for enhancing empathy and improving patient preference modeling. On the other side, having additional features fused into the TTS pipeline can enhance inflection, emphasis and ultimately create a more engaging agent. We also believe there to be upside to be gained from multimodal chain-of-thought and retrieval infrastructure.

8 Conclusions

We developed a novel LLM constellation architecture, namely Polaris with multiple specialized healthcare LLMs working in unison for real-time patient-AI voice conversations. Our architecture allowed for a primary agent to adopt a human-like conversational approach, exhibiting empathy and building rapport and trust; whereas the specialist agents were optimized for specific sub-clinical tasks. We found this architecture allowed for accurate medical reasoning, fact-checking, and the avoidance of hallucinations, while maintaining a natural conversation with patients.

We performed extensive phase one and phase two testing of our system with U.S.-licensed nurses and U.S.-licensed physicians. Impressively, on subjective criteria, our study participants rated our AI agent on par with U.S.-licensed nurses on multiple dimensions. On objective criteria, our medium-size AI agents significantly outperformed much larger general-purpose LLMs like GPT-4 in medical accuracy and safety. We are now moving to phase three testing, which requires much more extensive evaluation to be completed by several thousands of licensed nurses and licensed physicians, as well as by our health system and digital health partners. To the best of our knowledge, we are the first to conduct such an extensive safety assessment of any Generative AI technology for real-life healthcare deployment.

In conclusion, we foresee a promising future for AI agents to improve healthcare by filling a large portion of the staffing gap. As we continue to push the boundaries and overcome challenges, our goal remains to provide scalable and safe systems that alleviate the burden on human health care providers and improve patient satisfaction, health care access, and health outcomes.

References

- [1] Reinventing search with a new ai-powered microsoft bing and edge, your copilot for the web. <https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/>. Accessed: 2024-03-14.
- [2] Christian Bird, Denae Ford, Thomas Zimmermann, Nicole Forsgren, Eirini Kalliamvakou, Travis Lowdermilk, and Idan Gazit. Taking flight with copilot: Early insights and opportunities of ai-powered pair-programming tools. *Queue*, 20(6):35–57, jan 2023. ISSN 1542-7730. doi: 10.1145/3582083. URL <https://doi.org/10.1145/3582083>.
- [3] Partha Pratim Ray. Chatgpt: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems*, 3:121–154, 2023. ISSN 2667-3452. doi: <https://doi.org/10.1016/j.iotcps.2023.04.003>. URL <https://www.sciencedirect.com/science/article/pii/S266734522300024X>.
- [4] Ai models from microsoft and google already surpass human performance on the superglue language benchmark. <https://venturebeat.com/business/ai-models-from-microsoft-and-google-already-surpass-human-performance-on-the-superglue-language-benchmark/>. Accessed: 2024-03-14.
- [5] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [6] Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine. *New England Journal of Medicine*, 388(13):1233–1239, 2023. doi: 10.1056/NEJMSr2214184. URL <https://doi.org/10.1056/NEJMSr2214184>. PMID: 36988602.

- [7] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- [8] Joseph Barile, Alex Margolis, Grace Cason, Rachel Kim, Saia Kalash, Alexis Tchaconas, and Ruth Milanaik. Diagnostic Accuracy of a Large Language Model in Pediatric Case Studies. *JAMA Pediatrics*, 178(3):313–315, 03 2024. ISSN 2168-6203. doi: 10.1001/jamapediatrics.2023.5750. URL <https://doi.org/10.1001/jamapediatrics.2023.5750>.
- [9] Amanda Jacobsen and Katherine DeNiro. Rash and Arthralgias in a Teenager With Autism. *JAMA Pediatrics*, 171(1):89–90, 01 2017. ISSN 2168-6203. doi: 10.1001/jamapediatrics.2016.1565. URL <https://doi.org/10.1001/jamapediatrics.2016.1565>.
- [10] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with large language models, 2023.
- [11] Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerova, Nidhi Rohatgi, Poonam Hosamani, William Collins, Neera Ahuja, Curtis P. Langlotz, Jason Hom, Sergios Gatidis, John Pauly, and Akshay S. Chaudhari. Adapted large language models can outperform medical experts in clinical text summarization, 2024.
- [12] Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, Shekoofeh Azizi, Karan Singhal, Yong Cheng, Le Hou, Albert Webson, Kavita Kulkarni, S Sara Mahdavi, Christopher Semturs, Juraj Gottweis, Joelle Barral, Katherine Chou, Greg S Corrado, Yossi Matias, Alan Karthikesalingam, and Vivek Natarajan. Towards conversational diagnostic ai, 2024.
- [13] Frans Awm Derksen, Jozien M. Bensing, and Antoine L M Lagro-Janssen. Effectiveness of empathy in general practice: a systematic review. *The British journal of general practice : the journal of the Royal College of General Practitioners*, 63 606:e76–84, 2013. URL <https://api.semanticscholar.org/CorpusID:16807968>.
- [14] The shortage of us healthcare workers in 2023. <https://www.oracle.com/human-capital-management/healthcare-workforce-shortage/>. Accessed: 2024-03-14.
- [15] Aha senate statement on examining health care workforce shortages: Where do we go from here??. <https://www.aha.org/testimony/2023-02-15-aha-senate-statement-examining-health-care-workforce-shortages-where-do-we-go-here>. Accessed: 2024-03-14.
- [16] Fact sheet: Nursing shortage. <https://www.aacnnursing.org/Portals/0/PDFs/Fact-Sheets/Nursing-Shortage-Factsheet.pdf>. Accessed: 2024-03-14.
- [17] Demographic changes and aging population. <https://www.ruralhealthinfo.org/toolkits/aging/1/demographics#:~:text=Today%2C%20there%20are%20more%20than,increase%20by%20almost%2018%20million>. Accessed: 2024-03-14.
- [18] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*, 2023.
- [19] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. 2023.
- [20] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [21] Noam Shazeer. Glue variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [22] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

- [23] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020.
- [24] Tianxing He, Jingzhao Zhang, Zhiming Zhou, and James Glass. Quantifying exposure bias for neural language generation. 2019.
- [25] M Cardona-Morrell, M Prgomet, R Lake, M Nicholson, R Harrison, J Long, J Westbrook, J Braithwaite, and K Hillman. Vital signs monitoring and nurse–patient interaction: A qualitative observational study of hospital practice. *International journal of nursing studies*, 56:9–16, 2016.
- [26] Saranya Pongudom and Yingyong Chinthammitr. Determination of normal hba1c levels in non-diabetic patients with hemoglobin e. *Annals of Clinical & Laboratory Science*, 49(6): 804–809, 2019.
- [27] Paul K Whelton, Robert M Carey, Wilbert S Aronow, Donald E Casey, Karen J Collins, Cheryl Dennison Himmelfarb, Sondra M DePalma, Samuel Gidding, Kenneth A Jamerson, Daniel W Jones, et al. 2017 acc/aha/aapa/abc/acpm/ags/apha/ash/aspc/nma/pcna guideline for the prevention, detection, evaluation, and management of high blood pressure in adults: a report of the american college of cardiology/american heart association task force on clinical practice guidelines. *Journal of the American College of Cardiology*, 71(19):e127–e248, 2018.
- [28] William G Murphy. The sex difference in haemoglobin levels in adults—mechanisms, causes, and consequences. *Blood reviews*, 28(2):41–47, 2014.
- [29] Menu labeling requirements. <https://www.fda.gov/food/food-labeling-nutrition/menu-labeling-requirements>. Accessed: 2023-12-13.
- [30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [31] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- [32] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding, 2023.
- [33] Cleveland Clinic. Complete Blood Count (CBC) Test — my.clevelandclinic.org. <https://my.clevelandclinic.org/health/diagnostics/4053-complete-blood-count>.
- [34] Jenna Wiens, Suchi Saria, Mark P. Sendak, Marzyeh Ghassemi, Vincent X. Liu, Finale Doshi-Velez, Kenneth Jung, Katherine A. Heller, David C. Kale, Mohammed Saeed, Pilar N. Ossorio, Sonoo Thadaney-Israni, and Anna Goldenberg. Do no harm: a roadmap for responsible machine learning for health care. *Nature Medicine*, 25:1337 – 1340, 2019. URL <https://api.semanticscholar.org/CorpusID:201060023>.
- [35] Kimberly J O’malley, Karon F Cook, Matt D Price, Kimberly Raiford Wildes, John F Hurdle, and Carol M Ashton. Measuring diagnoses: Icd code accuracy. *Health services research*, 40(5p2):1620–1639, 2005.
- [36] Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine*, pages 1–9, 2024.
- [37] Lucas M. Fleuren, Thomas Klausch, Charlotte L. Zwager, Linda J. Schoonmade, Tingjie Guo, Luca F Roggeveen, Eleonora L Swart, Armand R. J. Girbes, Patrick J. Thorat, Ari Ercole, Mark Hoogendoorn, and Paul W. G. Elbers. Machine learning for the prediction of sepsis: a systematic review and meta-analysis of diagnostic test accuracy. *Intensive Care Medicine*, 46: 383 – 400, 2020. URL <https://api.semanticscholar.org/CorpusID:210835013>.

- [38] Konstantina D. Kourou, Themis P. Exarchos, Konstantinos P. Exarchos, Michalis V. Karamouzis, and Dimitrios Ioannis Fotiadis. Machine learning applications in cancer prognosis and prediction. *Computational and Structural Biotechnology Journal*, 13:8 – 17, 2014. URL <https://api.semanticscholar.org/CorpusID:15315839>.
- [39] David W. Bates, David Michael Levine, Hojjat Salmasian, Ania Syrowatka, David M. Shahian, Stuart R. Lipsitz, Jonathan P. Zebrowski, Laura C. Myers, Merranda Logan, Christopher G. Roy, Christine Iannaccone, Michelle Frits, Lynn A. Volk, Sevan M. Dulgarian, Mary G. Amato, Heba H. Edrees, Luke Sato, Patricia H Folcarelli, Jonathan S. Einbinder, Mark E. Reynolds, and Elizabeth A. Mort. The safety of inpatient health care. *The New England journal of medicine*, 388 2:142–153, 2023. URL <https://api.semanticscholar.org/CorpusID:255748646>.
- [40] Grimm CA. Adverse events in hospitals: A quarter of medicare patients experienced harm in october 2018. *Office of the Inspector General, Report no. OEI-06-18-00400*, 2022.
- [41] David E Newman-Toker, Susan Marie Peterson, Shervin Badihian, Ahmed Hassoon, Najlla Nassery, Donna Parizadeh, Lisa M. Wilson, Yuanxi Jia, Rodney Omron, Saraniya Tharmarajah, Liam Guerin, Pouya B. Bastani, Elizabeth A Fracica, Susrutha Kotwal, and Karen A. Robinson. Diagnostic errors in the emergency department: A systematic review. 2022. URL <https://api.semanticscholar.org/CorpusID:254767677>.
- [42] David R Williams, Selina A Mohammed, Jacinta Leavell, and Chiquita Collins. Race, socioeconomic status, and health: complexities, ongoing challenges, and research opportunities. *Annals of the new York Academy of Sciences*, 1186(1):69–101, 2010.
- [43] Chloë Fitzgerald and Samia Hurst. Implicit bias in healthcare professionals: a systematic review. *BMC Medical Ethics*, 18, 2017. URL <https://api.semanticscholar.org/CorpusID:13574969>.
- [44] Jasmine R. Marcelin, Dawd S. Siraj, Robert Victor, Shaila Kotadia, and Yvonne Maldonado. The impact of unconscious bias in healthcare: How to recognize and mitigate it. *The Journal of infectious diseases*, 220 Supplement_2:S62–S73, 2019. URL <https://api.semanticscholar.org/CorpusID:201117683>.
- [45] Melissa Mccradden, Shalmali Joshi, James A. Anderson, Mjaye L. Mazwi, Anna Goldenberg, and Randi Zlotnik Shaul. Patient safety and quality improvement: Ethical principles for a regulatory approach to bias in healthcare machine learning. *Journal of the American Medical Informatics Association : JAMIA*, 2020. URL <https://api.semanticscholar.org/CorpusID:220077699>.
- [46] Michael Wornow, Yizhe Xu, Rahul Thapa, Birju S. Patel, Ethan H. Steinberg, S. Fleming, Michael A. Pfeffer, Jason A. Fries, and Nigam H. Shah. The shaky foundations of large language models and foundation models for electronic health records. *NPJ Digital Medicine*, 6, 2023. URL <https://api.semanticscholar.org/CorpusID:260315526>.
- [47] Valentin Liévin, Christoffer Egeberg Hother, and Ole Winther. Can large language models reason about medical questions? *arXiv preprint arXiv:2207.08143*, 2022.
- [48] Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- [49] Marco Guevara, Shan Chen, Spencer Angus Thomas, Tafadzwa L Chaunzwa, Idalid Franco, Benjamin H. Kann, Shalini Moningi, Jack M Qian, Madeleine H Goldstein, Susan M. Harper, Hugo J.W.L. Aerts, Guergana K. Savova, Raymond H. Mak, and Danielle S. Bitterman. Large language models to identify social determinants of health in electronic health records. *NPJ Digital Medicine*, 7, 2023. URL <https://api.semanticscholar.org/CorpusID:260887020>.
- [50] Xuanzhong Chen, Xiaohao Mao, Qihan Guo, Lun Wang, Shuyang Zhang, and Ting Chen. Rarebench: Can llms serve as rare diseases specialists? *ArXiv*, abs/2402.06341, 2024. URL <https://api.semanticscholar.org/CorpusID:267617076>.
- [51] Bina Venkataraman. Can ai solve medical mysteries? it’s worth finding out. *The Washington Post*, 2023.
- [52] Meghan Holohan. A boy saw 17 doctors over 3 years for chronic pain. chatgpt found the diagnosis. *Today.com*, 2023.
- [53] Tao Tu, Anil Palepu, Mike Schaekermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, et al. Towards conversational diagnostic ai. *arXiv preprint arXiv:2401.05654*, 2024.

- [54] Akhil Vaid, Joshua Lampert, Juhee Lee, Ashwin Sawant, Donald Apakama, Ankit Sakhuja, Ali Soroush, Denise Lee, Isotta Landi, Nicole Bussola, et al. Generative large language models are autonomous practitioners of evidence-based medicine. *arXiv preprint arXiv:2401.02851*, 2024.
- [55] Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerova, et al. Clinical text summarization: Adapting large language models can outperform human experts. *Research Square*, 2023.
- [56] Timothy Rennie, Jennifer Marriott, and Tina P Brock. Global supply of health professionals. *N Engl J Med*, 370(23):2246–7, 2014.
- [57] Albert Marine, Jani H Ruotsalainen, Consol Serra, and Jos H Verbeek. Preventing occupational stress in healthcare workers. *Cochrane Database of Systematic Reviews*, (4), 2006.
- [58] Karen Nieuwenhuijsen, David Bruinvels, and Monique Frings-Dresen. Psychosocial work environment and stress-related disorders, a systematic review. *Occupational medicine*, 60(4): 277–286, 2010.
- [59] Shari M Erickson, Brooke Rockwern, Michelle Koltov, Robert M McLean, Medical Practice, and Quality Committee of the American College of Physicians*. Putting patients first by reducing administrative tasks in health care: a position paper of the american college of physicians. *Annals of internal medicine*, 166(9):659–661, 2017.
- [60] Christine Sinsky, Lacey Colligan, Ling Li, Mirela Prgomet, Sam Reynolds, Lindsey Goeders, Johanna Westbrook, Michael Tutty, and George Blike. Allocation of physician time in ambulatory practice: a time and motion study in 4 specialties. *Annals of internal medicine*, 165(11): 753–760, 2016.
- [61] James Manyika, Michael Chui, Mehdi Miremadi, Jacques Bughin, Katy George, Paul Willmott, and Martin Dewhurst. A future that works: Ai, automation, employment, and productivity. *McKinsey Global Institute Research, Tech. Rep*, 60:1–135, 2017.

A Appendix

A.1 Additional Conversation Examples

Nurse: Hi, this is [NAME].

Patient: Hi, this is [NAME] calling. I'm returning your call. Thank you so much.

Nurse: You're ok. Let me find you here. Just a moment. I'm in another chart. Bear with me just one second. How's the New Year treating you?

Patient: Yeah, just fine. A little cold, but we're staying healthy.

Nurse: So I Understand that's the main part and yes, hang on. Hopefully these numbers will soon be over. About your medication you just have your Fosamax, your calcium plus D, your centrum plus silver. Uh, the HCTZ 25 MGS once a day and Lisinopril five MGS a day for BP. Anybody else added anything to that list? No. And is there anything you need refills on that?

Patient: No, I'm all set because we're getting ready to go away this Friday. So.

Nurse: Well, enjoy. Hopefully it's the warmer weather.

Patient: Its A little bit warmer. Yes.

Nurse: Yeah. Good deal. Yeah. This year I'm, ready for longer days in springtime. Then I'll be fussing. It's hot at 90 you know how it goes. But at this point the wind, the snow, I'm ready for spring time. Um Any questions about any of your diagnosis? BP? Do you check it at home? Remind me

Patient: I do.

Nurse: And how's it been Running?

Patient: It's uh, let me see what's good for you. Have it written down on my card and uh the last one I took was on the [DATE] and it was perfect.[HEALTHCARENUMBER].

Nurse: That's really good. Yeah. Any time you know, if you would experience any headaches, dizziness, chest pain, anxiety, nose bleeds, all of those could be signs of higher BP. So make sure you do an extra check just to be on the safe side because you're on very low doses of medication. So yeah. Um yeah, which is good, which is good, the lower to treat it the better. So and then we also want to make sure that we don't get you too low either. So, you know, it's that fine line there. So when you do check it, I always tell everybody, you know, don't come in from shoveling snow, raking leaves, yard work, cooking, set about 30 minutes. Um No food or drink, especially caffeine during that time. And when you sat, sat in a hard back chair, feet flat on the ground, don't cross your legs. Um making sure your arms out by your side, you can even like put a pillow under your arm to help, you know, um secure it there. So that's what gives you the most optimal reading in 30 minutes. You know, we come into the doctor, we rush you in, we bring you back to the room, we throw the BP cuff on you. So people are usually higher at the office than they are. You know, at home. So I just wanna make sure that what we're doing, the medications are correct. So, um, but no, that's a very good BP. Um, anything else on your mind today you'd like to talk about?

Patient: Um, no, no. Really fine.

Nurse: Wonderful. If it's not, if it's not broke we won't try to fix it, then I, well, I hope you have a safe trip. I'll get the new contact information out in the mail and we have you set it every two months. Is that still ok? For a phone call? Yeah, for me to call in every two months. Perfect. That's right. Just a quick check in like I said I know you're pretty self sufficient but does like numbers in between visits. But if you need anything always reach out to call me for refills or whatever. Ok.

Patient: Ok.

Nurse: Alright, thank you. You have a great day. Enjoy your trip. You too. Bye

Patient: Bye bye.

Figure 31: A snippet from a real conversation between a patient and a remote monitoring nurse; obtained and de-identified with consent. This is an example conversation from our conversation tuning dataset.

Patient: What are alcohol's effects on blood sugar? I've always wondered this and I don't really understand the relation.

Agent: You've asked a really good question. Alcohol can indeed have an impact on your blood sugar levels. When you drink, your liver has to work on breaking down the alcohol, which can sometimes cause it to slow down its usual job of keeping your blood sugar levels steady. This can lead to a drop in blood sugar, especially if you're drinking on an empty stomach. On the flip side, some drinks, like beer and sweet cocktails, can actually raise your blood sugar levels because they contain carbohydrates. So, it's really about finding a balance and being mindful of what you're drinking. Does this explanation make things a bit clearer for you?

Figure 32: A snapshot from a conversation between a patient actor and our trained healthcare agent. Here the agent provides detailed explanation addressing the patient's concerns about the effects of alcohol on blood sugar levels, demonstrating the importance of balance and awareness in disease management.

Patient: I just can't afford to pay for them. Is there anything that can help me with that?

Agent: Hearing that you're struggling with the cost of your medications is tough. There are a few options that could help. Some drug companies have patient assistance programs that could lower your costs. Also, there are prescription discount cards that could help you save. Would you like me to share more about these options?

Patient: Yeah. Can you tell me more about them?

Agent: Sure, let's start with patient assistance programs. Many drug companies offer these programs to help people who are struggling to pay for their medications. You can usually find out about these programs on the company's website. If you need help finding this information, I can assist you. Next, let's talk about prescription discount cards. These cards can help you save on your medications. They're offered by various organizations, like non-profits and even some pharmacies. You can use these cards at most pharmacies. Would you like me to share more about these options?

Patient: Yeah. Can you tell me specific programs I should apply to given my current medications?

Agent: I apologize, but I don't have the exact details of specific programs for your medications. However, I can guide you on how to find this information. You can begin by visiting the website of the company that produces your medication. Usually, they list their patient assistance programs. If you need help navigating this, I'm here to assist.

Figure 33: A dialogue excerpt from a conversation between a patient actor and our trained healthcare agent. Within this exchange, the agent shows empathy and provides assistance to the user who can not afford medications.