



Campus Recife

Departamento Acadêmico dos Cursos Superiores

Curso de Tecnologia em Análise e Desenvolvimento de Sistemas

Júlia Didra Bezerra de Assis

**Agente Inteligente Baseado em LLM para Orientação Clínica
em Profilaxia Pós-Exposição a Materiais Biológicos**

Recife

2026

Júlia Didra Bezerra de Assis

**Agente Inteligente Baseado em LLM para Orientação Clínica em Profilaxia
Pós-Exposição a Materiais Biológicos**

Trabalho de Conclusão de Curso apresentado como requisito parcial para obtenção do título de Tecnólogo em Análise e Desenvolvimento de Sistemas, no Instituto Federal de Educação, Ciência e Tecnologia de Pernambuco (IFPE), Campus Recife.

Orientador: Prof. Me. Hilson Gomes Vilar de Andrade

Recife

2026

Ficha elaborada pela bibliotecária Maria do Perpétuo Socorro Cavalcante Fernandes CRB4/1666

A848a
2026

Assis, Júlia Didra Bezerra de

Agente Inteligente baseado em LLM para orientação clínica em profilaxia pós exposição a materiais biológicos / Júlia Didra Bezerra de Assis.---- Recife: A autora, 2026.

40f. il. Color.

Trabalho de Conclusão (Curso Superior de Tecnologia em Análise e Desenvolvimento de sistemas) – Instituto Federal de Pernambuco, Recife, 2026.

Inclui Referências.

Orientador: Prof. Dr. Ms. Hilson Gomes Vilar de Andrade

1. Software- desenvolvimento. 2. LLM. 3. HIV. 4. Profilaxia pós-exposição. 5. Agente Conversacional. 6. RAG. I. Título. II. Andrade, Hilson Gomes Vilar de (orientador). III. Instituto Federal de Pernambuco.

CDD 005.117 (21 ed.)

Júlia Didra Bezerra de Assis

**Agente Inteligente Baseado em LLM para Orientação Clínica em Profilaxia
Pós-Exposição a Materiais Biológicos**

Trabalho de Conclusão de Curso apresentado, como requisito parcial para obtenção do título de Tecnólogo em Análise e Desenvolvimento de Sistemas, do Instituto Federal de Educação, Ciência e Tecnologia de Pernambuco – Campus Recife.

Aprovado em: 09 de março de 2026.

Banca Examinadora:

Prof. Me. Hilson Gomes Vilar de Andrade (Orientador)
IFPE - Campus Recife

Prof^ª Dra. Aida Araújo Ferreira
IFPE - Campus Recife

Prof. Dr. Ricardo Luís Alves da Silva (Externo)
IFPE - Campus Recife

Recife

2025

AGRADECIMENTOS

Agradeço a Deus, em primeiro lugar, que sempre foi o meu alicerce e minha esperança.

À minha família, por não medir esforços para me ajudar, por acreditar em mim e por compreender minhas ausências. Em especial à minha mãe, Andrea Luiza, que, com muita garra, lutou para que eu pudesse chegar até aqui; e à minha irmã, Wilzandrea Bezerra, e ao meu sobrinho, Pedro Júlio, que sempre foram uma das minhas maiores alegrias.

Ao meu noivo, João Pedro, que muito me ensinou e que, com amor e paciência, me deu confiança para seguir em frente e acreditar em mim mesma.

Aos meus colegas de faculdade e amigos para a vida, Gustavo Melo e Raul Vitor, que fizeram parte dessa jornada e tornaram todo o processo mais leve.

Aos professores, deixo meu mais sincero agradecimento e reconhecimento pelos ensinamentos e pela dedicação, vocês são mais importantes do que imaginam. Em especial ao meu orientador, Hilson Gomes, cujo apoio e dedicação foram fundamentais para que eu me mantivesse firme nesse processo; à professora Aida Ferreira, que me deu a primeira oportunidade na área da pesquisa e me ensinou mais do que imagina; e à professora Renata Freire, que, mesmo sem saber, me inspirou a seguir na área de inteligência artificial. Vocês foram fundamentais nessa jornada.

Sê humilde para evitar o orgulho, mas voa alto para alcançar a sabedoria

Agostinho de Hipona

RESUMO

A exposição a materiais biológicos representa risco significativo para o HIV e outras infecções sexualmente transmissíveis, tornando essencial o início oportuno e o manejo adequado da profilaxia pós-exposição (PEP). A efetividade da PEP, contudo, depende de uma avaliação precisa do risco, de orientações claras baseadas em protocolos oficiais e de acompanhamento contínuo — processos frequentemente limitados por restrições nos serviços de saúde e por dificuldades de compreensão por parte dos pacientes. Esses desafios podem resultar em atraso no início da profilaxia ou em baixa adesão ao regime terapêutico de 28 dias. Com o objetivo de enfrentar essa lacuna, este estudo apresenta um Agente conversacional baseado em Modelo de Linguagem de Grande Escala (LLM), integrado a uma arquitetura de Geração Aumentada por Recuperação (Retrieval-Augmented Generation – RAG), desenvolvido para apoiar pacientes em uso de PEP no Brasil. O Agente foi avaliado por meio de um questionário em escala Likert de cinco pontos, aplicado a dezoito médicos especialistas, que analisaram critérios como acurácia clínica, alinhamento com protocolos oficiais, empatia e coerência conversacional. As respostas geradas pelo Agente obtiveram elevada avaliação, com mais de 80% de forte concordância em todas as dimensões analisadas. Na avaliação global, 77,8% dos especialistas classificaram o sistema como de alta qualidade, enquanto 16,7% o consideraram bom, com limitações menores, evidenciando sua viabilidade e potencial aplicabilidade como ferramenta de suporte à orientação clínica.

Palavras-chaves: pep; hiv; llm; agente conversacional; rag.

ABSTRACT

Exposure to biological materials poses significant risks for HIV and other sexually transmitted infections, making the timely initiation and proper management of post-exposure prophylaxis (PEP) essential. Effective PEP, however, depends on accurate risk assessment, clear protocol guidance, and continuous follow-up—processes often limited by healthcare constraints and patient misunderstanding. These challenges may lead to delayed initiation or poor adherence to the 28-day regimen. To address this gap, this study presents a Large Language Model-based conversational Agent integrated with a Retrieval-Augmented Generation (RAG) framework to support patients undergoing PEP in Brazil. The Agent was assessed using a five-point Likert-scale questionnaire applied to eighteen medical specialists, who evaluated clinical accuracy, alignment with official protocols, empathy, and conversational coherence. The Agent's responses were highly rated, with over 80% strong agreement across all evaluated dimensions. In the overall assessment, 77.8% of evaluators rated the system as high quality, while 16.7% considered it good with minor limitations, supporting its feasibility and potential applicability as a clinical guidance-support tool.

Keywords: pep; hiv; llm; conversational agent; rag.

LISTA DE FIGURAS

Figura 1	Fluxograma para indicação de PEP para HIV	20
Figura 2	Relação entre IA, Machine Learning, Deep Learning e IA Generativa .	21
Figura 3	Arquitetura do Agente Inteligente	25
Figura 4	Diagrama de sequência da interação com o agente	26
Figura 5	Avaliação Likert segundo os critérios centrais de avaliação clínica	35
Figura 6	Avaliação geral do Agente	35

LISTA DE TABELAS

Tabela 1	Casos simulados submetidos ao Agente.....	31
Tabela 2	Perguntas Avaliativas do Formulário	32

LISTA DE ABREVIATURAS E SIGLAS

AIDS	<i>Acquired Immune Deficiency Syndrome</i>
API	<i>Application Programming Interface</i>
DL	<i>Deep Learning</i>
GenAI	<i>Generative Artificial Intelligence</i>
HBV	<i>Hepatitis B Virus</i>
HCV	<i>Hepatitis C Virus</i>
HIV	<i>Human Immunodeficiency Virus</i>
IA	Inteligência Artificial
IHC	Interação Humano-Computador
ISTs	Infecções Sexualmente Transmissíveis
LGPD	Lei Geral de Proteção de Dados
LLMs	<i>Large Language Models</i>
ML	<i>Machine Learning</i>
NoSQL	<i>Not Only SQL</i>
PCDT	Protocolo Clínico e Diretrizes Terapêuticas
PEP	<i>Post-exposure prophylaxis</i>
PrEP	Profilaxia Pré-Exposição
RAG	<i>Retrieval Augmented Generation</i>
REST	<i>Representational State Transfer</i>
SBC	Sociedade Brasileira de Computação
SUS	Sistema Único de Saúde

SUMÁRIO

1	INTRODUÇÃO	13
1.1	Objetivos	14
1.1.1	Objetivo Geral	14
1.1.2	Objetivos Específicos	14
1.2	Organização do trabalho	15
2	TRABALHOS RELACIONADOS	16
3	FUNDAMENTAÇÃO TEÓRICA	18
3.1	Protocolo da PEP	18
3.2	Inteligência Artificial Generativa	19
3.3	Assistentes virtuais baseados em GenAI na telemedicina	22
4	METODOLOGIA	24
4.1	Arquitetura do Agente Inteligente Proposto	24
4.1.1	Visão Geral da Arquitetura	24
4.1.2	Camada de Comunicação e Backend	24
4.1.3	Orquestração Conversacional e sistema RAG	26
4.1.4	Segurança e Mecanismos de Controle	26
4.1.5	Armazenamento de Dados e Conformidade com a LGPD	27
4.1.6	Base de Conhecimento e Recuperação Semântica	27
4.1.7	LLM	29
4.1.8	Automação e Suporte à Adesão Terapêutica	29
4.2	Avaliação do Agente	29
5	RESULTADOS	33
5.1	Avaliação por escala Likert	33
5.2	Análise das respostas abertas	33
5.3	Síntese dos resultados	34
6	CONSIDERAÇÕES FINAIS	36
6.1	Principais Contribuições	36
6.2	Limitações e Trabalhos Futuros	36
6.3	Limitações e Trabalhos Futuros	37

REFERÊNCIAS.....

1 INTRODUÇÃO

A exposição a agentes biológicos representa uma preocupação significativa de saúde pública, afetando tanto profissionais de saúde quanto a população em geral, devido ao risco de transmissão de doenças infecciosas. Essa exposição ocorre por meio do contato com materiais biológicos contaminados — como sangue e outros fluidos corporais — seja por acidentes com perfurocortantes ou por contato com mucosas e pele não íntegra. Globalmente, estima-se que ocorram quase 3 milhões de exposições percutâneas por ano entre aproximadamente 35 milhões de profissionais de saúde, sendo o risco de soroconversão após uma única exposição percutânea a patógenos transmitidos pelo sangue estimado em aproximadamente 0.3% para o *Human Immunodeficiency Virus* (HIV), 1.8% para o *Hepatitis C Virus* (HCV) e entre 6% e 30% pelo vírus *Hepatitis B Virus* (HBV) em indivíduos não imunizados (Ministério da Saúde, 2010).

A *Post-exposure prophylaxis* (PEP) para o HIV e outras Infecções Sexualmente Transmissíveis (ISTs) consiste na administração de medicamentos com o objetivo de reduzir o risco de infecção após uma exposição de alto risco (Ministério da Saúde, 2024). O médico desempenha papel central nesse processo, incluindo o aconselhamento do paciente, o esclarecimento de dúvidas, o manejo das prescrições e o acompanhamento clínico (Lins, F. *et al.*, 2023). A PEP é composta por um regime antirretroviral com duração de 28 dias. Como intervenção profilática, a adesão adequada ao tratamento constitui fator determinante para sua eficácia, uma vez que impacta diretamente a prevenção da progressão da doença e o desenvolvimento de possíveis infecções. Nesse contexto, a não adesão à PEP configura um desafio que demanda atenção contínua, sendo a baixa adesão decorrente de múltiplas barreiras (Liu *et al.*, 2023).

A expansão das estratégias de saúde digital, especialmente da telemedicina, tem criado novas oportunidades para apoiar o acompanhamento de pacientes, ao mesmo tempo em que reduz a sobrecarga dos sistemas de saúde. No Brasil, aproximadamente 2,5 milhões de teleconsultas foram registradas em 2024, com a meta de alcançar 10 milhões até 2027 no âmbito da iniciativa do Sistema Único de Saúde (SUS) Digital (Ministério da Saúde, 2025). Esse cenário reforça a relevância de soluções digitais escaláveis para ampliar o acesso, a continuidade e a eficiência na prestação de serviços de saúde.

Paralelamente, a *Generative Artificial Intelligence* (GenAI) tem avançado como um

conjunto de técnicas computacionais capazes de produzir conteúdos semelhantes aos humanos, incluindo textos, imagens e áudios (Feuerriegel *et al.*, 2023). Entre essas técnicas, os *Large Language Models* (LLMs) destacam-se por serem modelos de aprendizado profundo treinados em grandes volumes de dados para compreender e gerar linguagem natural (Naveed *et al.*, 2025). A integração de LLMs à telemedicina possibilita o desenvolvimento de agentes conversacionais capazes de fornecer orientações estruturadas, responder a perguntas e apoiar o manejo terapêutico.

Diante desse contexto, este estudo propõe e avalia um agente conversacional baseado em GenAI, desenvolvido para apoiar pacientes em uso da PEP. O agente utiliza um LLM fundamentado em protocolos médicos estabelecidos para oferecer orientações alinhadas às diretrizes oficiais, esclarecer dúvidas relacionadas ao tratamento, fornecer suporte emocional e reforçar a adesão medicamentosa por meio de lembretes automatizados, com o objetivo de promover a continuidade do cuidado e apoiar os fluxos de decisão clínica.

1.1 Objetivos

1.1.1 Objetivo Geral

O presente estudo tem como objetivo desenvolver e avaliar um agente conversacional baseado em GenAI capaz de apoiar pacientes em tratamento da PEP, fornecendo orientação contínua alinhada a protocolos clínicos, suporte emocional, esclarecimento de dúvidas relacionadas ao tratamento e lembretes automatizados de medicação.

1.1.2 Objetivos Específicos

- Integrar ao agente técnicas de *Retrieval Augmented Generation* (RAG) somadas à LLM, mecanismos de chamada de ferramentas externas e reranking de respostas, garantindo que as informações fornecidas estejam alinhadas aos protocolos oficiais;
- Disponibilizar o agente no WhatsApp, uma plataforma de mensagens amplamente acessível, visando facilidade de uso e ampla adoção;
- Avaliar a efetividade do agente por meio da validação por profissionais de saúde;

1.2 Organização do trabalho

Este trabalho está organizado em seis capítulos. O Capítulo 2 apresenta estudos já existentes sobre o uso de ferramentas de telemedicina com foco em Inteligência Artificial Generativa para orientação e acompanhamento clínico da PEP devido a exposição a materiais biológicos; O Capítulo 3 apresenta o referencial teórico, onde serão abordados os principais assuntos discutidos; no Capítulo 4 descreve-se a metodologia, referindo-se a como foram executados os principais assuntos abordados no referencial teórico; e, por fim, nos Capítulos 5 e 6 são apresentados os resultados e as conclusões finais deste trabalho.

2 TRABALHOS RELACIONADOS

Considerando a aplicação da telemedicina no contexto das profilaxias relacionadas ao HIV, observa-se que a literatura apresenta iniciativas tanto voltadas ao acompanhamento de pacientes expostos a material biológico quanto à prevenção por meio da Profilaxia Pré-Exposição (PrEP). Entretanto, quando se restringe a análise ao suporte estruturado de pacientes em uso da PEP, especialmente com apoio de agentes conversacionais baseados em GenAI, os estudos ainda são limitados.

No que se refere especificamente ao acompanhamento de indivíduos expostos a material biológico, destacam-se dois estudos de Lins *et al.* No primeiro, é proposta a utilização de ferramentas digitais para garantir a continuidade do tratamento durante a pandemia de COVID-19. A solução foi estruturada com base em plataformas amplamente disponíveis e não customizadas, como WhatsApp Business, YouTube e Google Forms, sendo aplicada no acompanhamento de 742 casos de exposição entre julho de 2020 e julho de 2021, na cidade do Recife (Holanda Lins *et al.*, 2023). Os resultados demonstraram a viabilidade do uso de ferramentas digitais de uso geral para apoiar o monitoramento clínico em situações emergenciais, evidenciando o potencial da telemedicina para manutenção do cuidado mesmo em cenários de restrição presencial. No segundo estudo, os autores avaliaram o impacto da telemedicina na prestação do cuidado e nos indicadores de desfecho em um serviço de referência para exposição a materiais biológicos durante a pandemia de COVID-19, no Hospital Correia Picinço, em Recife (Lins, F. H. d. H. *et al.*, 2024).

Paralelamente, no contexto da PrEP, Hoagland *et al.* e Massa *et al.* investigaram estratégias de telemedicina voltadas à prevenção do HIV em grandes centros urbanos brasileiros. Na primeira etapa do estudo, conduzido durante a pandemia de COVID-19, pacientes foram selecionados com base em sintomas relacionados à COVID-19 e à *Acquired Immune Deficiency Syndrome* (AIDS) e acompanhados pelo Instituto Nacional de Infecologia Evandro Chagas (INI-Fiocruz), no Rio de Janeiro, por meio de monitoramento telefônico (Hoagland *et al.*, 2020). Já na segunda etapa do estudo, foi desenvolvido um Agente Inteligente denominado Amanda Selfie, concebido como educador virtual com o objetivo de promover a adesão à PrEP entre adolescentes (Massa *et al.*, 2023). O agente foi implementado na plataforma Facebook Messenger e aplicado nas cidades de Salvador, Belo Horizonte e São Paulo. Embora o agente Amanda Selfie represente um avanço em

termos de usabilidade e escalabilidade, por empregar processamento de linguagem natural em interações automatizadas, sua aplicação esteve direcionada especificamente à PrEP e a um público-alvo delimitado.

Em conjunto, esses estudos demonstram que soluções digitais podem ampliar o acesso à informação e apoiar estratégias profiláticas em diferentes contextos. No entanto, permanece uma lacuna quanto ao uso de agentes conversacionais baseados em GenAI especificamente direcionados ao acompanhamento estruturado de pacientes em uso de PEP. Diante desse cenário, o presente trabalho propõe o desenvolvimento de um agente inteligente voltado a esse público, com foco na oferta de orientação clínica qualificada e no fortalecimento da adesão ao tratamento.

3 FUNDAMENTAÇÃO TEÓRICA

A presente fundamentação teórica tem como objetivo contextualizar os principais conceitos, diretrizes clínicas e referenciais científicos que embasam a proposta deste trabalho. Considerando que a solução proposta situa-se na interface entre saúde e computação, torna-se necessário discutir, de forma articulada, os fundamentos clínicos relacionados à PEP e os aspectos conceituais associados ao uso de GenAI na telemedicina.

Inicialmente, são apresentados os aspectos normativos e clínicos que regem o protocolo da PEP, com ênfase nos critérios de indicação, no fluxo decisório e nos fatores que influenciam a efetividade do tratamento. Em seguida, discute-se o papel dos assistentes virtuais baseados em GenAI no contexto da saúde digital, abordando seus fundamentos conceituais, potenciais benefícios e limitações. Por fim, são analisadas soluções baseadas em telemedicina voltadas à orientação clínica e ao suporte à profilaxia em casos de exposição a material biológico.

3.1 Protocolo da PEP

A PEP consiste na utilização de medicamentos com o objetivo de reduzir o risco de aquisição do HIV, das hepatites virais, da sífilis e de outras ISTs após uma exposição potencialmente de alto risco (Ministério da Saúde, 2024). No contexto da PEP para o HIV, o tratamento deve ser iniciado o mais precocemente possível após a exposição — preferencialmente nas primeiras horas e, no máximo, em até 72 horas — e mantido por meio de um regime antirretroviral com duração de 28 dias (Ministério da Saúde, 2024). Quando corretamente prescrita e integralmente seguida, a PEP demonstra elevada efetividade na redução do risco de infecção pelo HIV (Nunes *et al.*, 2024).

A indicação da PEP depende de uma avaliação clínica estruturada, que considera o tipo de exposição, o tempo decorrido desde o evento e o estado sorológico da pessoa exposta e da pessoa-fonte. A Figura 1 sintetiza o fluxo decisório recomendado pelo Protocolo Clínico e Diretrizes Terapêuticas (PCDT), evidenciando as etapas que orientam a indicação ou não da profilaxia.

Inicialmente, verifica-se se houve exposição a material biológico com potencial risco de transmissão. Em caso negativo, a PEP não é indicada. Quando há exposição relevante — como em situações envolvendo via percutânea, mucosa ou pele não íntegra —, deve-se

avaliar o intervalo entre o evento e o atendimento, sendo a profilaxia recomendada apenas quando iniciada em até 72 horas após a exposição.

Superada essa etapa, procede-se à testagem da pessoa exposta. Caso o teste rápido para HIV seja reagente, a PEP não é indicada, devendo o indivíduo ser encaminhado para investigação diagnóstica confirmatória e início oportuno da Terapia Antirretroviral (TARV). Quando o resultado é não reagente, considera-se o estado sorológico da pessoa-fonte. Se esta apresentar teste positivo ou status desconhecido associado a risco recente, recomenda-se o início imediato da PEP. Por outro lado, quando a pessoa-fonte apresenta teste não reagente e ausência de exposição de risco recente, a profilaxia não é recomendada (Ministério da Saúde, 2024).

Dessa forma, observa-se que a indicação da PEP não se baseia apenas na ocorrência da exposição, mas em uma análise combinada de fatores clínicos e epidemiológicos, conforme representado no fluxo decisório apresentado.

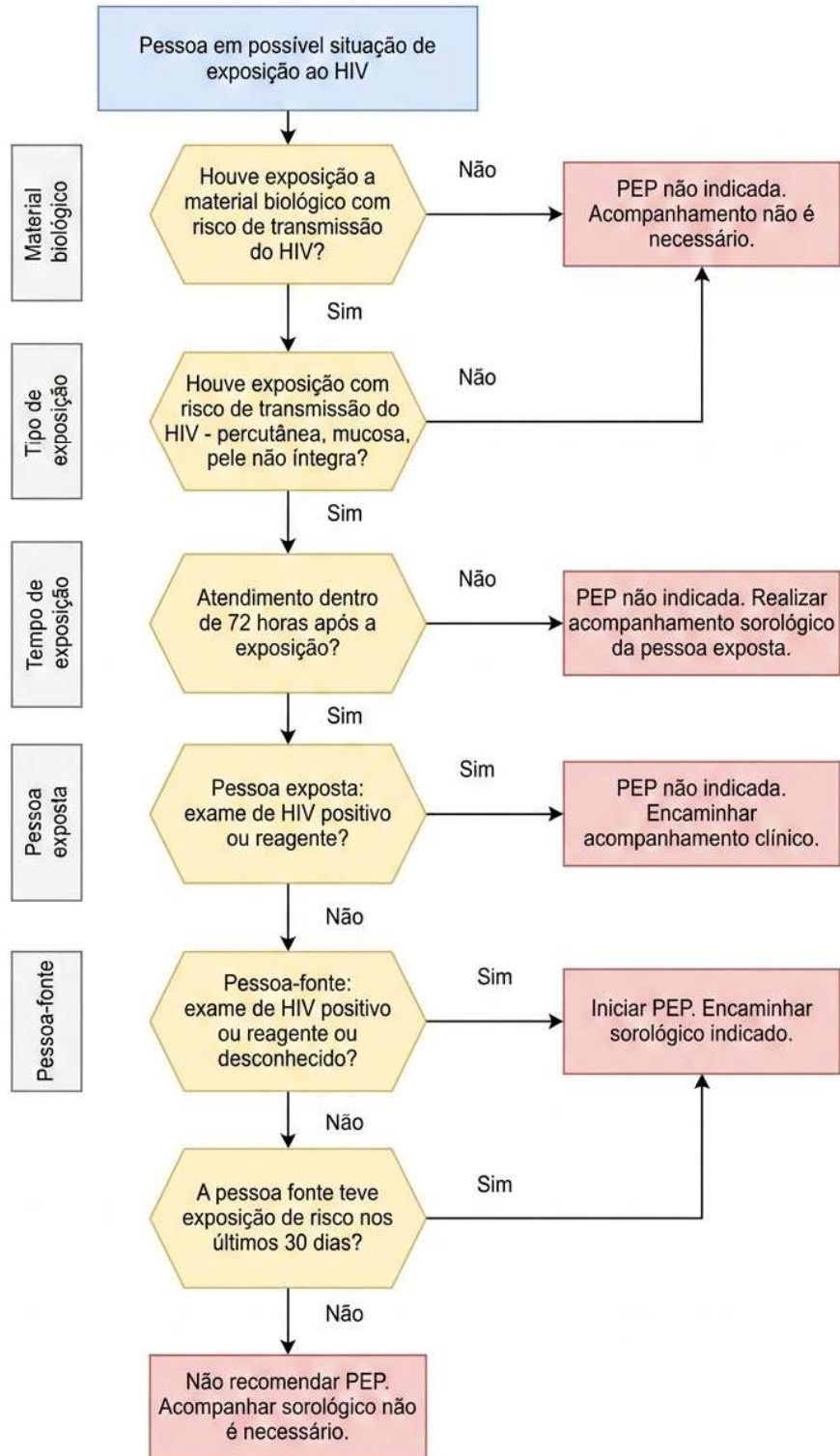
3.2 Inteligência Artificial Generativa

A GenAI representa o estágio mais recente de uma longa trajetória evolutiva no campo da Inteligência Artificial (IA) — área dedicada ao desenvolvimento de sistemas computacionais capazes de simular comportamentos inteligentes, como aprender, raciocinar e tomar decisões de forma autônoma (Chowdhury *et al.*, 2025). Ao longo das últimas décadas, esse campo evoluiu de sistemas baseados em regras fixas para abordagens cada vez mais sofisticadas, impulsionadas pelo avanço do *Machine Learning* (ML) e, posteriormente, do *Deep Learning* (DL).

O ML emergiu como um subconjunto fundamental da IA, caracterizado pela capacidade de algoritmos detectarem padrões em grandes volumes de dados sem depender de programação explícita para cada tarefa. Fatores como a crescente disponibilidade de grandes coleções de textos em formato digital, o aumento do poder computacional e os avanços em algoritmos de aprendizado estatístico impulsionaram essa transição, permitindo que os sistemas aprendessem diretamente a partir dos dados (Chowdhury *et al.*, 2025).

A partir do ML, o DL consolidou-se como uma abordagem especializada, fundamentada em redes neurais artificiais com múltiplas camadas de processamento. Esse paradigma viabilizou avanços expressivos em tarefas como reconhecimento de voz, visão

Figura 1: Fluxograma para indicação de PEP para HIV

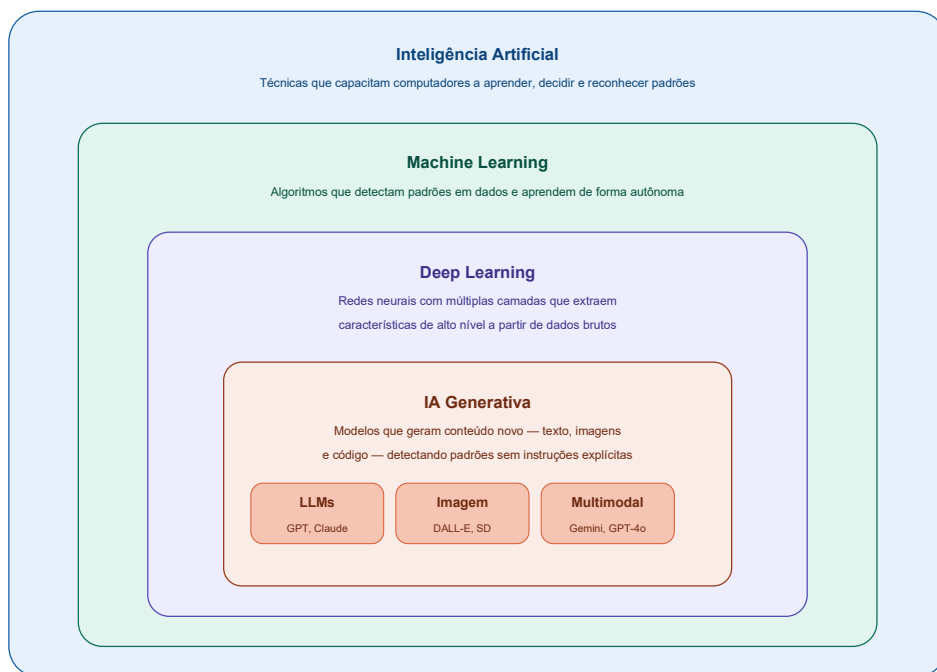


Fonte: Ministério da Saúde (2024)

computacional e compreensão de linguagem natural, ao possibilitar a extração automática de características de alto nível diretamente a partir de dados brutos (Chowdhury *et al.*, 2025).

É nesse contexto que a GenAI se estabelece como um subconjunto do DL, voltado especificamente para a geração autônoma de conteúdo novo — incluindo texto, imagens, áudio e código — com criatividade e complexidade comparáveis às produções humanas (Chowdhury *et al.*, 2025). Conforme ilustrado na Figura 2, essa relação é hierárquica: a GenAI está contida no DL, que por sua vez integra o ML, o qual compõe o campo mais abrangente da IA.

Figura 2: Relação entre IA, Machine Learning, Deep Learning e IA Generativa



Fonte: Autora (2026)

Um marco decisivo para a GenAI foi a introdução da arquitetura *Transformer*, apresentada em 2017 no trabalho *Attention Is All You Need* (Chowdhury *et al.*, 2025). Diferentemente das abordagens anteriores baseadas em redes recorrentes, os *Transformers* utilizam um mecanismo de autoatenção (*self-attention*) que permite capturar dependências de longo alcance em sequências de texto, conferindo maior coerência e fluidez às saídas geradas. Essa arquitetura tornou-se a base dos principais LLMs contemporâneos, como as sucessivas versões do GPT, desenvolvidas pela OpenAI (Chowdhury *et al.*, 2025).

Modelos generativos são pré-treinados em vastos conjuntos de dados e, a partir desse processo, desenvolvem a capacidade de gerar saídas originais sem depender de instruções explícitas para cada situação. Técnicas como o aprendizado por reforço com feedback humano (*Reinforcement Learning from Human Feedback*) foram introduzidas para alinhar o comportamento desses modelos a valores e expectativas humanas, resultando em sistemas mais seguros, úteis e contextualmente precisos (Chowdhury *et al.*, 2025).

Além dos modelos baseados em texto, a GenAI abrange arquiteturas voltadas à geração de imagens bem como modelos multimodais capazes de processar e integrar simultaneamente diferentes tipos de dados, como texto, imagem e áudio (Chowdhury *et al.*, 2025). Essa diversidade de abordagens amplia consideravelmente o espectro de aplicações da GenAI, que atualmente se estende por setores como saúde, educação, indústria e segurança cibernética.

No domínio da saúde, em particular, a GenAI tem sido empregada na construção de assistentes conversacionais, no suporte ao diagnóstico, na geração de resumos clínicos e no auxílio à tomada de decisão (Chowdhury *et al.*, 2025). Sua capacidade de produzir respostas dinâmicas e adaptadas ao contexto do usuário diferencia-a de abordagens baseadas em regras fixas, conferindo maior flexibilidade e potencial de personalização às soluções desenvolvidas. Não obstante, a ausência de mediação humana suscita preocupações quanto à precisão e à confiabilidade das informações geradas, uma vez que esses modelos estão sujeitos a limitações como a ocorrência de alucinações — isto é, a produção de conteúdo incorreto ou sem fundamentação factual (Chowdhury *et al.*, 2025).

3.3 Assistentes virtuais baseados em GenAI na telemedicina

A construção de assistentes virtuais baseados em GenAI fundamenta-se em princípios oriundos dos sistemas de diálogo e da Interação Humano-Computador (IHC), que orientam o design, a implementação e a avaliação de soluções conversacionais centradas no usuário. Esses princípios visam garantir interações eficazes, intuitivas e contextualizadas, promovendo maior engajamento e usabilidade (Oliveira; Oliveira, 2015). Paralelamente, avanços teóricos e tecnológicos no processamento de linguagem natural e de fala têm contribuído para o desenvolvimento de sistemas capazes de compreender e gerar linguagem de forma cada vez mais precisa, fluida e próxima à comunicação humana. (Akpan *et al.*, 2025).

Agentes conversacionais baseados em GenAI têm demonstrado benefícios relevantes nas últimas décadas, incluindo apoio ao diagnóstico, monitoramento clínico e suporte ao tratamento (Schachner; Keller; Wangenheim, 2020). Diferentemente de sistemas tradicionais baseados em regras fixas, os modelos generativos ampliam a capacidade de adaptação a diferentes contextos clínicos, permitindo interações mais dinâmicas e humanizadas.

Além de suas aplicações clínicas diretas, os assistentes virtuais têm sido apontados como ferramentas capazes de fortalecer a relação entre usuários, instituições e profissionais de saúde. Ao possibilitar a abordagem de temas sensíveis sem a necessidade de interação humana imediata, esses sistemas podem reduzir barreiras relacionadas ao constrangimento, ao medo de julgamento e ao estigma social (Massa *et al.*, 2023). Esse aspecto é particularmente relevante em contextos envolvendo infecções sexualmente transmissíveis e outras condições associadas à vulnerabilidade social.

No caso de tratamentos que exigem acompanhamento contínuo, como aqueles relacionados à profilaxia do HIV, o suporte ao paciente ao longo de todo o ciclo terapêutico pode representar um desafio significativo. Dificuldades na manutenção da adesão medicamentosa, dúvidas recorrentes sobre o tratamento e preocupações com possíveis efeitos adversos são fatores que podem comprometer sua efetividade (Holanda Lins *et al.*, 2023). Ademais, aspectos psicossociais — como falta de conhecimento, estigmatização, medo, vergonha e ansiedade — podem impactar negativamente a continuidade do cuidado, sobretudo em situações envolvendo vítimas de violência sexual (Pelegriño; Vioto; Kerche, 2022). Nesse cenário, assistentes virtuais baseados em GenAI podem contribuir como instrumentos de suporte contínuo, oferecendo orientação estruturada e apoio informacional, além de favorecer um ambiente de acolhimento que pode mitigar barreiras emocionais associadas ao tratamento.

Por outro lado, a ausência de mediação humana pode suscitar questionamentos quanto à precisão, à confiabilidade e à segurança das informações fornecidas. Modelos baseados em GenAI estão sujeitos a limitações inerentes, como a ocorrência de alucinações — isto é, a geração de informações incorretas ou não fundamentadas — o que representa um risco relevante no contexto da saúde, onde decisões clínicas exigem alto grau de exatidão. Ainda assim, quando comparados a outras ferramentas digitais, como aplicativos tradicionais de monitoramento, assistentes conversacionais tendem a apresentar maior potencial de engajamento, satisfação e retenção de usuários (Massa *et al.*, 2023).

4 METODOLOGIA

Este capítulo apresenta o delineamento metodológico adotado para o desenvolvimento do presente trabalho. A metodologia foi estruturada de forma a contemplar as etapas de concepção, implementação e avaliação do Agente Inteligente proposto, integrando fundamentos clínicos relacionados à PEP e técnicas contemporâneas de GenAI. Inicialmente, descreve-se a arquitetura e os componentes tecnológicos empregados na construção do sistema. Em seguida, detalham-se os procedimentos adotados para a validação qualitativa do Agente, incluindo a definição dos casos simulados, os instrumentos de avaliação utilizados e os critérios analisados pelos especialistas.

4.1 Arquitetura do Agente Inteligente Proposto

4.1.1 Visão Geral da Arquitetura

Conforme ilustrado na Figura 3, a arquitetura adotada segue o padrão de microserviços, no qual cada componente encapsula uma funcionalidade específica e se comunica por meio de interfaces de rede (Newman, 2021). Esse padrão possibilita a integração de um sistema de mensagens, de um banco de dados estruturado e de um banco de dados vetorial. A fim de garantir isolamento, portabilidade e facilidade de implantação, os serviços foram containerizados com o uso do Docker¹, incluindo tanto a aplicação backend quanto o serviço WAHA².

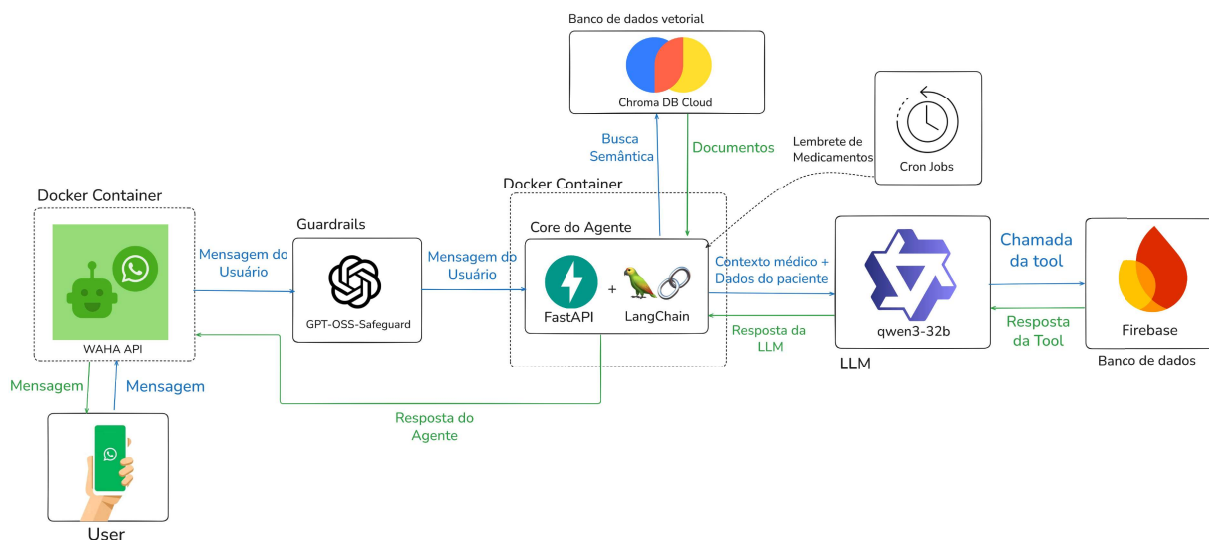
4.1.2 Camada de Comunicação e Backend

Considerando a ampla adoção do WhatsApp no Brasil, local onde se encontra o público-alvo da solução, o WhatsApp foi adotado como interface de interação do usuário com o Agente. Toda a comunicação entre o usuário e o sistema ocorre por meio dessa plataforma: o usuário envia mensagens pelo WhatsApp, e as respostas do Agente são devolvidas no mesmo canal. Para viabilizar essa integração, foi utilizada a *Application Programming Interface* (API) WAHA, que permite a comunicação programática por meio do WhatsApp Web, atuando como ponte entre as mensagens trocadas na plataforma e

¹<https://www.docker.com/>

²<https://waha.devlike.pro/>

Figura 3: Arquitetura do Agente Inteligente



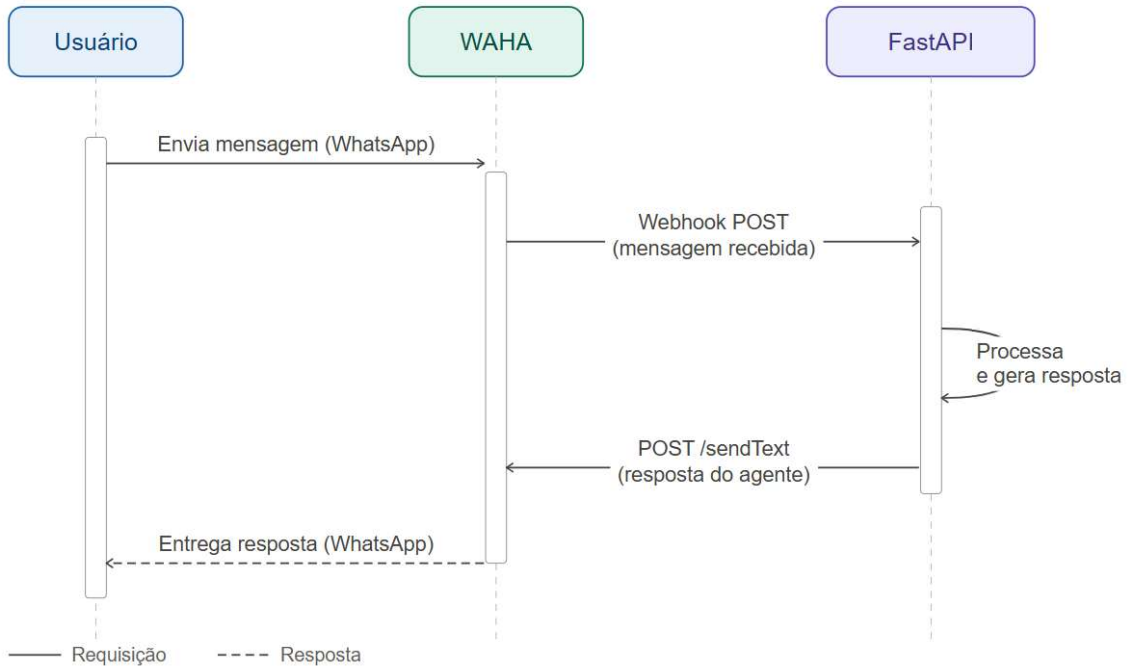
Fonte: Autora (2026)

a lógica do sistema. Ao receber uma nova mensagem, o WAHA a encaminha automaticamente para o backend por meio de um *webhook* — mecanismo pelo qual um serviço notifica outro em tempo real enviando uma requisição HTTP assim que um evento ocorre, eliminando a necessidade de consultas periódicas ao servidor.

Toda a lógica de funcionamento do Agente está centralizada na camada de backend, implementada como uma API REST com o framework FastAPI³. APIs *Representational State Transfer* (REST) são amplamente utilizadas em sistemas distribuídos e baseiam-se em operações padronizadas sobre recursos identificados por URLs (Lubanovic, 2023). É por meio dessa API que todas as operações do sistema são executadas: o recebimento e o processamento das mensagens encaminhadas pelo WAHA via *webhook*, a geração das respostas pelo modelo de linguagem, o gerenciamento do histórico conversacional e a consulta às camadas de armazenamento de dados. Ao concluir o processamento, o FastAPI envia a resposta gerada de volta ao WAHA, que a entrega ao usuário pelo WhatsApp. Dessa forma, o FastAPI concentra a inteligência do sistema, coordenando a comunicação entre todos os componentes da arquitetura. A Figura 4 ilustra esse fluxo de interação.

³<https://fastapi.tiangolo.com/>

Figura 4: Diagrama de sequência da interação com o agente



Fonte: Autora (2026)

4.1.3 Orquestração Conversacional e sistema RAG

As funcionalidades centrais do agente foram implementadas por meio da biblioteca LangChain⁴, que fornece abstrações para encadeamento de chamadas a LLMs, integração de ferramentas externas e implementação de estratégias como RAG (Oshin; Campos, 2025).

O LangChain foi utilizado para organizar o fluxo conversacional, integrar dinamicamente dados clínicos recuperados no RAG e coordenar a execução de ferramentas especializadas, permitindo que as respostas do modelo fossem fundamentadas em conteúdos previamente validados.

4.1.4 Segurança e Mecanismos de Controle

Para assegurar segurança e alinhamento ao domínio clínico, foram implementados mecanismos de *guardrails* com o objetivo de controlar o comportamento do LLM e prevenir respostas inadequadas ou potencialmente inseguras (Ayyamperumal; Ge, 2024). De

⁴<https://www.langchain.com/>

maneira geral, *guardrails* consistem em mecanismos de controle e validação aplicados ao fluxo de interação com o modelo, estabelecendo restrições quanto ao escopo temático e ao tipo de conteúdo processado.

No contexto deste trabalho, os *guardrails* foram implementados exclusivamente na etapa de entrada do sistema, atuando como um filtro prévio às requisições encaminhadas ao LLM. Esse mecanismo realiza a análise das mensagens do usuário com o objetivo de identificar conteúdos maliciosos, tentativas de *prompt injection*, solicitações fora do domínio clínico ou mensagens potencialmente perigosas. Caso seja detectada alguma violação das regras definidas, a requisição não é encaminhada ao modelo. Essa estratégia reduz riscos de uso indevido do sistema e contribui para manter a interação estritamente restrita ao escopo da PEP.

Para essa finalidade, o modelo *GPT-oss-safeguard* foi empregado como componente de validação da entrada, reforçando a confiabilidade e a segurança das interações no contexto médico.

4.1.5 Armazenamento de Dados e Conformidade com a LGPD

O armazenamento de dados foi realizado utilizando o Firebase⁵, um banco de dados *Not Only SQL* (NoSQL) orientado a documentos, adotado por sua facilidade de configuração e por oferecer infraestrutura nativamente hospedada na nuvem, dispensando a necessidade de gerenciamento de servidores.

Em conformidade com a Lei Geral de Proteção de Dados (LGPD) — Lei nº 13.709/2018, que regula o tratamento de dados pessoais no Brasil com o objetivo de garantir a privacidade e a proteção dos indivíduos —, o sistema não armazena informações pessoalmente identificáveis. São registrados apenas dados estritamente necessários ao acompanhamento terapêutico, incluindo identificador do chat, datas de início e término do tratamento, medicamento prescrito, status do tratamento e tipo de profilaxia.

4.1.6 Base de Conhecimento e Recuperação Semântica

A base de conhecimento utilizada pelo sistema foi composta por duas fontes principais: (i) o PCDT para PEP e (ii) um conjunto de dados em formato de perguntas e respostas, produzido e validado por profissionais de saúde, contendo 84 pares de questões

⁵<https://firebase.google.com/>

e respostas. Essas fontes foram selecionadas com o objetivo de assegurar alinhamento às diretrizes clínicas nacionais e contemplar dúvidas recorrentes dos usuários.

No contexto da recuperação semântica, os documentos textuais precisam ser convertidos em representações numéricas denominadas *embeddings*. *Embeddings* são vetores em um espaço de alta dimensionalidade que capturam propriedades semânticas do texto, permitindo representar computacionalmente relações de significado entre diferentes trechos textuais (Osinga, 2018).

Para viabilizar essa etapa, os documentos médicos foram segmentados em unidades textuais menores (trechos ou segmentos), procedimento necessário para permitir indexação mais precisa e recuperação contextual mais específica. Para essa finalidade, foi utilizada a biblioteca GROBID⁶, que possibilita a extração estruturada e o processamento de textos científicos. Esse procedimento permitiu converter os conteúdos do PCDT e do conjunto de perguntas e respostas em segmentos textuais coerentes, preservando a organização documental e favorecendo maior granularidade na etapa de indexação.

Em seguida, os segmentos textuais foram vetorizados utilizando o modelo all-MiniLM-L6-v2 da HuggingFace⁷, gerando os *embeddings* correspondentes. Essas representações permitem mapear semanticamente o conteúdo textual, capturando relações de significado entre consultas e documentos e viabilizando a busca baseada em similaridade.

Para armazenar e indexar essas representações vetoriais, foi adotado o Chroma DB como banco de dados vetorial. Bancos vetoriais são sistemas especializados no armazenamento e na recuperação eficiente de embeddings, possibilitando a chamada recuperação semântica (Ma *et al.*, 2023). A recuperação semântica consiste na busca de informações com base na similaridade de significado entre textos, em vez de depender exclusivamente da correspondência literal de palavras-chave.

O Chroma DB foi configurado com similaridade do cosseno como métrica de comparação vetorial e com o algoritmo HNSW (Hierarchical Navigable Small World) para busca aproximada de vizinhos mais próximos (Lee *et al.*, 2025). Essa configuração permite identificar, de forma escalável e com elevado desempenho de *recall*, os segmentos semanticamente mais próximos à consulta do usuário.

Após a etapa inicial de recuperação, foi aplicada uma estratégia de *reranking* com

⁶<https://github.com/kermitt2/grobid>

⁷<https://huggingface.co/>

o modelo rerank-multilingual-v3.0 da Cohere⁸, que reordena os documentos recuperados e prioriza os trechos semanticamente mais relevantes. O conjunto final de segmentos selecionados compõe o contexto fornecido ao LLM no processo de geração de respostas (Abdallah *et al.*, 2025).

4.1.7 LLM

O LLM selecionado foi o modelo Qwen3-32B, que ocupa atualmente a segunda posição no ranking Arena Hard, benchmark projetado para avaliar o desempenho de modelos de linguagem (Li *et al.*, 2024). A escolha do modelo fundamentou-se em seu desempenho consistente em cenários com ampla janela de contexto e natureza de código aberto, características particularmente relevantes para aplicações clínicas que demandam processamento confiável de informações extensas e integração com ferramentas externas.

4.1.8 Automação e Suporte à Adesão Terapêutica

Por fim, foram configurados *cron jobs* — tarefas agendadas executadas em intervalos de tempo predefinidos (Foote, 2014). Na proposta apresentada, implementou-se um *cron job* diário responsável pelo envio automático de lembretes de medicação a usuários com tratamento ativo. Esse mecanismo visa apoiar a adesão terapêutica e automatizar parte do acompanhamento clínico realizado pelo agente.

4.2 Avaliação do Agente

A avaliação do Agente Inteligente foi conduzida no contexto de uma pesquisa de abordagem qualitativa, com caráter exploratório e descritivo. Como procedimento metodológico, foi empregada a técnica de entrevista com um médico especialista, cujo objetivo foi subsidiar a elaboração de um conjunto de perguntas representativas das principais dúvidas apresentadas por indivíduos expostos a material biológico no início do atendimento ambulatorial.

A partir dessa etapa, foi construído um conjunto de dezoito perguntas, estruturadas como casos simulados, entendidos neste estudo como interações representativas de

⁸<https://cohere.com/>

situações reais de atendimento. Esses casos foram utilizados como base para a avaliação do desempenho do Agente.

O processo de avaliação foi conduzido em duas etapas distintas. Na primeira etapa, cada um dos dezoito casos simulados foi submetido diretamente ao Agente pela própria pesquisadora, e as respostas geradas foram coletadas individualmente. Na segunda etapa, essas respostas foram organizadas em um formulário estruturado de avaliação, desenvolvido especificamente para este estudo. O formulário continha o enunciado de cada caso simulado acompanhado da respectiva resposta gerada pelo Agente.

Esse formulário foi então encaminhado a profissionais de saúde, que atuaram como avaliadores especialistas. Ressalta-se que os avaliadores não interagiram diretamente com o Agente em nenhum momento, limitando-se à análise dos pares (caso simulado, resposta do Agente) previamente registrados no instrumento de avaliação. O formulário foi disponibilizado mediante aceite do Termo de Consentimento Livre e Esclarecido (TCLE), garantindo que os participantes estivessem cientes dos objetivos da pesquisa e concordassem com sua participação.

A adequação de cada resposta do Agente em relação ao respectivo caso simulado foi avaliada por meio de uma escala Likert de cinco pontos, instrumento amplamente utilizado para mensurar percepções e níveis graduais de concordância ou avaliação, organizados em categorias ordinais que variam de avaliações muito negativas a avaliações muito positivas (Joshi *et al.*, 2015). Neste estudo, a escala variou de 1 a 5, em que 1 indicava uma resposta completamente inadequada e 5 indicava uma resposta totalmente adequada.

A Tabela 1 apresenta a distribuição dos casos simulados submetidos ao Agente. Cada uma das respostas fornecidas pelo Agente foi avaliada pelos especialistas segundo quatro aspectos: qualidade da resposta, aderência ao PCDT, adequação da conduta recomendada para o paciente e presença de empatia e acolhimento na comunicação. Adicionalmente, foi incluído um campo aberto para que os avaliadores pudessem registrar comentários, sugestões e críticas.

Além da avaliação dos 18 casos simulados, foram incluídas duas questões adicionais ao final do formulário: uma questão fechada, avaliada por meio da escala Likert ("Qual é a sua avaliação geral do Agente Inteligente?"), e uma questão aberta ("Você acredita que o Agente é adequado para uso em Unidades de Saúde? Por quê?"). A Tabela 2 apresenta o resumo dos itens do questionário submetido aos avaliadores, com o objetivo

de mensurar a qualidade das respostas geradas pelo Agente.

Tabela 1: Casos simulados submetidos ao Agente

Nº	Caso Simulado
1	Encostou sangue na minha pele intacta. Preciso tomar PEP?
2	Sangue caiu em um machucado pequeno no meu braço. Ainda conta como risco?
3	Fui mordido, mas não saiu sangue. Preciso de PEP?
4	Cheguei ao hospital 2 horas depois da exposição. Isso melhora a eficácia?
5	Já se passaram 73 horas desde a exposição. Ainda vale iniciar a PEP?
6	Qual é o esquema padrão de PEP para HIV no SUS?
7	A PEP dura quantos dias?
8	A PEP muda se foi sexo, perfurocortante ou respingo em mucosa?
9	Se eu começar a PEP, isso significa que já peguei HIV?
10	A PEP funciona 100%?
11	Sou técnica de enfermagem. Me furei com uma agulha usada em um paciente durante a coleta de sangue. O acidente ocorreu há 1 hora. Não sei o status sorológico do paciente. Preciso iniciar PEP agora?
12	Tive relação sexual sem camisinha com uma pessoa que não conheço bem. Só consegui procurar ajuda depois de 4 dias. Mesmo assim posso tomar a PEP?
13	Tive relação sem camisinha com meu parceiro fixo, que vive com HIV, mas está em tratamento e com carga viral indetectável há mais de 6 meses. Preciso tomar PEP?
14	Sou estudante de odontologia, nunca me vacinei contra hepatite B. Me furei com um instrumental usado há poucas horas. O que precisa ser feito além da PEP para HIV?
15	Sou uma mulher de 22 anos, sofri violência sexual há cerca de 10 horas. Quais cuidados preciso receber agora?
16	Turno 1: Estou com muito medo, comecei a PEP ontem. Turno 2: Hoje acordei enjoado. Posso parar de tomar? Turno 3: Se eu parar agora, ainda estou protegido?
17	Turno 1: Esqueci de tomar um comprimido do meu tratamento da PEP ontem. Turno 2: Então vou tomar dois comprimidos hoje, certo?
18	Turno 1: Acho que alguém me espetou no metrô, mas não vi agulha. Turno 2: O médico disse que eu podia escolher tomar a PEP ou não. Isso é normal? Turno 3: Tenho muito medo, mesmo achando que o risco é baixo.

Fonte: Autora (2026)

Tabela 2: Perguntas Avaliativas do Formulário

ID	Perguntas	Escala Likert
P1	Avalie se a resposta do agente está de acordo com a pergunta	1 – Resposta totalmente incorreta 5 – Resposta Totalmente correta
P2	A resposta está de acordo com o PCDT e protocolos oficiais?	1 – Não, totalmente em desacordo 5 – Sim, totalmente de acordo
P3	A conduta indicada é apropriada para o caso descrito?	1 – Inadequada 5 – Totalmente adequada
P4	A resposta demonstra empatia e acolhimento adequados?	1 – Inadequado 5 – Excelente
P5	Campo livre, fique a vontade para fazer comentários ou sugerir melhorias	N/A
P6	Considerando todos os aspectos avaliados, qual sua avaliação geral para o agente inteligente?	1 – Inadequado para uso clínico 5 – Alta qualidade
P7	Você acha que o agente é confiável para ser usado em unidades de saúde? Porque?	N/A

Fonte: Autora (2026)

5 RESULTADOS

Para validar a viabilidade do uso do Agente como ferramenta de apoio a pacientes durante o tratamento da PEP, esta avaliação contou com a participação de 18 profissionais de saúde, obtidos por meio da divulgação de um formulário estruturado de avaliação em canais de comunicação dos profissionais da área. Esse quantitativo corresponde ao número de avaliadores que responderam ao instrumento no período de coleta de dados, realizado ao longo de aproximadamente dez dias.

Apresentam-se, a seguir, os resultados da avaliação conduzida por esses especialistas da área. Os achados estão organizados conforme a estrutura do instrumento de avaliação, incluindo tanto os itens em escala Likert quanto as respostas abertas. As médias das pontuações dos itens em escala Likert, correspondentes aos principais critérios de avaliação clínica do Agente proposto, estão ilustradas na Figura 5. Esses itens avaliam a aderência ao PCDT, a correção clínica, a adequação da conduta recomendada e o nível de empatia demonstrado nas respostas.

5.1 Avaliação por escala Likert

A Figura 6 apresenta a distribuição das respostas em escala Likert para a questão Q6, que avaliou a percepção global sobre o Agente Inteligente. Entre os 18 profissionais de saúde participantes, 14 atribuíram nota 5, três atribuíram nota 4 e um atribuiu nota 3. Nenhuma pontuação inferior a 3 foi registrada.

5.2 Análise das respostas abertas

Nove respostas abertas foram submetidas para a questão Q7 (“Você acredita que o Agente é adequado para uso em Unidades de Saúde? Por quê?”). Os comentários foram agrupados de acordo com a pontuação atribuída em Q6. Entre os respondentes que atribuíram nota 5, os comentários destacaram a confiabilidade do Agente, sua aplicabilidade em diferentes pontos de atenção à saúde e sua utilidade para esclarecer dúvidas relacionadas a tratamentos preventivos e pós-exposição. Houve ainda menções à relevância dos protocolos estabelecidos, como PEP, PrEP e tratamento do HIV, bem como à sensibilidade do tema abordado. Um dos comentários ressaltou que, apesar de

fornecer orientações práticas, certas decisões exigem encaminhamento a um profissional de saúde ou serviço especializado, especialmente devido à impossibilidade de realizar exame físico ou clínico.

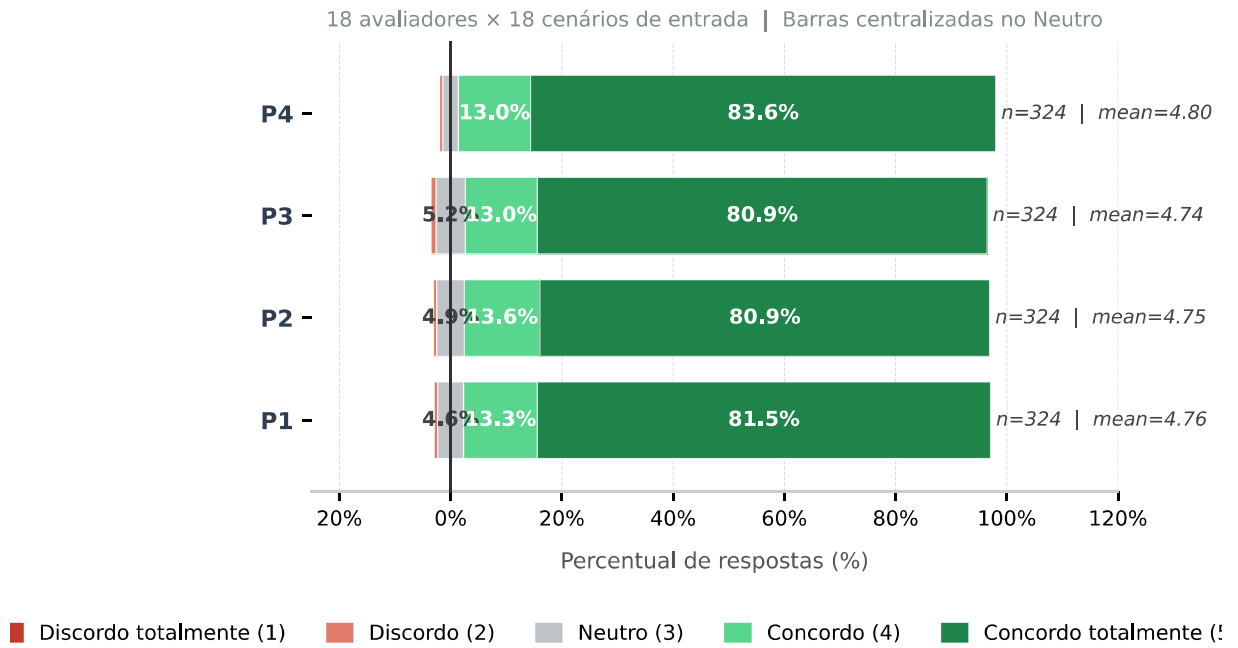
Entre aqueles que atribuíram nota 4, os comentários indicaram que o Agente seria valioso com pequenos ajustes e enfatizaram a importância de confirmar informações diretamente com um profissional de saúde quando necessário. O respondente que atribuiu nota 3 mencionou que o questionário era muito longo e continha textos extensos.

5.3 Síntese dos resultados

Os resultados demonstram elevado nível de aceitação das respostas do Agente entre os especialistas médicos, com mais de 80% de concordância plena em todas as quatro dimensões avaliadas: alinhamento com a questão clínica (P1), aderência aos protocolos clínicos estabelecidos (P2), adequação da conduta recomendada com base no caso descrito (P3) e demonstração de empatia e comunicação centrada no paciente (P4). A consistência nas altas pontuações tanto em domínios técnicos quanto relacionais reforça a confiabilidade e a coerência clínica das respostas geradas.

Quando avaliado sob uma perspectiva mais ampla (P6) — incluindo usabilidade do sistema e desempenho geral —, o Agente também obteve avaliações altamente favoráveis. A maioria substancial de especialistas (77,8%) classificou o sistema como de alta qualidade geral, enquanto 16,7% o consideraram bom, embora com pequenas limitações. Esses achados fornecem evidências preliminares robustas que sustentam a viabilidade, relevância clínica e aplicabilidade prática do Agente como ferramenta de apoio durante o tratamento da PEP.

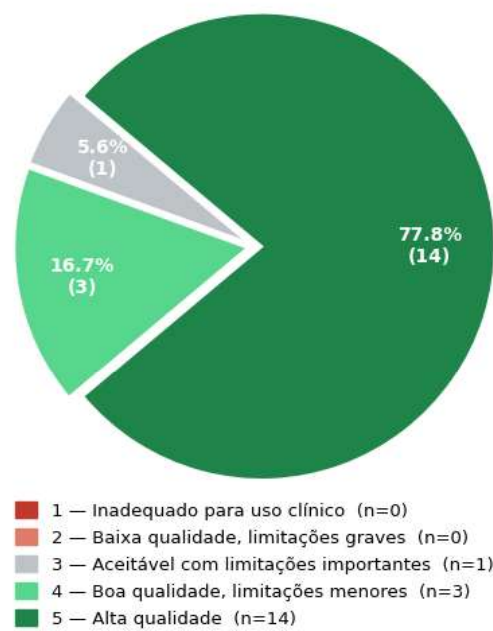
Figura 5: Avaliação Likert segundo os critérios centrais de avaliação clínica



Fonte: Autora (2026)

Figura 6: Avaliação geral do Agente

P6 — Considerando todos os aspectos avaliados, qual sua avaliação geral para o agente inteligente?



Fonte: Autora (2026)

6 CONSIDERAÇÕES FINAIS

Esta seção apresenta uma síntese das principais contribuições do estudo, destacando os resultados obtidos com a avaliação do Agente , bem como as limitações identificadas e as perspectivas para trabalhos futuros.

6.1 Principais Contribuições

As elevadas taxas de aprovação observadas em todas as dimensões de avaliação indicam que o Agente Inteligente proposto foi percebido como clinicamente confiável e contextualmente adequado pelos especialistas médicos. As pontuações consistentemente altas em alinhamento, aderência a diretrizes, adequação de conduta juntamente com empatia e acolhimento sugerem que a arquitetura baseada em RAG em conjunto com a LLM foi eficaz em fundamentar as respostas do agente nos protocolos clínicos oficiais, traduzindo o conhecimento médico estruturado em orientações coerentes e aplicáveis.

As avaliações gerais predominantemente positivas reforçam a viabilidade de implantação de tal sistema em ambientes de saúde reais, particularmente no contexto da PEP para HIV, nos quais informações precisas e oportunas são críticas. Os comentários abertos indicam ainda que os especialistas reconhecem o agente não como substituto do julgamento clínico, mas como uma ferramenta valiosa de suporte à tomada de decisão — especialmente em cenários nos quais o acesso a um especialista ou a um centro especializado possa ser limitado.

De modo geral, os resultados sustentam a conclusão de que a combinação de RAG e LLM com conhecimento clínico específico de domínio pode gerar um sistema capaz de fornecer suporte significativo e alinhado aos protocolos em contextos médicos.

6.2 Limitações e Trabalhos Futuros

O presente estudo apresenta algumas limitações. O conjunto de dados utilizado na avaliação foi relativamente pequeno devido à necessidade de validação conduzida por especialistas, o que pode restringir a generalização estatística dos achados. A expansão do conjunto de testes em investigações futuras permitirá uma avaliação mais robusta e conclusões estatísticas mais sólidas.

Trabalhos futuros concentrar-se-ão na ampliação do conjunto de avaliação, tanto em termos de cenários clínicos quanto do número de especialistas participantes, a fim de viabilizar conclusões estatísticas mais consistentes. Do ponto de vista técnico, melhorias planejadas incluem otimizações no componente de busca semântica — como estratégias refinadas de seleção de fragmentos e técnicas de redução de contexto —, bem como avaliação sistemática da qualidade das respostas geradas por meio de métricas automatizadas. Além disso, pretende-se avaliar outros componentes arquiteturais do Agente, incluindo uma análise estruturada de suas capacidades de chamada de ferramentas externas e da efetividade dos mecanismos de segurança (*guardrails*) na garantia de confiabilidade e segurança do sistema. A validação controlada em ambientes reais de saúde permanece um objetivo de longo prazo.

6.3 Limitações e Trabalhos Futuros

O presente estudo apresenta algumas limitações. O conjunto de dados utilizado na avaliação foi relativamente pequeno devido à necessidade de validação conduzida por especialistas, o que pode restringir a generalização estatística dos achados. A expansão do conjunto de testes em investigações futuras permitirá uma avaliação mais robusta e conclusões estatísticas mais sólidas.

Adicionalmente, ressalta-se que não houve restrição quanto à especialidade dos profissionais participantes, de modo que a amostra foi composta por profissionais de saúde de diferentes áreas de atuação, não se limitando a médicos especialistas diretamente relacionados ao contexto da PEP, como infectologistas ou médicos do trabalho. Esse aspecto pode influenciar a interpretação dos resultados, considerando possíveis variações no nível de familiaridade com os protocolos clínicos avaliados.

Trabalhos futuros concentrar-se-ão na ampliação do conjunto de avaliação, tanto em termos de cenários clínicos quanto do número de especialistas participantes, a fim de viabilizar conclusões estatísticas mais consistentes. Adicionalmente, pretende-se aprimorar o processo de caracterização dos participantes, incluindo o mapeamento das especialidades dos profissionais de saúde e de seu tempo de experiência, de modo a possibilitar a estratificação e filtragem das avaliações conforme a área de atuação e a vivência clínica, especialmente em relação a especialistas diretamente envolvidos no contexto da PEP.

Do ponto de vista técnico, melhorias planejadas incluem otimizações no compo-

nente de busca semântica — como estratégias refinadas de seleção de fragmentos e técnicas de redução de contexto —, bem como avaliação sistemática da qualidade das respostas geradas por meio de métricas automatizadas. Além disso, pretende-se avaliar outros componentes arquiteturais do Agente, incluindo uma análise estruturada de suas capacidades de chamada de ferramentas externas e da efetividade dos mecanismos de segurança (*guardrails*) na garantia de confiabilidade e segurança do sistema. A validação controlada em ambientes reais de saúde permanece um objetivo de longo prazo.

REFERÊNCIAS

ABDALLAH, Abdelrahman *et al.* **Rankify: A Comprehensive Python Toolkit for Retrieval, Re-Ranking, and Retrieval-Augmented Generation.** [*S. l.*]: arXiv, 2025. DOI: 10.48550/ARXIV.2502.02464. Disponível em: <https://arxiv.org/abs/2502.02464>.

AKPAN, Ikpe Justice *et al.* Conversational and generative artificial intelligence and human–chatbot interaction in education and research. **International Transactions in Operational Research**, 2025.

AYYAMPERUMAL, Suriya Ganesh; GE, Limin. **Current state of LLM Risks and AI Guardrails.** [*S. l.*]: arXiv, 2024. DOI: 10.48550/ARXIV.2406.12934. Disponível em: <https://arxiv.org/abs/2406.12934>.

CHOWDHURY, Abdullahi *et al.* Generative AI: A survey of historical development, emerging trends, and future outlook. **Computer Science and Engineering Research**, Genesis Publishing Consortium Limited, v. 2, n. 1, p. 19–31, ago. 2025. DOI: 10.69517/cser.2025.02.01.0004. Disponível em: <https://doi.org/10.69517/cser.2025.02.01.0004>.

FEUERRIEGEL, Stefan *et al.* Generative AI. **Business & Information Systems Engineering**, Springer Science e Business Media LLC, v. 66, n. 1, p. 111–126, set. 2023. ISSN 1867-0202. DOI: 10.1007/s12599-023-00834-7. Disponível em: <https://doi.org/10.1007/s12599-023-00834-7>.

FOOTE, S. **Learning to Program.** [*S. l.*]: Pearson Education, 2014. (Learning). ISBN 9780133795226. Disponível em: <https://books.google.com.br/books?id=XHnbBAAAQBAJ>.

HOAGLAND, Brenda *et al.* Telemedicine as a tool for PrEP delivery during the COVID-19 pandemic in a large HIV prevention service in Rio de Janeiro-Brazil. **Brazilian Journal of Infectious Diseases**, SciELO Brasil, v. 24, n. 4, p. 360–364, 2020.

HOLANDA LINS, Flávio Henrique de *et al.* Produto mínimo viável de telemedicina para profilaxia pós-exposição a material biológico na pandemia da covid-19. **Rev Bras Med**, v. 2, p. 9, 2023.

JOSHI, Ankur *et al.* Likert Scale: Explored and Explained. **British Journal of Applied Science & Technology**, Sciencedomain International, v. 7, n. 4, p. 396–403, jan. 2015. ISSN 2231-0843. DOI: 10.9734/bjast/2015/14975. Disponível em: <https://doi.org/10.9734/bjast/2015/14975>.

LEE, Haena *et al.* P-HNSW: Crash-Consistent HNSW for Vector Databases on Persistent Memory. **Applied Sciences**, MDPI AG, v. 15, n. 19, p. 10554, set. 2025.

ISSN 2076-3417. DOI: 10.3390/app151910554. Disponível em:
<https://doi.org/10.3390/app151910554>.

LI, Tianle *et al.* **From Crowdsourced Data to High-Quality Benchmarks: Arena-Hard and BenchBuilder Pipeline.** [*S. l.*]: arXiv, 2024. DOI: 10.48550/ARXIV.2406.11939. Disponível em: <https://arxiv.org/abs/2406.11939>.

LINS, Flávio *et al.* Produto mínimo viável de telemedicina para profilaxia pós-exposição a material biológico na pandemia da covid-19. **Revista Brasileira de Medicina do Trabalho**, 2023.

LINS, Flávio Henrique de Holanda *et al.* Telemedicine in Post-Exposure Prophylaxis to Biological Material During the COVID-19 Pandemic: Impact on Care and Outcome Indicators. **Telemedicine and e-Health**, SAGE Publications Sage CA: Los Angeles, CA, v. 30, n. 9, p. 2445–2455, 2024.

LIU, Shanshan *et al.* Adherence, adverse drug reactions, and discontinuation associated with adverse drug reactions of HIV post-exposure prophylaxis: a meta-analysis based on cohort studies. **Annals of Medicine**, Informa UK Limited, v. 55, n. 2, dez. 2023. ISSN 1365-2060. DOI: 10.1080/07853890.2023.2288309. Disponível em: <https://doi.org/10.1080/07853890.2023.2288309>.

LUBANOVIC, B. **FastAPI: Modern Python Web Development.** [*S. l.*]: O'Reilly Media, 2023. ISBN 9781098135461. Disponível em: <https://books.google.com.br/books?id=XJHhEAAAQBAJ>.

MA, Le *et al.* **A Comprehensive Survey on Vector Database: Storage and Retrieval Technique, Challenge.** [*S. l.*]: arXiv, 2023. DOI: 10.48550/ARXIV.2310.11703. Disponível em: <https://arxiv.org/abs/2310.11703>.

MASSA, Paula *et al.* A Transgender Chatbot (Amanda Selfie) to Create Pre-exposure Prophylaxis Demand Among Adolescents in Brazil: Assessment of Acceptability, Functionality, Usability, and Results. **Journal of Medical Internet Research**, JMIR Publications Inc., v. 25, e41881, jun. 2023. ISSN 1438-8871. DOI: 10.2196/41881. Disponível em: <https://doi.org/10.2196/41881>.

MINISTÉRIO DA SAÚDE. **Manual A B C D E das Hepatites Virais para Cirurgiões Dentistas.** Brasília: Ministério da Saúde, 2010.

MINISTÉRIO DA SAÚDE. **Ministério da Saúde entrega 3 mil kits de telessaúde e lança edital inédito para expandir atendimento a distância com hospitais privados.** [*S. l.*: *s. n.*], 2025. Disponível em: <https://www.gov.br/saude/pt-br/assuntos/noticias/2025/agosto/ministerio-da-saude-entrega-3-mil-kits-de-telessaude-e-lanca-edital-inedito-para-expandir-atendimento-a-distancia-com-hospitais-privados>. Acesso em: 28 jan. 2026.

MINISTÉRIO DA SAÚDE. **Protocolo Clínico e Diretrizes Terapêuticas para Profilaxia Pós-Exposição (PEP) de Risco à Infecção pelo HIV, IST e Hepatites Virais.** [S. l.]: Ministério da Saúde, 2024. p. 80. Disponível em: https://www.gov.br/aids/pt-br/central-de-conteudo/pcdts/2021/hiv-aids/prot_clinico_diretrizes_terap_pep_risco_infeccao_hiv_ist_hv_2021.pdf/view. Acesso em: 20 jan. 2025.

NAVEED, Humza *et al.* A Comprehensive Overview of Large Language Models. **ACM Trans. Intell. Syst. Technol.**, Association for Computing Machinery, New York, NY, USA, v. 16, n. 5, ago. 2025. ISSN 2157-6904. DOI: 10.1145/3744746. Disponível em: <https://doi.org/10.1145/3744746>.

NEWMAN, S. **Building Microservices.** [S. l.]: O'Reilly Media, 2021. ISBN 9781492033998. Disponível em: <https://books.google.com.br/books?id=aPM5EAAAQBAJ>.

NUNES, Danielly Brito Guidobono *et al.* A DESMITIFICAÇÃO NA UTILIZAÇÃO DA PROFILAXIA PÓS-EXPOSIÇÃO (PEP) E OS ENTRAVES MEDIANTE CONSTÂNCIA PELO MEDICAMENTO. **Revista Políticas Públicas & Cidades**, Editoriais Iberoamericanos, v. 13, n. 2, e1060, out. 2024. ISSN 2359-1552. DOI: 10.23900/2359-1552v13n2-163-2024. Disponível em: <https://doi.org/10.23900/2359-1552v13n2-163-2024>.

OLIVEIRA, Francisco Carlos de Mattos Brito;
OLIVEIRA, Fernando Antonio de Mattos Brito. **Interação Humano Computador.** 2. ed. Fortaleza, CE: EdUECE, 2015. (Computação). ISBN 978-85-7826-565-6.

OSHIN, M.; CAMPOS, N. **Learning LangChain.** [S. l.]: O'Reilly Media, 2025. ISBN 9781098167257. Disponível em: https://books.google.com.br/books?id=_3VGEQAAQBAJ.

OSINGA, Douwe. **Deep Learning Cookbook.** 1st. Sebastopol, CA: O'Reilly Media, Inc., jun. 2018.

PELEGRINO, Gabriela da Silva; VIOTO, Jéssica Maria; KERCHÉ, Leandra Ernst. Adesão à profilaxia pós-exposição utilizada para HIV: uma revisão integrativa. **Brazilian Journal of Health Review**, 2022.

SCHACHNER, T.; KELLER, R.; WANGENHEIM, F. von. Artificial Intelligence-Based Conversational Agents for Chronic Conditions: Systematic Literature Review. **Journal of Medical Internet Research**, 2020.