



## Projeto de Pesquisa

### Dados do Projeto Pesquisa

Código:	PVM2264-2025
Titulo do Projeto:	YARA Your Assistant for Results Analysis "Inteligência artificial para geração automática de relatórios bioinformáticos na Amazônia"
Tipo do Projeto:	INTERNO ( Projeto Novo)
Categoria do Projeto:	Iniciação Científica
Situação do Projeto:	SUBMETIDO
Unidade:	DEPARTAMENTO DE ENSINO - CSGC (11.01.16.01.05.09)
Centro:	CAMPUS MANAUS CENTRO - CMC (11.01.03)
Palavra-Chave:	Agente Conversacional; Inteligência Artificial; Relatórios Automatizados; Bioinformática; Amazônia; 16S rRNA.
E-mail:	diego.rios@ifam.edu.br
Edital:	Edital do Programa de Bolsas de Iniciação Científica - Graduação
Cota:	Edital N° 001/2025/PPGI/IFAM - IC Graduação (25/03/2025 a 01/10/2026)

### Área de Conhecimento, Grupo e Linha de Pesquisa

Área de Conhecimento:	Sistemas de Informação
Grupo de Pesquisa:	Ciências da Natureza
Linha de Pesquisa:	Genética, Genômica e Bioinformática

### Comitê de Ética

Nº do Protocolo:	Não possui protocolo de pesquisa em Comitê de Ética.
------------------	--

### Resumo

A análise de dados metagenômicos via sequenciamento 16S rRNA tem se consolidado como uma das principais abordagens para o estudo da diversidade microbiana. No entanto, a interpretação dos resultados gerados por ferramentas como o QIIME 2 ainda representa um desafio significativo para pesquisadores que não possuem formação em bioinformática, especialmente em regiões com infraestrutura limitada como a Amazônia. Essa barreira técnica limita a autonomia científica de grupos locais e reforça a dependência de instituições externas, o que contribui para práticas de extrativismo biotecnológico. Neste contexto, este projeto propõe o desenvolvimento da YARA Your Automated Reporting Agent, uma Inteligência Artificial baseada em agentes conversacionais, capaz de interpretar os resultados das análises geradas pelo QIIME 2 e gerar relatórios comprehensíveis automaticamente. Inspirada em modelos avançados de Processamento de Linguagem Natural (PLN), como o ChatGPT, a YARA será acessível via interface de chat, permitindo que pesquisadores interajam em linguagem natural para obter interpretações, gráficos e relatórios técnicos formatados. Essa solução visa promover acessibilidade computacional, democratizar a bioinformática e fortalecer a soberania científica da região amazônica. Para avaliação da eficácia e usabilidade da ferramenta, serão conduzidos testes com pesquisadores da região, utilizando métodos padronizados como o SUS (System Usability Scale). Com isso, espera-se que a YARA atue como um agente de transformação na forma como a ciência é feita e compreendida localmente, promovendo a inclusão digital, a autonomia na análise de dados biológicos e a preservação do conhecimento na Amazônia.

### Introdução/Justificativa

(incluindo os benefícios esperados no processo ensino-aprendizagem e o retorno para os cursos e para os professores da IFAM em geral)

Os avanços na ciência e na tecnologia vêm possibilitando inovações promissoras, especialmente no contexto da bioinformática e da análise de dados genômicos, transformando a forma como o conhecimento é acessado e utilizado por pesquisadores de diferentes áreas (SADAVARTE; BODANESE, 2019). Pode-se afirmar que a história da ciência é marcada pela evolução das tecnologias computacionais, uma vez que os cientistas estão constantemente em busca de ferramentas mais eficazes para processar e interpretar grandes volumes de dados. Na era da informação, tecnologias emergentes como a Inteligência Artificial (IA) passaram a desempenhar papel fundamental na democratização do acesso ao conhecimento e à análise de dados complexos, como os oriundos do sequenciamento de DNA ambiental (SGARBOSA; VECCHIO, 2020). Segundo Kurzweil e Goldberger (2019), os avanços tecnológicos ocorrem em ciclos cada vez mais curtos, caracterizando uma evolução exponencial. Termos como Inteligência Artificial (IA), Aprendizado de Máquina (Machine Learning - ML) e Processamento de Linguagem Natural (PLN) deixaram de ser conceitos futuristas para se tornarem tecnologias incorporadas em diversas áreas do conhecimento. Na bioinformática, tais ferramentas vêm sendo utilizadas para otimizar desde a triagem de dados até a interpretação de resultados. Esse processo de inovação contínua tem favorecido o surgimento de soluções acessíveis a públicos com pouca ou nenhuma formação em programação ou ciência de dados (TEFFÉ; MEDON, 2020).

A IA é um campo multidisciplinar em constante aperfeiçoamento e desempenha um papel crucial na interação acessível entre humanos e máquinas. Com sua capacidade de aprender, raciocinar e tomar decisões com base em dados e experiências anteriores, a IA simula aspectos do pensamento humano e contribui significativamente para tarefas como extração de padrões, reconhecimento de contexto e geração automatizada de relatórios (RODRIGUES, 2021). Essa flexibilidade e versatilidade tornam a IA uma tecnologia estratégica, com aplicação desde áreas industriais até campos altamente especializados como a bioinformática, facilitando o acesso de pesquisadores a análises antes restritas a especialistas (LIMA et al., 2022).

Nos últimos anos, a aplicação de IA em bioinformática tem sido intensificada pelo crescimento de dados oriundos de tecnologias de sequenciamento de nova geração (NGS), especialmente na análise de comunidades microbianas por meio do gene 16S rRNA. Esse contexto exige soluções automatizadas para a análise, organização e interpretação dos dados, e a IA surge como uma alternativa eficiente para lidar com essa complexidade. Áreas como visão computacional, robótica e, especialmente, o Processamento de Linguagem Natural (PLN), têm sido adaptadas para a extração de informações de grandes volumes de texto e dados biomoleculares (MARTINS, 2010).

O PLN é uma subárea que combina computação e linguística para permitir a comunicação entre humanos e sistemas computacionais (CASELI et al., 2022). Na bioinformática, ele pode ser usado para interpretar descrições taxonômicas, gerar resumos automáticos de resultados analíticos e oferecer feedback textual aos usuários, promovendo um uso mais acessível dos dados gerados por ferramentas como o QIIME 2. O PLN também é fundamental para a criação de agentes conversacionais capazes de entender perguntas e fornecer respostas claras, interpretando e traduzindo os resultados bioinformáticos de forma compreensível por não especialistas (SOUZA; FELIPE, 2021).

Com o aumento do uso de plataformas de mensagens e assistentes virtuais, os agentes conversacionais tornaram-se soluções atrativas por sua capacidade de oferecer suporte personalizado e contínuo ao usuário (TEFFE; MEDON, 2020). Esses sistemas utilizam IA e PLN para responder em linguagem natural, oferecendo uma experiência intuitiva e acessível. Modelos como ChatGPT, Siri, Alexa e Google Assistant são exemplos amplamente reconhecidos (SMUTNY; SCHREIBEROVA, 2020). No contexto científico, os agentes conversacionais podem ser aplicados para fornecer suporte técnico, interpretar resultados experimentais e auxiliar pesquisadores em tempo real.

A proposta do projeto YARA Your Assistant for Research in the Amazon nasce da convergência dessas tecnologias com o objetivo de preencher uma lacuna crítica na análise e interpretação dos resultados bioinformáticos gerados por ferramentas como o QIIME 2. Especialmente na região amazônica, onde há carência de infraestrutura computacional e de especialistas em bioinformática, YARA será uma ferramenta capaz de interagir com pesquisadores de forma acessível, traduzindo saídas analíticas complexas em relatórios interpretáveis, claros e prontos para uso científico ou educacional.

Por meio de uma arquitetura baseada em IA generativa e PLN, integrada com os dados oriundos da interface FLORA e do pipeline SPARTA (software criado por nosso grupo), YARA promoverá a soberania científica regional, reduzindo a dependência de consultorias externas e fortalecendo a autonomia de grupos de pesquisa locais. A solução terá como base modelos de linguagem natural treinados especificamente com vocabulário técnico da área de microbioma, permitindo geração de relatórios contextualizados, consistentes e explicativos.

Este trabalho será desenvolvido com base no projeto YARA Your Automated Reporting Agent, um agente inteligente que integra IA conversacional e Processamento de Linguagem Natural (PLN) para gerar relatórios explicativos a partir dos resultados do QIIME 2. Inspirada na arquitetura do Rasa Open Source, a proposta será adaptada ao contexto amazônico, com foco em acessibilidade computacional, soberania de dados e democratização da bioinformática. A ferramenta visa apoiar pesquisadores na interpretação de análises metagenómicas, ampliando a autonomia científica de instituições da região.

Trata-se de uma iniciativa interdisciplinar e interinstitucional que contará com a colaboração de pesquisadores da EMBRAPA Amazônia Ocidental, do Instituto Nacional de Pesquisas da Amazônia (INPA), do Instituto Federal do Amazonas (IFAM), além de estudantes e docentes das áreas de Biotecnologia, Computação, Linguística, Ecologia Microbiana e Ciências Biológicas. O desenvolvimento será conduzido em articulação com pesquisadores envolvidos na execução de pipelines QIIME 2, de modo a garantir fidelidade aos desafios enfrentados na prática. Assim, o projeto YARA representa um novo patamar na democratização da bioinformática, unindo acessibilidade computacional, inteligência artificial e responsabilidade científica, com impacto direto no fortalecimento da pesquisa genômica na Amazônia e na prevenção de práticas de extrativismo biotecnológico.

## Objetivos

### OBJETIVO GERAL

Desenvolver um agente conversacional inteligente capaz de interpretar os resultados das análises metagenómicas realizadas no QIIME 2 e gerar relatórios automatizados, acessíveis e compreensíveis para pesquisadores da Amazônia, com foco na democratização da bioinformática e na autonomia científica regional.

### Objetivos Específicos

Realizar uma análise bibliográfica sobre arquiteturas de agentes conversacionais baseados em IA e sua aplicabilidade na interpretação de dados científicos;

Desenvolver um protótipo funcional do agente YARA, utilizando modelos de linguagem natural semelhantes ao ChatGPT, com integração direta aos resultados gerados por pipelines 16S do QIIME 2;

Projetar e implementar um fluxo de conversação que permita ao usuário interagir com os resultados bioinformáticos por meio de linguagem natural, obtendo relatórios personalizados, explicações e sugestões de análise;

Validar a efetividade e clareza dos relatórios gerados pelo agente junto a grupos de pesquisa da região amazônica, por meio de testes de usabilidade como o SUS (System Usability Scale) e entrevistas semiestruturadas;

Avaliar o impacto do agente na autonomia de pesquisadores locais, medindo a redução na dependência de especialistas externos para interpretação de dados microbiológicos.

## Metodologia

## Cenários de Interação

Os cenários de interação do agente YARA foram desenhados com base em situações reais enfrentadas por pesquisadores durante análises com o QIIME 2. Esses cenários simulam demandas típicas de interpretação de resultados, como explicação de gráficos de diversidade, taxonomias detectadas e exportação de tabelas. A definição dos fluxos conversacionais será feita com o apoio de especialistas da EMBRAPA e do INPA, os quais contribuirão para selecionar os pontos críticos de interpretação nos pipelines do QIIME 2.

Foram definidos os seguintes cenários iniciais:

1. Explicação de métricas de diversidade alfa (Shannon, Simpson, Observed Features);
2. Interpretação de análises de diversidade beta (PCoA, UniFrac);
3. Visualização e descrição de gráficos de barras taxonômicas;
4. Esclarecimento sobre rarefação e curvas de amostragem;
5. Alertas sobre erros comuns (classificações truncadas, baixa cobertura);
6. Exportação de tabelas e gráficos em formatos diversos.

Esses cenários serão enriquecidos com respostas baseadas em bibliografia atualizada sobre ecologia microbiana e técnicas de metagenômica, com linguagem adaptada para usuários com diferentes níveis de familiaridade com o QIIME 2.

## Escolha da Ferramenta

O Rasa foi escolhido como framework principal por ser open source, transparente e altamente customizável. Como é baseado em Python, permite integração direta com bibliotecas essenciais para análise de dados metagenômicos, como pandas, matplotlib, seaborn, scikit-learn, biom-format e plotly. Sua arquitetura modular facilita a adaptação a vocabulários técnicos específicos da bioinformática, respeitando o contexto regional e científico da Amazônia.

Diferente de soluções comerciais que funcionam como caixas pretas, o Rasa oferece total controle sobre cada fase do ciclo de diálogo. Isso permite a criação de um agente confiável, auditável e customizado para os fluxos de trabalho de pesquisa científica, atendendo aos princípios de soberania de dados e transparência científica.

## Arquitetura Técnica

A arquitetura do YARA será estruturada em três camadas principais:

NLU (Natural Language Understanding): tokenização com WhitespaceTokenizer, vetorização com CountVectorsFeaturizer, extração de padrões via RegexFeaturizer, mapeamento de sinônimos com EntitySynonymMapper e classificação por DIETClassifier. As frases de treinamento serão compostas por perguntas reais de pesquisadores amazônicos, coletadas em oficinas com os parceiros institucionais.

Core: responsável pela política de diálogo, utilizará TEDPolicy com max\_history de 10 turnos e 100 epochs, além de MemoizationPolicy para reaproveitamento de interações comuns e RulePolicy para fluxos determinísticos como comandos de exportação e explicações fixas.

Actions Server: módulo externo em Python que se conectarão aos arquivos de resultados do QIIME 2 para gerar respostas enriquecidas com gráficos, descrições técnicas, explicações sobre ASVs/OTUs, interpretação de gráficos de dispersão e comparação entre grupos. A exportação poderá ser feita como PDF, PNG ou HTML, conforme solicitado pelo usuário.

O fluxo de diálogo ocorrerá da seguinte forma: uma mensagem do usuário no Telegram será interceptada por um webhook (via Ngrok durante o desenvolvimento), enviada ao Rasa, processada para reconhecimento de intenção e entidade, e então uma resposta será calculada com base no modelo treinado e enviada de volta pela API do Telegram.

## Ambiente de Desenvolvimento

O ambiente será configurado com as seguintes etapas:

Instalação do Python 3.8 ou superior;

Uso do Anaconda para gerenciamento de ambientes e pacotes;

Criação de ambiente virtual específico para o YARA;

Instalação do Rasa (pip install rasa) e bibliotecas auxiliares (pandas, numpy, scikit-learn, biom-format, matplotlib, seaborn, plotly, click, rich, jinja2);

Utilização do Ngrok para testes seguros de integração com o Telegram.

As etapas de desenvolvimento seguirão um modelo incremental e participativo, com validação contínua feita junto a grupos de pesquisadores usuários do QIIME 2 na Amazônia.

## Estrutura do Projeto

O projeto será iniciado com rasa init, que criará a estrutura básica contendo:

domain.yml: com intenções como "explicar diversidade", "mostrar taxonomia", "exportar resultados";

nlu.yml: frases de exemplo com marcação de entidades e intenções;

stories.yml: fluxos de diálogo com múltiplas etapas de interação científica;

rules.yml: respostas fixas e mensagens de fallback.

Além disso, será incluído um diretório de actions/ com scripts Python que realizam o parsing de arquivos .qzv, .tsv e .biom produzidos pelo QIIME 2.

## Integração com Telegram

A integração com o Telegram será feita via token gerado no BotFather, configurado em credentials.yml. Durante a fase de desenvolvimento, o webhook será testado com Ngrok para garantir a comunicação segura HTTPS. A API do Telegram será utilizada para ativar o webhook com comandos curl ou requests em Python.

O agente YARA responderá a comandos como:

"Explique a curva de rarefação"

"Qual a diversidade alfa das minhas amostras?"

"Me mostre os grupos mais abundantes"

"Exportar gráfico de taxonomia"

Além disso, serão implementados botões interativos e menus em carrossel para facilitar a navegação por usuários com pouca experiência técnica. A interface será bilíngue (PT/EN), e preparada para funcionar com mensagens curtas e comandos estruturados.

## Referências

1. BAHJA, M. et al. An antenatal care awareness prototype chatbot application using a user-centric design approach. In: SPRINGER. International Conference on Human-Computer Interaction. [S.I.], 2020. p. 2031.
2. BHARTI, U. et al. Medbot: Conversational artificial intelligence powered chatbot for delivering tele-health after covid-19. In: IEEE. 2020 5th International Conference on Communication and Electronics Systems (ICCES). [S.I.], 2020. p. 870875.
3. CASELI, H.; FREITAS, C.; VIOLA, R. Processamento de linguagem natural. Sociedade Brasileira de Computação, 2022.
4. KURZWEIL, R.; GOLDBERGER, A. A singularidade está próxima: quando os humanos transcendem a biologia. [S.I.]: Itaú Cultural, 2019.
5. LIMA, B. N. et al. Inteligência artificial (IA), prototipagem e aplicações da acelerometria controlada por k-nn para análise do movimento humano: Uma revisão bibliográfica. Revista CPAQV - Centro de Pesquisas Avançadas em Qualidade de Vida - CPAQV Journal, v. 14, n. 3, 2022.
6. MARTINS, A. L. Potenciais aplicações da inteligência artificial na ciência da informação. Informação & Informação, v. 15, n. 1, p. 116, 2010.
7. RESENDE, F. K. S. et al. Impactos da inteligência artificial na tomada de decisão médica: Um mapeamento sistemático. In: SBC. Anais da XXI Escola Regional de Computação Bahia, Alagoas e Sergipe. [S.I.], 2021. p. 4150.
8. RODRIGUES, G. V. J. O uso da inteligência artificial na triagem e seleção de processos para conciliação. Revista Consultor Jurídico, 2021.
9. SADAVARTE, S. S.; BODANESE, E. Pregnancy companion chatbot using Alexa and Amazon Web Services. In: IEEE. 2019 IEEE Pune Section International Conference (PuneCon). [S.I.], 2019. p. 15.
10. SANTOS JUNIOR, J. B. et al. Uma proposta de chatbot para apoio a gestantes no contexto do sistema de saúde brasileiro. Revista Ibérica de Sistemas e Tecnologias de Informação, n. E42, p. 344352, 2021.
11. SGARBOSA, P.; VECCHIO, G. H. D. Inteligência artificial e suas implicações: como os dispositivos inteligentes e assistentes virtuais influenciam o cotidiano das pessoas. Revista Interface Técnologica, v. 17, n. 2, p. 193205, 2020.
12. SMUTNY, P.; SCHREIBEROVA, P. Chatbots for learning: A review of educational chatbots for the Facebook Messenger. Computers & Education, Elsevier, v. 151, p. 103862, 2020.
13. SOUZA, A. D. de; FELIPE, E. R. Processamento de linguagem natural aplicado à anamneses do domínio da ginecologia. Fronteiras da Representação do Conhecimento, v. 1, n. 2, p. 5169, 2021.
14. SWICK, R. K. The accuracy of artificial intelligence (AI) chatbots in telemedicine. Journal of the South Carolina Academy of Science, v. 19, n. 2, p. 17, 2021.
15. TEFFÉ, C. S. de; MEDON, F. Responsabilidade civil e regulação de novas tecnologias: questões acerca da utilização de inteligência artificial na tomada de decisões empresariais. REI - Revista Estudos Institucionais, v. 6, n. 1, p. 301333, 2020.

## Membros do Projeto

<b>CPF</b>	<b>Nome</b>	<b>Categoria</b>	<b>CH Dedicada</b>	<b>Tipo de Participação</b>
012.170.506-47	DIEGO LISBOA RIOS	DOCENTE	8	COORDENADOR(A)



- IMPLEMENTAÇÃO DE RESPOSTAS TÉCNICAS COM APOIO DE BIBLIOTECAS PYTHON (PANDAS, SEABORN, BIOM-FORMAT)						
- INTEGRAÇÃO COM MODELOS TREINADOS PARA ENTENDIMENTO DE INTENÇÕES						
- INTEGRAÇÃO COM API DO TELEGRAM E TESTES DE COMUNICAÇÃO VIA WEBHOOK (NGROK)						
- IMPLEMENTAÇÃO DE RESPOSTAS DINÂMICAS COM GERAÇÃO DE GRÁFICOS E LINKS DE EXPORTAÇÃO						
- INSERÇÃO DE BOTÕES INTERATIVOS E MENU DE COMANDOS ESTRUTURADOS NO BOT						
- ADIÇÃO DE FALBACK MESSAGES E TRATAMENTO DE ERROS/AMBIGUIDADE DE INTENÇÕES						
- REFINO DOS FLUXOS DE CONVERSAÇÃO COM BASE EM TESTES INTERNOS						
- OTIMIZAÇÃO DE PERFORMANCE DO MODELO E INTEGRAÇÃO COM DIFERENTES TIPOS DE OUTPUT (PNG, PDF, HTML)						
- VALIDAÇÃO TÉCNICA COM PESQUISADORES DA AMAZÔNIA						
- COLETA DE FEEDBACK QUALITATIVO E AJUSTES FINOS NAS RESPOSTAS E LINGUAGEM UTILIZADA						
- APLICAÇÃO DE TESTE DE USABILIDADE (SYSTEM USABILITY SCALE - SUS)						
- DOCUMENTAÇÃO DE INSTALAÇÃO, TREINAMENTO E USO DO SISTEMA						
- TESTES COM RESULTADOS REAIS DE PROJETOS DA EMBRAPA/INPA						
- ESCRITA DO MANUAL DO USUÁRIO E TUTORIAIS DE USO PARA CIENTISTAS NÃO PROGRAMADORES						
- ESCRITA DE ARTIGO TÉCNICO-CIENTÍFICO DETALHANDO A ARQUITETURA, METODOLOGIA E IMPACTO DO YARA						
- PUBLICAÇÃO DO REPOSITÓRIO GITHUB COM LICENÇA LIVRE						
- ELABORAÇÃO E SUBMISSÃO DO RELATÓRIO FINAL DE INICIAÇÃO CIENTÍFICA						
- PREPARAÇÃO DE APRESENTAÇÃO EM EVENTOS CIENTÍFICOS E DIVULGAÇÃO EM CANAIS INSTITUCIONAIS						

#### Histórico do Projeto

Data	Situação	Usuário
08/04/2025	SUBMETIDO	DIEGO LISBOA RIOS / diego.rios

**Relatório Emitido por:** DIEGO LISBOA RIOS

SIGAA | Copyright © 2014 - DGTI - IFAM