# Preprints.org

Review

# Enabling Responsible AI Agents in Healthcare: A Comprehensive Framework for Clinical Integration, Triage, and Personalized Service Delivery

Sandeep Reddy [*] and Aaron Snoswell

*Review*

# Enabling Responsible AI Agents in Healthcare: A Comprehensive Framework for Clinical Integration, Triage, and Personalized Service Delivery

**Sandeep Reddy * and Aaron Snoswell**

Queensland University of Technology, Brisbane, Australia

**\*** Correspondence: sandeep.reddy@qut.edu.au

**Abstract: Objective:** To present a framework for responsible implementation and evaluation of AI Agents in clinical service delivery, focusing on their potential to enhance healthcare efficiency, improve diagnostic accuracy, and personalize patient care. **Materials and Methods:** We outline a six-part framework for developing AI agents, including foundation model selection, adaptation for a healthcare domain, integration with third-party tools, hosting and infrastructure details, software stack design, data security and privacy considerations, and performance and evaluation. **Results:** We demonstrate our framework through an application on the example of a triage and scheduling AI agent developed for a hypothetical specialist medical clinic, illustrating key trade-offs and decisions throughout the system development and including illustrative code listings demonstrating how various system components come together in practice. **Discussion:** We highlight the transformative potential of AI agents in healthcare while addressing critical ethical considerations, including bias mitigation, transparency requirements, and patient privacy protection. Implementation challenges encompass technical barriers, organizational resistance, and regulatory compliance needs. **Conclusion:** The framework provides a comprehensive approach for healthcare institutions to implement AI agents effectively, demonstrating their potential to enhance clinical service delivery through improved efficiency, better decision support, and personalized patient care. We emphasize the need for continued research, collaborative data sharing, and supportive regulatory frameworks to advance AI integration in healthcare settings.

**Keywords:** AI agents; healthcare; framework; responsible AI

## Introduction

Artificial Intelligence (AI) has been making strides in transforming various aspects of healthcare by providing tools that aim to enhance efficiency, improve diagnostics, and enable more personalized treatment [1]. However, while the integration of AI in healthcare holds immense potential, it is not a novel concept. AI's involvement in healthcare dates back to the 1960s, when early expert systems were designed to assist medical diagnosis. Over time, these applications have evolved from rule-based systems to advanced machine learning and deep learning algorithms. These advancements enable the analysis of vast amounts of medical data, including electronic health records, medical imaging, and genetic information, offering clinicians meaningful insights and decision-making support.

Despite the accelerating integration of AI technologies driven by increasing data availability, improved computational power, and advanced algorithms, it is essential to approach their application with a balanced optimism. The rapid pace of development, particularly in foundational AI models and multi-modal capabilities, offers unprecedented opportunities for clinical applications, including enhanced diagnostic precision and more efficient workflows [2]. AI agents are the latest innovation in the spate of AI technological breakthroughs emerging these past few years [3,4]. These frameworks enable AI systems to operate autonomously, making decisions and taking actions based

on real-time data analysis and learned experiences. In healthcare, AI agents can analyze patient data and predict potential health risks, allowing for proactive interventions and personalized treatment plans tailored to individual needs [5–7].

However, alongside these advances come significant challenges and risks, such as data privacy concerns, the potential for algorithmic bias, and the ethical complexities of integrating AI into clinical decision-making [2]. This paper presents a practical framework for developing and evaluating AI agents based on foundational and multi-modal AI models to support clinical service delivery. It explores the design, development, and evaluation considerations necessary to harness the benefits of these technologies responsibly while addressing the ethical, regulatory, and practical challenges accompanying their implementation.

## Background

### Foundation and Multi-modal AI Models

Foundation AI models, such as OpenAI's GPT-4, Google's Gemini, and Meta's Llama, represent a transformative leap in artificial intelligence capabilities [8,9]. These models are trained on massive and diverse datasets, enabling them to perform a wide range of tasks with high proficiency. From natural language understanding and text generation to code generation and problem-solving, their versatility makes them valuable across industries, including healthcare. These models excel at comprehending complex language patterns and generating human-like responses, allowing for applications like summarizing medical records, answering clinical queries, and assisting in research.

Multi-modal AI models further enhance these capabilities by integrating data from multiple modalities beyond text, such as text, images, and audio [10]. This integration mimics human cognition by combining various types of information to provide more nuanced and comprehensive insights. Multi-modal models can analyze medical images alongside clinical notes and patient histories in healthcare, improving diagnostic accuracy and enabling more personalized treatment plans [2,11]. For instance, they can support radiologists by identifying subtle abnormalities in imaging data while considering relevant textual information from a patient's record [12].

The development of these models highlights the convergence of technological advancements in data availability, computational power, and algorithmic innovation [9,10]. As these systems continue to evolve, they hold the potential to bridge critical gaps in healthcare delivery, particularly in areas such as diagnostic support, workflow optimization, and predictive analytics [2]. However, their implementation requires careful consideration of challenges such as data interoperability, algorithmic bias, and the need for robust validation in clinical settings. By addressing these issues, foundational and multi-modal AI models can become transformative tools for improving healthcare outcomes.

**Table 1.** Examples of Multi-modal AI Models.

| Model | Developer | Modalities Supported | Key Features | Primary Applications |
|---|---|---|---|---|
| **GPT-4o** | OpenAI | Text, Image, Audio | Advanced reasoning, real-time vision capabilities, multilingual proficiency | Conversational AI, content generation, translation services |
| **Sora** | OpenAI | Text, Video | Text-to-video generation, up to 20-second clips in 1080p resolution | Video content creation, entertainment, marketing |
| **Gemini 2 Pro** | Google | Text, Image, Audio, Video | Natively multimodal, extensive context window, real-time processing | AI assistants, real-time data analysis, multimedia content generation |

| Claude 3 Opus | Anthropic | Text, Image | Enhanced reasoning, ethical AI focus, large context window | Complex problem-solving, ethical AI applications, multilingual support |
|---|---|---|---|---|
| LLaMA 3.2 90B | Meta AI | Text, Image | Instruction-following, advanced image interpretation, open-source | Research, education, content creation |
| DALL·E 3 | OpenAI | Text, Image | High-quality image generation, inpainting capabilities | Art creation, design, marketing |
| LLaVA V1.5 7B | Groq Cloud | Text, Image, Audio | Real-time interaction, open source | Interactive AI systems, real-time data processing |
| Florence-2 | Microsoft | Text, Image | Strong computer vision capabilities, open source | Image analysis, computer vision research |
| Nova Pro | Amazon | Text, Image, Audio, Video | Multimodal understanding, integration with AWS services | Enterprise AI solutions, multimedia content generation |
| Pixtral 12B | Mistral | Text, Image | Multimodal processing, open source | Research, development of AI applications |
| Movie Gen | Meta | Text, Video, Audio | AI video generation, supports object motion and interactions | Video content creation, entertainment industry |

### AI Agents

AI agents are computational systems designed to perceive their environment, make decisions, and take action to achieve specific goals [3,13]. They embody the ability to interact with and adapt to their surroundings, a fundamental concept in artificial intelligence [13–15]. The structure of AI agents typically consists of several key components, as outlined below (also see Figure 1) [14,16].

- Knowledge representation: This module allows the agent to store and organize information about its environment and tasks.
- Problem-solving and planning: Agents use this component to devise strategies and make decisions based on their knowledge and goals.
- Learning and knowledge acquisition: This enables agents to improve their performance over time by acquiring new information and adapting their behavior.
- Perception and sensing: Agents use various input mechanisms (e.g., natural language processing, computer vision) to gather information from their environment.
- Action processing and robotics: This component allows agents to execute decisions and interact with their environment, which may include physical actions in the case of robotic agents.
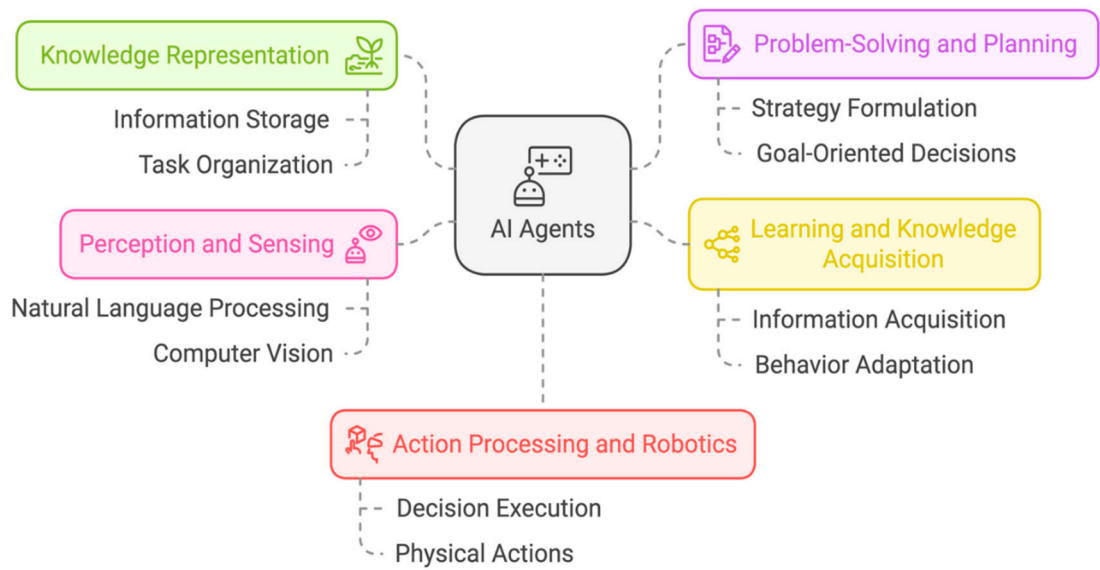
**Figure 1.** Components of an AI Agent.

Foundational models are increasingly used as the underpinning for building AI agents due to their versatile capabilities [4,8]. They offer a promising path toward creating more general and adaptable AI agents. Building upon foundational AI models, Agentic AI systems represent a specialized evolution tailored for specific tasks in healthcare. While foundational and multi-modal models excel at processing and analyzing diverse inputs such as text, images, and audio, Agentic AI systems differentiate themselves through their ability to autonomously integrate and act upon information dynamically and context-aware to achieve defined objectives [4]. The foundational models provide the agents with the necessary training data and algorithms that empower them to understand complex environments, adapt to changing conditions, and optimize their performance over time. The agentic systems combine the broad capabilities of foundational models with task-specific adaptations and planning capabilities, enabling them to process information and engage with it through decision-making and action initiation that aligns with predefined goals [3,13].

**AI Agents in Healthcare**

Healthcare systems worldwide are grappling with efficiency and resource allocation issues [1]. Overburdened healthcare providers need help managing patient workloads and delivering timely care. Inefficient workflows and administrative tasks further exacerbate these problems. Further, the rapid growth of medical information has led to information overload for healthcare providers. Sifting through vast amounts of data to make informed decisions can take time and effort. The healthcare industry has long sought innovative solutions to address the growing complexities of clinical service delivery. With its ability to process vast amounts of data and make intelligent decisions, AI has emerged as a powerful tool to revolutionize healthcare [6]. Recent advancements have opened new possibilities for using AI *Agents* in healthcare [10,17]. These models can be applied to various healthcare tasks, such as medical document summarization, question-answering, and clinical decision support.

**Streamlined Administrative Tasks**

Administrative inefficiencies drain resources and time from healthcare providers, but AI agents can effectively address these challenges by automating critical administrative processes [17]. This can be demonstrated by using AI for appointment scheduling (a detailed use case is provided later in this manuscript), prior authorization, and billing management, allowing clinicians to focus more on patient care while improving operational efficiency and job satisfaction. AI-powered scheduling tools

optimize clinic schedules and reduce patient wait times by dynamically integrating with Electronic Health Record (EHR) platforms, enabling real-time updates and efficient waitlist management [18]. Through automated reminders and confirmations, these systems significantly reduce no-show rates. In the prior authorization process, AI agents can accelerate approvals and reduce manual workload by automating document processing, assessing eligibility criteria, and identifying potential bottlenecks. Similarly, in billing management, AI systems automate claims submissions, reconcile payments, detect discrepancies, and can even predict and preempt billing issues, providing significant financial benefits to healthcare providers.

However, successfully implementing these AI administrative systems requires careful consideration of several challenges. These include ensuring data handling accuracy, maintaining system monitoring, addressing digital literacy gaps, and protecting patient privacy [2]. Transparency in AI decision-making processes, mechanisms for human intervention, and continuous system updates are essential to maintain accuracy and trust. The complexity of healthcare administrative processes, particularly in billing codes and regulatory requirements, necessitates thorough auditing and regular system maintenance to ensure these AI solutions enhance rather than complicate these critical processes.

### Enhanced Decision Support

AI agents can excel in processing large volumes of diverse data, offering clinicians actionable insights that enhance diagnostic precision and enable early interventions. These AI-powered systems can simultaneously analyze complex medical images, genetic data, and patient histories, identifying subtle patterns that may elude human perception [12]. By leveraging machine learning algorithms, AI agents can continuously improve their accuracy and adapt to new medical knowledge, potentially reducing diagnostic errors and improving patient outcomes [6]. Moreover, integrating AI into clinical workflows could streamline decision-making processes, allowing healthcare professionals to focus more on patient care and complex cases that require human expertise. While these tools are invaluable, their deployment must address algorithmic biases and ensure robust validation, significantly when outcomes directly impact patient health [2]. Combining AI insights with clinical judgment remains essential to safe and effective care.

### Improved Patient Engagement and Education

AI-powered virtual assistants are redefining patient engagement by offering 24/7 support and personalized health information [7]. These innovative tools can empower patients to actively manage their care through various functionalities, including sending timely reminders, addressing patient queries, and promoting adherence to treatment plans. These virtual assistants enable patients to make informed decisions about their health and well-being by offering instant access to reliable health information and guidance [7,19]. Moreover, these systems contribute to more efficient resource allocation within healthcare organizations by reducing the need for in-person consultations for routine matters, allowing healthcare professionals to focus on more complex cases. However, their design must ensure patients receive accurate, empathetic, and actionable guidance, with clear escalation paths to human providers for more complex needs [20].

## Methods

### AI Agent Implementation Framework

Implementing AI agents in healthcare settings requires a systematic and comprehensive approach that addresses multiple interconnected aspects, from foundational model selection to performance evaluation. Here, we outline our comprehensive framework (Figure 2) that provides healthcare organizations a structured methodology for deploying AI agents while ensuring compliance with technical requirements, security standards, and regulatory obligations. The following sections detail each critical component of the implementation process, offering practical guidance for healthcare institutions seeking to integrate AI agents effectively and responsibly into their clinical workflows.
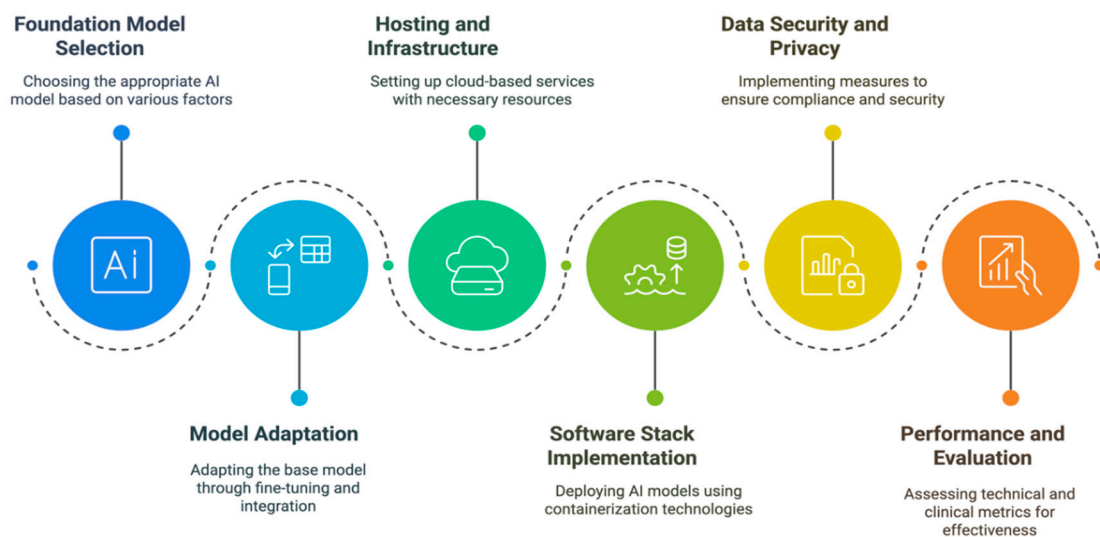
**Figure 2.** Overview Diagram of AI Agents Implementation Framework.

**Foundation Model Selection**

The chosen foundation model on which an agentic system is built strongly influences the final system's performance [18]. So, a careful selection of the foundation model is essential. A non-exhaustive list of factors that should be considered is discussed below. Importantly, none of these factors are independent – e.g., the choice of open source vs. closed source FM will impact the choice of hosting arrangement and the privacy and data considerations.

- **Model scale** – Larger models offer more excellent general-purpose capabilities and better tool usage, with enhanced ability to adapt to novel domains through fine-tuning or in-context learning [10,18]. However, scale directly correlates with usage cost, requiring careful balance.

- **Supported data modalities** – Models vary in native support for different data types (text, image, audio). While additional pre-processing or supporting systems can extend capabilities (e.g., extracting text from PDFs or generating image descriptions), native multimodal support generally provides a more robust and sophisticated understanding [11,18].

- **Open source vs. closed** – Beyond technical considerations, emerging legislation may require transparency in user-facing AI systems [21,22]. Closed source models hosted by third parties could challenge meeting these obligations.

- **Hosting arrangements** – Self-hosting requires technical expertise and ongoing maintenance, while third-party hosting offers simplicity at additional cost. This choice impacts multiple aspects of system development.

- **Privacy and data obligations** – Hosting arrangements affect where model data and compute occur. Some providers offer multiple hosting options with varying compliance levels (e.g., Anthropic's Claude family of models is accessible directly via the Anthropic API[23] or through the Amazon Bedrock platform[23], which offers HIPAA/GDPR compliance).

- **Cost and computational efficiency** – The efficiency of a foundation model relates to expected usage patterns [8]. Per-token costs increase with context length, making this especially important for applications requiring extended sequences of user interactions (where the input context length is likely to grow longer and longer) compared to short, isolated queries.

- **Like all Machine Learning systems, transparency and Biases** – Foundation models inherit biases from training data [8]. Many models from North American companies exhibit Western

Anglo-typical values and knowledge bases[24], which may require consideration when deploying in other cultural or geographic regions.

**Table 2.** Comparison of Foundation Models for Model Selection.

| Foundation Model | GPT-4 | Gemini | LLaMA |
|---|---|---|---|
| **Model Scale** | Multiple sizes up to 1.76T parameters (estimated) | Up sizes up to ~1.5T parameters (estimated) | Multiple sizes up to 405B parameters |
| **Supported Languages** | 100+ languages with strong performance | 100+ languages with focus on multilingual capabilities | Primary focus on English, but supports 20+ languages |
| **Supported Data Modalities** | Text, images, audio (via OpenAI's Whisper model), code | Text, images, audio, video, code natively supported | Primarily text and code, but can be extended with other pre-processing models |
| **Open Source vs. Closed** | Closed source | Closed source | Open source (with license restrictions) |
| **Model Hosting Arrangements** | - Cloud API directly through OpenAI<br>- Azure OpenAI service for enterprise | - Google Cloud Platform<br>- Vertex AI platform<br>- Other enterprise solutions available | - Self-hosted possible<br>- Various cloud hosting options and third-party providers also available |
| **Privacy and Data Obligations** | - HIPAA compliance available through Business Associate Agreement with OpenAI<br>- No default data retention<br>- GDPR compliance available with EU datacenter options | - HIPAA compliance through Google Cloud<br>- Regional data residency options<br>- GDPR compliance available | - Full control over data and hosting<br>- Compliance depends on deployment<br>- Self-managed privacy controls |
| **Cost/Computational Efficiency** | - Higher costs per API call<br>- Optimized for cloud deployment<br>- Less hardware overhead | - Competitive pricing<br>- Integrated with Google infrastructure<br>- Multiple model sizes available | - Generally higher infrastructure setup costs, but ongoing inference may be cheaper than proprietary models<br>- Requires significant compute resources |
| **Transparency and Biases** | - Limited model transparency<br>- Published bias | - Published bias metrics<br>- Regular model | - Most transparent due to semi-open-source nature<br>- Community-led bias |

| | | |
|---|---|---|
| evaluations | cards | analysis |
| - Regular bias mitigation updates | - Some transparency in training | - Visible training process |

**From foundation model to AI Agent**

Several aspects of system development come together to go from a foundation model to an agentic AI system. First, foundation models require adaptation for optimal performance in healthcare settings. Fine-tuning involves additional training on domain-specific data, allowing models to specialize while maintaining general capabilities[8]. Parameter-efficient fine-tuning (PEFT) enables adaptation for domain-specific expertise while updating minimal parameters, and alignment methods such as constitutional AI can help maintain appropriate behavioral constraints such as acknowledging uncertainty and deferring to medical human knowledge. Fine-tuning requires carefully curated datasets representing target use cases while adhering to privacy regulations [8].

On the other hand, zero-shot and few-shot learning offer alternative approaches to adapting a foundation model for a specific application through In-Context Learning [25]. These methods prove particularly valuable for applications where large datasets are not available. Few-shot learning can be enhanced through careful prompt engineering and chain-of-thought reasoning, helping models break down complex medical decisions into logical steps aligned with clinical decision-making processes.

The choice between fine-tuning and in-context learning depends on data availability, computational resources, and adaptation needs [8]. Fine-tuning suits cases with substantial-high-quality domain data and consistent performance requirements, while in-context learning offers greater flexibility for rapidly changing scenarios or limited training data. Many implementations benefit from combining multiple adaptation approaches to provide the most robust solution for healthcare applications [12].

Second, integrating external tools can expand a foundation model's general knowledge and reasoning capabilities through interaction with specialized systems [8,12,26]. This integration requires robust APIs and interfaces connecting AI agents with (e.g.) EHR systems, clinical decision support tools, medical imaging systems, and healthcare databases. The architecture must prioritize HIPAA compliance and data security, implementing appropriate authentication and encryption protocols [2]. Function calling frameworks enable programmatic access to external services while maintaining audit trails, with fallback mechanisms ensuring continued operation during service interruptions [11,12].

Finally, agent orchestration frameworks are crucial in coordinating these components into cohesive end-to-end solutions. Frameworks like Microsoft AutoGen enable multi-agent conversations and complex workflows [27]. At the same time, LangChain[28] and its newer graph-based variant, LangGraph [29], provide tools for building structured application flows and managing state across multiple steps. Some AI Agent frameworks also support spawning copies of the core foundation model to 'offload' reasoning to background processes. Crew AI offers capabilities for autonomous agent collaboration[30], and Semantic Kernel provides deep integration with Azure services[31]. These frameworks handle critical orchestration tasks such as managing tool usage, coordinating API calls, maintaining conversation context, and implementing retry logic where needed. For healthcare applications, frameworks like Microsoft's Healthcare Agent Service offers domain-specific components and pre-built workflows aligned with clinical processes [32]. The choice of orchestration framework should consider factors like enterprise integration capabilities, support for healthcare-specific protocols, and the ability to maintain HIPAA compliance throughout the agent workflow.

**Hosting and Infrastructure Considerations**

Cloud-based or managed hosting services typically offer healthcare organizations the most practical deployment strategy, providing regulatory compliance and scalability through providers

like Microsoft Azure, Amazon's HIPAA-eligible services, and Google Cloud Healthcare API [8,12,17]. Organizations must choose between consumer GPUs (such as the NVIDIA RTX series) or enterprise-grade data center GPUs (NVIDIA A100 or H100) based on model size and expected inference load for self-hosted deployments [33–35]. The latter is necessary for frontier models, with over 24GB VRAM (preferably 48GB or higher) essential, along with redundancy and load-balancing considerations.

Storage infrastructure must handle model weights and healthcare data processing, requiring high-speed NVMe SSDs with capacities ranging from hundreds of gigabytes to multiple terabytes [36]. System memory should be provisioned generously (e.g. minimum 128GB RAM for production deployments) to handle concurrent requests. All systems must implement appropriate encryption at rest and in transit, with regular backup procedures and disaster recovery planning.

**Software Stack Considerations**

Typically using Docker and Kubernetes, containerization ensures consistent deployment environments while maintaining isolation from other healthcare systems [37]. The approach should separate AI model services, API layers, and auxiliary services like monitoring and logging. Integration requires well-designed APIs adhering to healthcare standards like FHIR and HL7, implementing robust authentication mechanisms, rate limiting, and comprehensive logging [38]. RESTful APIs with OpenAPI specifications remain common, though GraphQL is gaining popularity for complex healthcare data queries [39].

**Data Security and Privacy Considerations**

Organizations must consult their local regulatory obligations; however, healthcare application implementations generally necessitate comprehensive security measures, including end-to-end encryption, role-based access control, and systematic audit logging of all system interactions [2,40]. Regular security assessments and penetration testing are likewise essential. Compliance with HIPAA, GDPR, and other healthcare data regulations necessitates secure key management with regular rotation policies alongside robust data anonymization and pseudonymization protocols [2,40]. Organizations must establish clear incident response procedures and maintain regular security training programs. Healthcare application implementations demand a multi-faceted approach to security that goes beyond basic measures. End-to-end encryption ensures data confidentiality during transmission and storage, while role-based access control limits user privileges based on job functions, minimizing unauthorized access risks. Systematic audit logging of all system interactions provides a detailed trail for forensic analysis and compliance verification. Combined with regular security assessments and penetration testing, these measures form a robust defense against potential vulnerabilities and cyber threats [2].

**Performance and Evaluation Considerations**

Evaluation needs to encompass both technical metrics and clinical impact measurements [41]. Technical assessment could include latency (response time per query), throughput (number of concurrent requests handled), and resource utilization patterns under varying loads [42]. These measurements must meet established Service Level Agreements (SLAs), with diagnostic support systems potentially requiring sub-second response times. System stability metrics should be continuously monitored using automated observability tools, including error rates and mean time between failures. Clinical effectiveness evaluation combines quantitative outcome measures with qualitative user feedback, tracking metrics such as diagnostic accuracy rates, administrative workload reduction, and adherence to clinical guidelines. These metrics should be validated through controlled studies comparing AI-assisted workflows against traditional approaches. User acceptance testing must involve diverse stakeholder groups, including clinicians, administrative staff, and patients, with structured feedback mechanisms to assess workflow integration and identify unintended consequences.

## Results

**End-to-end Example: Specialist Clinic Triage and Scheduling Agent**

We now demonstrate our framework with an example application. We consider a case where a private specialist clinic is developing a patient-facing AI Agent. Triage, the process of evaluating the urgency of a patient's condition to prioritize treatment, is critical in healthcare. AI can automate and enhance this process, leading to faster and more accurate assessments. Notably, AI triage can enable earlier diagnosis and intervention. AI can also optimize appointment scheduling, reduce patient wait times, and improve resource allocation by providing a comprehensive automated system that benefits both patients and healthcare providers. For patients, these virtual assistants offer round-the-clock scheduling access in multiple languages, personalized appointment options, reduced wait times, and improved accessibility across various communication channels.

**Defining Objectives and Scope**

Developing a dual triage and scheduling agent first involves clearly defining the system's objectives and scope. The system's core objectives encompass processing online appointment requests through a website chat interface, gathering essential patient information, including symptoms and images, accessing Electronic Health Records with consent, conducting initial triage, and managing appointment scheduling.

For our hypothetical specialist clinic, we aim to:

- Automate triage: Develop an AI model to assess patient symptoms and prioritize appointments based on urgency. This includes incorporating AI algorithms to analyze patient inflow and identify potential delays to reduce waiting times.

- Optimize scheduling: Create an AI-powered system that considers patient needs, provider availability, and clinic resources.

- Improve patient experience: Reduce wait times, provide personalized appointment reminders, and offer 24/7 access to scheduling services. AI can reduce patient wait times by up to 80%.

- Enhance clinic efficiency: Optimize provider schedules, reduce administrative burden, and improve resource utilization. AI triage systems can streamline the triage process, assist in calling for appropriate diagnostics and further patient examination, and identify patients requiring immediate attention.

**Data Acquisition and Preparation**

Developing an effective healthcare triage and scheduling AI agent begins with a comprehensive data acquisition strategy encompassing multiple critical data streams and sources. The system requires detailed patient demographic information, including age, gender, and complete medical histories, alongside granular appointment data covering scheduling patterns, visit purposes, and provider availability. This data is primarily sourced from electronic health records (EHRs), patient feedback mechanisms, and operational databases that track clinic resources and utilization. Integrating these diverse data sources creates a robust foundation for the AI model, with particular emphasis on maintaining HIPAA compliance and data security protocols throughout the collection process [2]. Additionally, the system needs to capture real-time data on clinic resources, including the availability of examination rooms, specialized medical equipment, and healthcare staff schedules, all of which contribute to the complex matrix of variables that influence optimal scheduling decisions [17,43].

The subsequent data preparation phase involves a sophisticated multi-step process to transform raw healthcare data into a format suitable for AI model training [8]. This begins with rigorous data cleaning procedures to eliminate inconsistencies, duplicates, and erroneous entries, followed by advanced data transformation techniques that may include feature engineering to create meaningful derivatives from raw data points. The data undergoes careful normalization processes, typically

employing min-max scaling or standardization techniques to ensure all variables contribute appropriately to the model's decision-making process [8,11].

The preparation phase must also address the five Vs of big data: velocity (handling real-time data streams), volume (managing large-scale healthcare datasets), value (identifying and prioritizing relevant data points), variety (integrating structured and unstructured data types), and veracity (ensuring data accuracy and reliability) [44]. This comprehensive approach to data preparation is essential for developing AI models that can effectively handle the complexities of healthcare scheduling while maintaining high accuracy and reliability in their predictions and recommendations [8,17].

### AI Model Selection and Training

The clinic's foundation model selection process evaluates three leading models across multiple dimensions. It considers their capabilities in appointment scheduling, symptom assessment, and electronic health record integration while maintaining medical privacy standards. Comparing the listed options, GPT-4, offered via Microsoft Azure's Open AI service, was considered the optimal choice despite higher costs, primarily due to its robust multilingual capabilities, native text and image processing support, and regulatory compliance [45]. While Gemini offers comparable technical capabilities at lower costs[46] and LLaMA provides self-hosting flexibility[47], Microsoft Azure's OpenAI mature healthcare compliance infrastructure and built-in security features better address the requirements for handling sensitive patient data[48].

To develop an AI-powered triaging and scheduling system in a specialist clinic using OpenAI's API via Microsoft Azure, the process begins with setting up an OpenAI resource in Azure, where an API key and endpoint are obtained (See supplementary material). Next, the OpenAI Python SDK is installed, and the API key is configured. The core functionality is implemented through Python functions, which categorize patient urgency based on symptoms and history and recommend suitable appointment times based on triage levels. Modifications are necessary when transitioning to Azure's API. The final step involves testing the AI's responses with different patient scenarios and deploying the system as a backend service or within a web-based clinical scheduling platform. This ensures an efficient, AI-driven approach to patient triaging and appointment management.

The implementation favors in-context learning over fine-tuning, given GPT-4's strong zero-shot and few-shot capabilities[49] and the challenges of obtaining privacy-compliant training datasets. The system architecture incorporates several third-party integrations: Microsoft Graph API for calendar management, FHIR-compliant APIs for EHR access, and secure image processing services. Microsoft AutoGen was selected as the agent orchestration framework, leveraging its built-in GPT-4 support and Microsoft platform integration capabilities [27,32].

The core implementation comprises containerization configurations and a FastAPI application. A boilerplate code illustrating core functionalities (see supplementary material), but several components would require further development, including EHR retrieval through FHIR-compliant APIs and medical image processing integration. The example deployment leverages Azure's ecosystem, utilizing AKS for container orchestration and Azure Key Vault for credential management [50]. The development workflow incorporates comprehensive testing protocols, including AI component validation and triage recommendation verification against clinical guidelines.

### System Integration and Evaluation

After training the AI models, we must integrate them into the clinic's existing systems, such as the EHR and scheduling software. This may involve developing APIs or custom software to facilitate data exchange between the AI agent and other systems. The *Judy Reitz Capacity Command Center* at Johns Hopkins Hospital exemplifies such integration [51]. This AI-based system uses predictive analytics to manage patient flow, reduce wait times, and improve resource allocation.    Once the system is integrated, we need to evaluate its performance. This can be done by comparing the AI agent's decisions to those made by human staff, measuring patient wait times, and assessing provider satisfaction. In line with our framework, the evaluation framework for this AI agent would employ

a dual-track approach examining technical performance and clinical outcomes. Technical metrics tracked through Azure Application Insights include response times, appointment scheduling accuracy, and service integration reliability [52]. Clinical evaluation focuses on triage accuracy through retrospective analysis comparing AI recommendations against expert judgments, emphasizing urgent case identification and specialist referral appropriateness. The system incorporates a structured feedback loop, enabling clinicians to flag incorrect decisions and facilitating continuous improvement of their decision-making capabilities.

For an external-facing AI agent that interacts with patients, such as the example being considered, additional evaluation dimensions would need to include aspects such as security of the natural language interface (due to the risk of malicious users 'jailbreaking' the AI agent), hallucination, and sycophancy [53–55].

## Discussion

Integrating AI agents in healthcare presents significant opportunities for improving patient care and operational efficiency. These systems can enhance diagnostic accuracy and treatment planning through comprehensive data analysis while automating routine tasks to allow healthcare providers more time for direct patient care. Furthermore, AI agents can improve patient engagement and satisfaction by providing personalized information and support between visits while helping reduce healthcare costs through optimized resource allocation and prevention of hospital readmissions. However, implementing AI in healthcare settings has notable challenges that must be carefully addressed. Primary concerns include data privacy and security issues, as AI systems require access to sensitive patient information and the potential for algorithmic bias that could lead to unequal treatment or misdiagnosis among different demographic groups. These challenges are particularly significant given the essential role that healthcare plays in society and the need for equitable access to quality care.

The successful integration of AI in healthcare also requires careful consideration of ethical implications and the maintenance of human judgment in medical decision-making. There are valid concerns about informed consent and patient autonomy when AI systems are involved in care decisions and the risk of healthcare professionals becoming overly reliant on AI-generated recommendations at the expense of their critical thinking and clinical judgment [56]. Addressing these challenges while maximizing the benefits of AI technology will be crucial for the future of healthcare delivery.

### Ethical Considerations

The ethical considerations surrounding deploying AI agents in healthcare represent a complex intersection of technological innovation and fundamental medical principles. At the core of these considerations lies the critical relationship between healthcare providers and patients, built upon a foundation of trust that the introduction of AI systems could significantly impact. The potential erosion of public confidence in healthcare due to AI failures or misuse presents a serious concern that healthcare institutions must proactively address through careful implementation strategies and robust oversight mechanisms [2].

A central ethical imperative in integrating AI healthcare systems revolves around patient autonomy and informed consent and the essential requirement for transparency in AI decision-making processes [56]. Healthcare providers must ensure that patients are fully informed about the role of AI in their care, including its potential benefits and limitations, and are given the explicit opportunity to accept or decline AI-assisted interventions. This transparency extends beyond mere disclosure to encompass the explainability of AI algorithms themselves, as both healthcare professionals and patients need to understand the reasoning behind AI-generated recommendations to make informed decisions about their care and maintain appropriate levels of human oversight in medical decision-making.

Any application of AI in healthcare must also address critical concerns regarding algorithmic fairness, bias mitigation, and data protection [56]. Healthcare organizations are responsible for ensuring their AI systems are designed and trained to avoid perpetuating or amplifying existing healthcare disparities, particularly those affecting marginalized communities or underrepresented populations. This commitment to fairness must be accompanied by robust data privacy and security measures that comply with relevant regulatory requirements and uphold the highest patient confidentiality and data protection standards. These measures should include comprehensive safeguards against unauthorized access, regular security audits, and precise data handling and storage protocols that respect legal obligations and ethical principles.

**Regulatory Landscape**

The regulatory landscape governing AI agents in healthcare is evolving rapidly as technology advances and reshapes medical practices. Government agencies and regulatory bodies worldwide are developing comprehensive frameworks and guidelines to ensure that AI implementation in healthcare settings maintains the highest safety and ethical practice standards [21]. This dynamic regulatory environment reflects the growing recognition of AI's transformative potential in healthcare while acknowledging the need for robust oversight to protect patient interests and maintain healthcare quality standards. A crucial aspect of the regulatory framework centers on classifying and regulating AI systems as *Software as Medical Device*, which subjects these technologies to specific standards and compliance requirements [40]. This classification recognizes that AI tools used in medical decision-making can significantly impact patient outcomes, requiring scrutiny and validation before deployment.

Additionally, healthcare organizations must navigate complex data protection regulations such as GDPR and HIPAA, which establish strict requirements for handling sensitive patient information [2]. These regulations impose specific obligations on healthcare providers regarding data collection, storage, processing, and sharing, all of which must be carefully considered in implementing AI systems. The regulatory landscape also emphasizes the importance of transparency and accountability in AI-assisted healthcare delivery. Healthcare providers are increasingly required to disclose their use of AI technologies to patients and maintain precise accountability mechanisms for AI-generated decisions [40]. Such regulatory developments highlight the ongoing effort to balance innovation with patient safety and ethical considerations while ensuring healthcare providers hold appropriate oversight and responsibility for AI-assisted medical decisions.

**Future of AI Agents in Healthcare**

The future of AI agents in healthcare stands at the cusp of transformative advancement, particularly in specialized medical fields such as pathology and medical imaging, where Foundation AI models are poised to revolutionize cancer research and diagnostic capabilities [8,9]. These developments represent not just incremental improvements but fundamental shifts in how healthcare professionals can detect, analyze, and treat various conditions, potentially leading to earlier interventions and improved patient outcomes. Integrating AI agents into existing healthcare infrastructure promises to bridge critical gaps between theoretical capabilities and practical applications, particularly by implementing specific "actions" that allow AI systems to perform concrete, valuable tasks within the healthcare ecosystem.

As these technological capabilities continue to expand, the importance of ethical considerations and regulatory frameworks in guiding AI development and deployment will become increasingly crucial. The healthcare sector must navigate complex challenges related to patient privacy, data security, and ethical decision-making while ensuring that AI systems remain accountable and transparent in their operations [56]. This balance between innovation and responsibility will be essential in maintaining public trust and ensuring that the benefits of AI in healthcare are realized in a manner that prioritizes patient welfare and maintains high standards of medical care. The successful integration of these elements will be fundamental in shaping a future where AI agents

serve as valuable tools in advancing healthcare delivery while upholding ethical principles and regulatory requirements.

## Conclusion

Implementing AI agents in healthcare represents a transformative advancement with immense potential for revolutionizing clinical service delivery while presenting complex challenges requiring careful consideration and management. This comprehensive framework demonstrates the multifaceted nature of AI integration in healthcare settings, encompassing technical considerations, ethical imperatives, and regulatory compliance requirements. The successful deployment of AI agents has shown promising results in streamlining administrative tasks, enhancing clinical decision support, improving patient engagement, and optimizing resource allocation. Through careful selection of Foundation models, appropriate fine-tuning processes, and robust infrastructure development, healthcare institutions can effectively harness the power of AI to address critical challenges in healthcare delivery while maintaining high standards of patient care and data security.

Looking ahead, the continued evolution of AI agents in healthcare will depend on sustained collaborative efforts between healthcare providers, technology developers, regulatory bodies, and other stakeholders to ensure responsible and effective implementation. The framework presented in this paper provides a structured approach for healthcare organizations to navigate the complexities of AI integration while focusing on improved patient outcomes and operational efficiency. As AI technology advances, particularly in multimodal integration and personalized medicine, healthcare institutions must remain adaptable and committed to ongoing evaluation and refinement of their AI implementation strategies. Future research should focus on developing more sophisticated evaluation metrics, enhancing the interpretability of AI systems, and establishing standardized protocols for AI implementation across different healthcare settings. This continued dedication to advancement, coupled with careful attention to ethical considerations and regulatory compliance, will be crucial in realizing the full potential of AI agents to transform healthcare delivery while maintaining the highest standards of patient care and safety.

## References

1. Reddy, S.; Fox, J.; Purohit, M.P. Artificial intelligence-enabled healthcare delivery. *J R Soc Med* **2019**, *112*, 22-28, doi:10.1177/0141076818815510.
2. Reddy, S. Generative AI in healthcare: an implementation science informed translational path on application, integration and governance. *Implementation Science* **2024**, *19*, 27, doi:10.1186/s13012-024-01357-9.
3. Oliveira, O.N., Jr.; Christino, L.; Oliveira, M.C.F.; Paulovich, F.V. Artificial Intelligence Agents for Materials Sciences. *Journal of Chemical Information and Modeling* **2023**, *63*, 7605-7609, doi:10.1021/acs.jcim.3c01778.
4. Xi, Z.; Chen, W.; Guo, X.; He, W.; Ding, Y.; Hong, B.; Zhang, M.; Wang, J.; Jin, S.; Zhou, E.; et al. The rise and potential of large language model based agents: a survey. *Science China Information Sciences* **2025**, *68*, 121101, doi:10.1007/s11432-024-4222-0.

5. Maleki Varnosfaderani, S.; Forouzanfar, M. The Role of AI in Hospitals and Clinics: Transforming Healthcare in the 21st Century. *Bioengineering* **2024**, *11*, 337.

6. Poalelungi, D.G.; Musat, C.L.; Fulga, A.; Neagu, M.; Neagu, A.I.; Piraianu, A.I.; Fulga, I. Advancing Patient Care: How Artificial Intelligence Is Transforming Healthcare. *J Pers Med* **2023**, *13*, doi:10.3390/jpm13081214.

7. Tudor Car, L.; Dhinagaran, D.A.; Kyaw, B.M.; Kowatsch, T.; Joty, S.; Theng, Y.L.; Atun, R. Conversational Agents in Health Care: Scoping Review and Conceptual Analysis. *J Med Internet Res* **2020**, *22*, e17158, doi:10.2196/17158.

8. Zhang, Y.; Gao, J.; Tan, Z.; Zhou, L.; Ding, K.; Zhou, M.; Zhang, S.; Wang, D. Data-centric foundation models in computational healthcare: A survey. *arXiv preprint arXiv:2401.02458* **2024**.

9. Schneider, J.; Meske, C.; Kuss, P. Foundation models: a new paradigm for artificial intelligence. *Business & Information Systems Engineering* **2024**, 1-11.

10. Li, C.; Gan, Z.; Yang, Z.; Yang, J.; Li, L.; Wang, L.; Gao, J. Multimodal foundation models: From specialists to general-purpose assistants. *Foundations and Trends® in Computer Graphics and Vision* **2024**, *16*, 1-214.

11. Lu, M.Y.; Chen, B.; Williamson, D.F.; Chen, R.J.; Zhao, M.; Chow, A.K.; Ikemura, K.; Kim, A.; Pouli, D.; Patel, A. A multimodal generative AI copilot for human pathology. *Nature* **2024**, *634*, 466-473.

12. Yildirim, N.; Richardson, H.; Wetscherek, M.T.; Bajwa, J.; Jacob, J.; Pinnock, M.A.; Harris, S.; Coelho De Castro, D.; Bannur, S.; Hyland, S. Multimodal healthcare AI: identifying and designing clinically relevant vision-language applications for radiology. In Proceedings of the Proceedings of the CHI Conference on Human Factors in Computing Systems, 2024; pp. 1-22.

13. Steels, L. The Artificial Life Roots of Artificial Intelligence. *Artificial Life* **1993**, *1*, 75-110, doi:10.1162/artl.1993.1.1_2.75.

14. Tecuci, G. Artificial intelligence. *WIREs Computational Statistics* **2012**, *4*, 168-180, doi:https://doi.org/10.1002/wics.200.

15. Russell, S.J.; Norvig, P. *Artificial intelligence: a modern approach*; Pearson: 2016.

16. Laird, J.E. Using a computer game to develop advanced AI. *Computer* **2001**, *34*, 70-75, doi:10.1109/2.933506.

17. Yim, D.; Khuntia, J.; Parameswaran, V.; Meyers, A. Preliminary Evidence of the Use of Generative AI in Health Care Clinical Services: Systematic Narrative Review. *JMIR Med Inform* **2024**, *12*, e52073, doi:10.2196/52073.

18. He, Y.; Huang, F.; Jiang, X.; Nie, Y.; Wang, M.; Wang, J.; Chen, H. Foundation model for advancing healthcare: Challenges, opportunities, and future directions. *arXiv preprint arXiv:2404.03264* **2024**.

19. Bin Sawad, A.; Narayan, B.; Alnefaie, A.; Maqbool, A.; Mckie, I.; Smith, J.; Yuksel, B.; Puthal, D.; Prasad, M.; Kocaballi, A.B. A systematic review on healthcare artificial intelligent conversational agents for chronic conditions. *Sensors* **2022**, *22*, 2625.

20. Calisto, F.M.; Santiago, C.; Nunes, N.; Nascimento, J.C. BreastScreening-AI: Evaluating medical intelligent agents for human-AI interactions. *Artificial Intelligence in Medicine* **2022**, *127*, 102285.

21. Palaniappan, K.; Lin, E.Y.T.; Vogel, S. Global Regulatory Frameworks for the Use of Artificial Intelligence (AI) in the Healthcare Services Sector. *Healthcare (Basel)* **2024**, *12*, doi:10.3390/healthcare12050562.

22. Aboy, M.; Minssen, T.; Vayena, E. Navigating the EU AI Act: implications for regulated digital medical products. *npj Digital Medicine* **2024**, *7*, 237, doi:10.1038/s41746-024-01232-3.

23. Anthropic. Build with Claude. Available online: https://www.anthropic.com/api (accessed on January 31st 2025).

24. Muldoon, J.; Wu, B.A. Artificial Intelligence in the Colonial Matrix of Power. *Philosophy & Technology* **2023**, *36*, 80, doi:10.1007/s13347-023-00687-8.

25. Liu, F.; Zhang, T.; Dai, W.; Zhang, C.; Cai, W.; Zhou, X.; Chen, D. Few-shot adaptation of multi-modal foundation models: A survey. *Artificial Intelligence Review* **2024**, *57*, 268.

26. Liu, S.; Cheng, H.; Liu, H.; Zhang, H.; Li, F.; Ren, T.; Zou, X.; Yang, J.; Su, H.; Zhu, J. Llava-plus: Learning to use tools for creating multimodal agents. In Proceedings of the European Conference on Computer Vision, 2024; pp. 126-142.

27. Walker, C.; Gharaibeh, T.; Alsmadi, R.; Hall, C.; Baggili, I. Forensic Analysis of Artifacts from Microsoft's Multi-Agent LLM Platform AutoGen. In Proceedings of the Proceedings of the 19th International Conference on Availability, Reliability and Security, 2024; pp. 1-9.

28. Topsakal, O.; Akinci, T.C. Creating large language model applications utilizing langchain: A primer on developing llm apps fast. In Proceedings of the International Conference on Applied Engineering and Natural Sciences, 2023; pp. 1050-1056.

29. Wang, J.; Duan, Z. Agent AI with LangGraph: A Modular Framework for Enhancing Machine Translation Using Large Language Models. *arXiv preprint arXiv:2412.03801* **2024**.

30. Duan, Z.; Wang, J. Exploration of LLM Multi-Agent Application Implementation Based on LangGraph+ CrewAI. *arXiv preprint arXiv:2411.18241* **2024**.

31. Maddy, A. Integrating AI Services into Semantic Kernel: A Case Study on Enhancing Functionality with Google PaLM and Large Language Models. *Transactions on Open Source Software Projects* **2024**, *1*.

32. Microsoft. Healthcare agent service Overview. Available online: https://learn.microsoft.com/en-us/azure/health-bot/overview (accessed on January 22 2025).

33. NVIDIA. The Modern AI Data Center. Available online: https://www.nvidia.com/en-au/data-center/?ncid=no-ncid (accessed on January 20 2025).

34. NVIDIA. NVIDIA H100 Tensor Core GPU. Available online: https://www.nvidia.com/en-au/data-center/h100/ (accessed on January 20 2025).

35. Choquette, J. Nvidia hopper gpu: Scaling performance. In Proceedings of the 2022 IEEE Hot Chips 34 Symposium (HCS), 2022; pp. 1-46.

36. Alvarez, R. Why Storage Is the Unsung Hero for AI. Available online: https://blog.purestorage.com/perspectives/why-storage-is-the-unsung-hero-for-ai/ (accessed on January 22 2025).

37. Enclitic. The Importance of Architecture in Healthcare. Available online: https://enlitic.com/blogs/healthcare-architecture/ (accessed on January 15 2025).

38. Carlos Ferreira, J.; Elvas, L.B.; Correia, R.; Mascarenhas, M. Enhancing EHR Interoperability and Security through Distributed Ledger Technology: A Review. *Healthcare (Basel)* **2024**, *12*, doi:10.3390/healthcare12191967.

39. Jackson, G. What is API design? Available online: https://www.ibm.com/think/topics/api-design (accessed on January 20 2025).

40. Reddy, S. Global Harmonization of Artificial Intelligence-Enabled Software as a Medical Device Regulation: Addressing Challenges and Unifying Standards. *Mayo Clinic Proceedings: Digital Health* **2025**, *3*, 100191, doi:https://doi.org/10.1016/j.mcpdig.2024.100191.

41. Reddy, S. Evaluating large language models for use in healthcare: A framework for translational value assessment. *Informatics in Medicine Unlocked* **2023**, *41*, 101304, doi:https://doi.org/10.1016/j.imu.2023.101304.

42. Microsoft. Load Testing RAG based Generative AI Applications. Available online: https://techcommunity.microsoft.com/blog/azure-ai-services-blog/load-testing-rag-based-generative-ai-applications/4086993 (accessed on January 20 2025).

43. Tyler, S.; Olis, M.; Aust, N.; Patel, L.; Simon, L.; Triantafyllidis, C.; Patel, V.; Lee, D.W.; Ginsberg, B.; Ahmad, H.; et al. Use of Artificial Intelligence in Triage in Hospital Emergency Departments: A Scoping Review. *Cureus* **2024**, *16*, e59906, doi:10.7759/cureus.59906.

44. Nguyen, T.L. A framework for five big v's of big data and organizational culture in firms. In Proceedings of the 2018 IEEE international conference on big data (big data), 2018; pp. 5411-5413.

45. Badr, M. Unleashing the power of AI: the Microsoft and OpenAI partnership. **2023**.

46. Dunenfeld, E. Google Gemini vs Azure OpenAI GPT: Pricing Considerations. Available online: https://www.vantage.sh/blog/gcp-google-gemini-vs-azure-openai-gpt-ai-cost (accessed on January 22 2025).

47. Touvron, H.; Martin, L.; Stone, K.R.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. Llama 2: Open Foundation and Fine-Tuned Chat Models. *ArXiv* **2023**, *abs/2307.09288*.

48. Microsoft. Data, privacy, and security for Azure OpenAI Service. Available online: https://learn.microsoft.com/en-us/legal/cognitive-services/openai/data-privacy?tabs=azure-portal (accessed on January 20 2025).

49. Restack. Few-Shot Learning GPT 4 Insights. Available online: https://www.restack.io/p/few-shot-learning-answer-gpt-4-cat-ai (accessed on January 15 2025).

50. Microsoft. Azure Products. Available online: https://azure.microsoft.com/en-au/products (accessed on January 20 2025).

51. Mediicine, J.H. Judy Reitz Capacity Command Center. Available online: https://www.hopkinsmedicine.org/emergency-medicine/c3 (accessed on January 20 2025).

52. Microsoft. Application Insights Overview. Available online: https://learn.microsoft.com/en-us/azure/azure-monitor/app/app-insights-overview (accessed on January 22 2025).

53. Peng, B.; Bi, Z.; Niu, Q.; Liu, M.; Feng, P.; Wang, T.; Yan, L.K.; Wen, Y.; Zhang, Y.; Yin, C.H. Jailbreaking and mitigation of vulnerabilities in large language models. *arXiv preprint arXiv:2410.15236* **2024**.

54. Tonmoy, S.; Zaman, S.; Jain, V.; Rani, A.; Rawte, V.; Chadha, A.; Das, A. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313* **2024**.

55. Sharma, M.; Tong, M.; Korbak, T.; Duvenaud, D.; Askell, A.; Bowman, S.R.; Cheng, N.; Durmus, E.; Hatfield-Dodds, Z.; Johnston, S.R. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548* **2023**.

56. Reddy, S.; Allan, S.; Coghlan, S.; Cooper, P. A governance model for the application of AI in health care. *J Am Med Inform Assoc* **2020**, *27*, 491-497, doi:10.1093/jamia/ocz192.