



UNIVERSIDADE FEDERAL DO ESTADO DO RIO DE JANEIRO
CENTRO DE CIÊNCIAS EXATAS E TECNOLÓGICAS
ESCOLA DE INFORMÁTICA APLICADA

Retrieval Augmented Generation Aplicada à Bibliotecas

Breno Costa da Silva Filgueiras

Orientador

Pedro Nuno de Souza Moura

RIO DE JANEIRO, RJ – BRASIL

DEZEMBRO, 2024

Retrieval Augmented Generation Aplicada à Bibliotecas

Breno Costa da Silva Filgueiras

Projeto de graduação apresentado à Escola de Informática Aplicada da Universidade Federal do Estado do Rio de Janeiro (UNIRIO) como cumprimento de requerimento parcial para obtenção título de Bacharel em Sistemas de Informação.

Approved by:

Pedro Nuno de Souza Moura – UNIRIO

Supervisor 2, D.Sc. – UNIRIO

Supervisor 3, D.Sc. – XXXX

RIO DE JANEIRO, RJ – BRASIL

DEZEMBRO, 2024

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Resumo

Em uma parceria entre Seagate e a International Data Corporation (IDC) foi realizado o estudo *The Digitization of the World From Edge to Core*, nele a IDC fala sobre diversos aspectos referentes aos dados presentes no mundo digital e um dos tópicos abordados no estudo é “Mankind is on a quest to digitize the world” e neste mesmo tópico eles explicam que os dados que geramos no dia a dia está em constante crescimento, ou seja, estamos gradualmente produzindo mais dados.

Com um volume cada vez maior de dados, uma busca por informação otimizada é essencial, dado que são necessárias ferramentas que nos garantam confiança e precisão da informação adquirida. Com isso em mente, este trabalho visa o desenvolvimento de um sistema capaz de ler, processar e armazenar documentos diversos de determinada biblioteca (conjunto de documentos) para que possamos utilizar um *Large Language Model* (LLM) para responder perguntas que os usuários possam ter acerca dos documentos.

A ideia é conseguir processar documentos de diferentes épocas, temas, formatos e conseguir responder o maior número possível de perguntas dos usuários com a melhor confiança possível.

Palavras-chave: retrieval, augmented, generation, inteligência, artificial.

Abstract

In a partnership between Seagate and the International Data Corporation (IDC), the study *The Digitization of the World From Edge to Core* was conducted. In it, IDC discusses various aspects related to data present in the digital world and one of the topics covered in the study is “Mankind is on a quest to digitize the world”. In this same topic, they explain that the data we generate on a daily basis is constantly growing, that is, we are gradually producing more data.

With an ever-increasing volume of data, an optimized search for information is essential, given that tools are needed that guarantee reliability and accuracy of the information acquired. With this in mind, this work aims to develop a system capable of reading, processing and storing various documents from a given library (set of documents) so that we can use a *Large Language Model* (LLM) to answer questions that users may have about the documents.

The idea is to be able to process documents from different periods, themes and formats and to be able to answer as many user questions as possible with the greatest possible confidence.

Keywords: retrieval, augmented, generation, artificial, intelligence.

Conteúdo

1	Introdução	1
1.1	Motivação	2
1.2	Objetivos	2
1.3	Organização do Texto	3
1.4	Metodologia	4
2	Conceitos Fundamentais	5
2.1	IA Generativa	5
2.1.1	Uma Breve História da IA Generativa	5
2.1.2	Mais à Frente	7
2.2	Large Language Models	7
2.2.1	Como funcionam?	8
2.3	O Problema	9
2.4	Retrieval Augmented Generation (RAG)	9
2.5	Ecossistema RAG	11
2.5.1	Retriever (D)	12
2.5.2	Generator (G)	14
2.5.3	Evaluator (E)	15
2.5.4	Trainer (T)	15
3	Modelagem	16
3.0.1	Pipeline de Processamento	16
3.0.2	Interface de Programação de Aplicações (API)	16
3.0.3	Interface Gráfica de Usuário (GUI)	17
3.1	Modelo Conceitual	17

3.1.1	Ingestão e Preparo	18
3.1.2	Geração de Embeddings	18
3.1.3	Armazenamento dos Embeddings	19
3.1.4	Query	19
3.1.5	Busca por informações relevantes	19
3.1.6	Envio de Query, Prompt e Contexto	19
3.1.7	Retorno de Resposta Humanizada	19
3.2	Tecnologias	20
3.2.1	Python e Ambientes Virtuais	20
3.2.2	Docling	21
3.2.3	Llama Index	21
3.2.4	Elastic Search	21
3.2.5	Ollama	22
3.2.6	Fast API	22
3.2.7	Streamlit	23
4	Solução Desenvolvida	24
4.1	Pipeline de Ingestão	24
4.2	Interface de Programação de Aplicações	26
4.3	Interface Gráfica de Usuário	27
5	Conclusão	28
5.1	Considerações Finais	28
5.2	Limitações	28
5.3	Trabalhos Futuros	29

Lista de Figuras

1	Total de dados produzidos por ano, [Duarte, 2024]	1
2	Funcionamento geral de uma estrutura RAG, [Rothman, 2024]	10
3	Arquitetura da estrutura RAG, [Rothman, 2024]	12
4	Visão arquitetural da solução RAG proposta para o trabalho	18

Lista de Tabelas

1 Introdução

Com a expansão contínua da Internet das Coisas (IoT), o cenário se transforma em um redemoinho de informações. Chegará, ou talvez já tenha chegado, o momento em que será impossível para qualquer ser humano consumir tudo o que criou em um único dia.

Durante estudo, a *International Data Cooperation* (IDC) previu que a *Global Datasphere* cresceria de 45 zettabytes em 2019 para 175 zettabytes em 2025 [David Reinsel e Rydning, 2018]. Um crescimento de aproximadamente 380% em 6 anos, porém essa previsão foi feita em 2018 e atualmente já existem estudos que projetam números ainda maiores para a produção de dados. Como é o exemplo do estudo *Amount of Data Created Daily*, onde foram calculados um total de 147 *zettabytes* produzidos em 2024, com previsão de 181 *zettabytes* para 2025, um crescimento de 23.12% [Duarte, 2024].

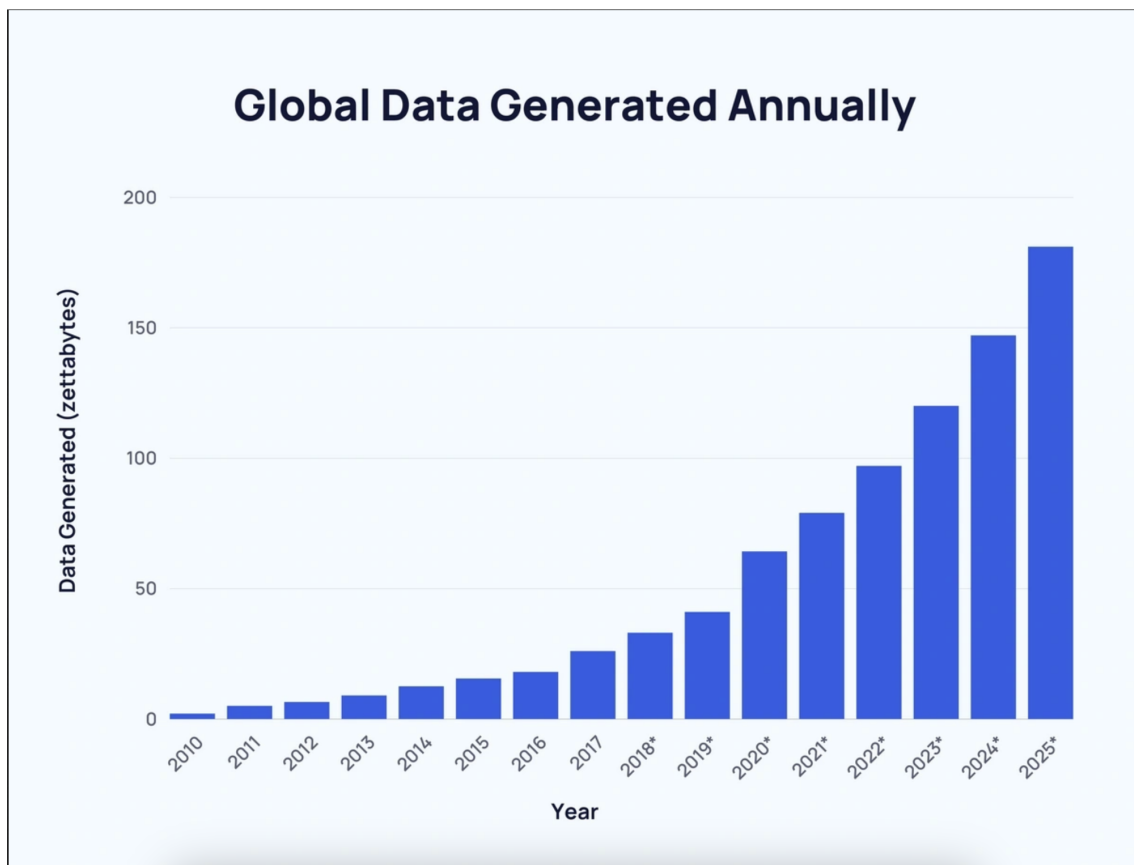


Figura 1: Total de dados produzidos por ano, [Duarte, 2024]

1.1 Motivação

No estudo *The Digitization of the World From Edge to Core*, a Seagate, gigante do armazenamento de dados, uniu forças com a *International Data Corporation* (IDC) para conduzir uma análise dos dados presentes na *Global Datasphere*, que quantifica e analisa o total de dados criados, capturados e replicados no mundo inteiro. A IDC destacou: “*Mankind is on a quest to digitize the world.*” Essa frase encapsula a era em que vivemos, marcada por um crescimento incessante no volume de dados que produzimos diariamente.

Cada clique, pagamento por aproximação ou uso de *wearables* adiciona mais um fragmento ao vasto oceano digital. Nesse turbilhão de dados, buscar uma matéria ou reportagem torna-se uma tarefa semelhante a encontrar uma agulha em um palheiro digital, um desafio tão fascinante quanto frustrante.

E esse contexto de imensidão de dados onde a busca por informações é cada vez mais difícil, é o berço deste projeto. O objetivo é implementar uma solução que processe bibliotecas de documentos e aplique o conceito de *Retrieval Augmented Generation* (RAG). Com uma interface de chat simples, o usuário poderá fazer perguntas e receber respostas humanizadas, geradas por um *Large Language Model* (LLM), com referências claras aos documentos de origem.

O desafio de buscar informações relevantes é significativo. A internet ainda abriga dados sem referência ou apresentados de formas variadas, como gráficos e textos, dificultando a assimilação. Além disso, interfaces pouco intuitivas e mecanismos de busca ineficazes consomem tempo valioso. Para estudantes e pesquisadores, essa batalha constante com a desorganização digital pode transformar o simples ato de encontrar informações em um verdadeiro labirinto.

1.2 Objetivos

O objetivo principal deste trabalho é implementar uma solução baseada em *Retrieval Augmented Generation* (RAG) para bibliotecas de documentos específicos, a fim de viabilizar consultas que retornem dados pertinentes junto com suas referências. Os objetivos específicos são:

1. Escrever uma introdução acessível ao conceito de RAG aos alunos do BSI.
2. Produzir um documento que instrua a implementação de um ecossistema RAG aos alunos do BSI.
3. Realizar a implementação de RAG para documentos, que seja agnóstica tanto ao LLM quanto embedding utilizados.
4. Executar uma validação sobre a solução gerada, de maneira que o resultado seja relevante.

1.3 Organização do Texto

Este trabalho está organizado em capítulos, com o objetivo de apresentar os processos, métodos, análises e descobertas de forma clara e coerente. A estrutura do documento foi elaborada para facilitar a compreensão do leitor sobre a complexidade do tema e os resultados obtidos, conduzindo-o até a implementação da solução final. Os capítulos estão estruturados da seguinte forma:

- **Introdução:** Apresenta o contexto do trabalho, destacando o problema do crescente volume de dados. Discute a motivação para a solução proposta, sua relevância no cenário atual e os objetivos estabelecidos.
- **Conceitos Fundamentais:** Dedicar-se à fundamentação teórica, abordando os conceitos essenciais para a compreensão da solução e sua implementação. Inclui uma introdução ao conceito de *Large Language Models* (LLM) e uma análise detalhada do *Retrieval Augmented Generation* (RAG).
- **Modelagem:** Descreve a composição da solução, explicando os artefatos envolvidos e suas responsabilidades. Também aborda o funcionamento e o papel de cada componente na solução final. Ao final, apresenta as tecnologias utilizadas, incluindo descrições breves sobre as ferramentas, suas versões e funções.
- **Solução Desenvolvida:** Detalha os artefatos implementados, explicando como a solução cumpre suas funções. Apresenta os resultados obtidos e discute as respostas fornecidas para algumas das questões propostas, avaliando a eficácia da solução.

- **Conclusão:** O capítulo final resume as considerações sobre os resultados alcançados, destacando tanto os aspectos positivos quanto as limitações da solução implementada. Além disso, discute possíveis trabalhos futuros ou aplicações derivadas da solução, encerrando com as referências utilizadas no desenvolvimento do trabalho.

1.4 Metodologia

Este trabalho adotará a abordagem de *Design Science Research* (DSR) para garantir que ao final do trabalho, o artefato modelado esteja implementado e funcionando conforme planejado.

O DSR, possui raízes na engenharia e nas ciências do artificial [Simon, 2019], é uma metodologia voltada para a resolução de problemas. Seu objetivo é aprimorar o conhecimento humano por meio da criação de artefatos inovadores e da geração de conhecimento de design, oferecendo soluções práticas para problemas do mundo real [Alan Hevner, março de 2004].

Assim, ao utilizar o *Design Science Research* (DSR), este trabalho resultará em um artefato produzido com base nas análises e discussões realizadas ao longo das próximas seções deste trabalho.

2 Conceitos Fundamentais

Para que este trabalho seja compreendido e os próximos capítulos possam ser apresentados com maior clareza, é necessário passarmos por alguns conceitos. Antes de nos aprofundarmos no contexto de um ecossistema de *retrieval augmented generation* (RAG), é necessário compreender um pouco a IA generativa e os *Large Language Models* (LLMs) que contribuem tanto para a interpretação das perguntas feitas durante as interações com o usuário, quanto na geração de uma resposta mais humana.

2.1 IA Generativa

A inteligência artificial generativa, às vezes chamada de *gen AI*, é um tipo de inteligência artificial (IA) capaz de criar conteúdo original — como texto, imagens, vídeos, áudio ou código de software — em resposta a um comando ou solicitação do usuário. [Mark Scapicchio, março de 2024]

A inteligência artificial generativa baseia-se em modelos avançados de *machine learning* (aprendizado de máquina) chamados modelos de *deep learning* (aprendizagem profunda) — algoritmos que simulam os processos de aprendizado e tomada de decisão do cérebro humano. Esses modelos trabalham identificando e codificando padrões e relações em grandes volumes de dados.

A partir dessas informações, a IA generativa é capaz de compreender solicitações ou perguntas feitas em linguagem natural pelos usuários, respondendo com conteúdos novos e relevantes. Essa capacidade permite a criação de textos, imagens, vídeos, áudios e até códigos de software, de forma original e adaptada ao pedido do usuário.

2.1.1 Uma Breve História da IA Generativa

O termo “IA generativa” explodiu na consciência pública na década de 2020, mas a IA generativa já faz parte de nossas vidas há décadas, e a tecnologia de IA generativa atual se baseia em avanços de aprendizado de máquina que remontam ao início do século 20. Uma história representativa não exaustiva da IA generativa pode incluir algumas das seguintes datas:

- **1964:** O cientista da computação do *Massachusetts Institute of Technology* (MIT), Joseph Weizenbaum, desenvolve o ELIZA, uma aplicação de processamento de linguagem natural baseada em texto. Essencialmente o primeiro *chatbot* (chamado de "*chatterbot*" na época), o ELIZA usava *scripts* de correspondência de padrões para responder a entradas de linguagem natural digitadas com respostas empáticas em texto.
- **1999:** A Nvidia lança o GeoForce, a primeira unidade de processamento gráfico (GPU). Originalmente desenvolvida para fornecer gráficos de movimento suave para videogames, as GPUs se tornaram a plataforma padrão para o desenvolvimento de modelos de IA e mineração de criptomoedas.
- **2004:** O Google *autocomplete* aparece pela primeira vez, gerando palavras ou frases potenciais à medida que os usuários digitam seus termos de busca.
- **2013:** Aparecem os primeiros *autoencoders* variacionais (VAEs).
- **2014:** Surgem as primeiras redes adversariais generativas (GANs) e modelos de difusão.
- **2017:** Ashish Vaswani, uma equipe do Google Brain e um grupo da Universidade de Toronto publicam o artigo *Attention is All you Need*, [Ashish Vaswani, junho de 2017], um artigo que documenta os princípios dos modelos de transformadores, amplamente reconhecidos como os responsáveis por permitir os modelos de fundação mais poderosos e as ferramentas de IA generativa que estão sendo desenvolvidas hoje.
- **2019-2020:** O OpenAI lança seus modelos de linguagem GPT (*Generative Pre-trained Transformer*), o GPT-2 e o GPT-3.
- **2022:** O OpenAI apresenta o ChatGPT, uma interface do GPT-3 que gera frases complexas, coerentes e contextuais, além de conteúdo de longo formato em resposta a comandos dos usuários.
- **2023-2024:** Com a notoriedade e popularidade do ChatGPT, que efetivamente abriu as portas para uma onda de desenvolvimentos, os avanços e lançamentos

de produtos em IA generativa têm ocorrido a um ritmo acelerado, incluindo lançamentos do Google Bard (agora Gemini), Microsoft Copilot, IBM watsonx.ai e o modelo de linguagem Llama-2 de código aberto da Meta.

2.1.2 Mais à Frente

A inteligência artificial tem sido um tema relevante na tecnologia, mas foi a IA generativa, especialmente com o lançamento do ChatGPT em 2022, que a destacou globalmente, gerando inovação e adoção. Ela oferece grandes benefícios de produtividade para indivíduos e organizações, e, apesar dos desafios e riscos, as empresas exploram como melhorar fluxos de trabalho e enriquecer produtos e serviços. De acordo com uma pesquisa da consultoria *McKinsey*, mais de 65% das empresas usam Gen AI no mundo. [Yran Bartolomeu Dias, julho de 2024].

2.2 Large Language Models

Os *Large Language Models* (LLMs) são uma categoria de modelos fundamentais treinados em grandes volumes de dados para oferecer capacidades versáteis, atendendo a diversos casos de uso e tarefas. Diferentemente dos modelos específicos para determinados domínios, que exigem treinamentos separados para cada aplicação—geralmente com altos custos e demandas significativas de infraestrutura—os LLMs promovem uma aplicação mais ampla, gerando sinergias entre diferentes áreas e, muitas vezes, alcançando um desempenho superior. [IBM, novembro de 2023]

Os LLMs representam um avanço significativo em *Natural Language Processing* (NLP) e inteligência artificial. Esses modelos estão amplamente acessíveis ao público por meio de interfaces como o *ChatGPT-3* e *GPT-4* da *OpenAI*, apoiados pela *Microsoft*. Outros exemplos incluem os modelos *Llama* da Meta, os modelos *BERT/RoBERTa* e *PaLM* do Google, e a série *Granite* lançada pela IBM.

Os LLMs são projetados para compreender e gerar texto de forma similar à humana, além de produzir outros tipos de conteúdo. Com base nos extensos volumes de dados em que foram treinados, conseguem traduzir idiomas, resumir textos, responder perguntas, auxiliar na redação e até mesmo na geração de código.

Essas capacidades são viabilizadas por bilhões de parâmetros que capturam padrões complexos da linguagem. Como resultado, os LLMs estão transformando áreas como chatbots, assistentes virtuais, geração de conteúdo, suporte à pesquisa e tradução de idiomas.

2.2.1 Como funcionam?

Os *Large Language Models* (LLMs) operam utilizando técnicas de aprendizagem profunda e grandes volumes de dados textuais. Baseados na arquitetura de transformadores [Ashish Vaswani, junho de 2017], como o *Generative Pre-trained Transformer* (GPT), esses modelos são especialmente eficazes em lidar com dados sequenciais, como entradas de texto. Compostos por várias camadas de redes neurais, os LLMs empregam um mecanismo de atenção para focar em partes específicas dos conjuntos de dados.

Durante o treinamento, os modelos aprendem a prever a próxima palavra em uma sentença com base no contexto fornecido pelas palavras anteriores. Isso é feito atribuindo probabilidades à recorrência de palavras que foram tokenizadas (divididas em sequências menores) e transformadas em embeddings, representações numéricas do contexto.

O treinamento envolve o uso de corpora massivas, com bilhões de páginas de texto, permitindo que os LLMs aprendam gramática, semântica e relações conceituais por meio de aprendizado auto-supervisionado e técnicas de zero-shot. Uma vez treinados, os modelos geram texto prevendo autonomamente a próxima palavra com base na entrada recebida e nos padrões adquiridos, resultando em uma produção linguística coerente e relevante para diversas tarefas de compreensão e geração de linguagem.

O desempenho dos LLMs pode ser aprimorado por meio de técnicas como *prompt engineering*, *fine-tuning* (ajuste fino) e aprendizado por reforço com feedback humano (*RLHF*). Essas abordagens ajudam a mitigar problemas como vieses, discurso de ódio e respostas incorretas ou ilusórias (*hallucinations*), que podem surgir devido ao treinamento em dados não estruturados. Garantir que os LLMs estejam prontos para uso em nível corporativo é crucial para evitar riscos à reputação e responsabilidades indesejadas.

2.3 O Problema

No livro *RAG-Driven Generative AI*, é dito que “mesmo o modelo mais avançado de Inteligência Artificial (IA) generativa é limitado a responder sobre dados nos quais ela foi treinada.” [Rothman, 2024] o que nos chama a atenção a um problema em especial, como fazer para que uma IA saiba responder perguntas referentes a um conjunto específico de dados?

Uma IA não tem como saber o que ela não sabe, ou seja, não existe conhecimento fora do conjunto de dados nos quais ela foi treinada. Perguntas fora do contexto do treinamento de uma IA geralmente levam a halucinações, viés e respostas sem sentido. Para isso, foi implementado um framework, ou estrutura, que combina abordagens baseadas em recuperação com modelos generativos, esta estrutura é a Retrieval Augmented Generation (RAG).

A RAG recupera dados relevantes de fontes externas em tempo real e usa esses dados para gerar respostas contextualmente relevantes. Uma de suas principais vantagens é a adaptabilidade, tendo em vista que a estrutura pode ser aplicada independente do tipo de dado abordado na solução, seja texto, imagens, áudios ou documentos diversos.

2.4 Retrieval Augmented Generation (RAG)

Quando um modelo de IA generativa não sabe responder determinada pergunta com precisão, diz-se que ele está alucinando ou apresentando viés, mas, na prática, está apenas gerando respostas sem sentido. Isso ocorre porque o modelo não foi treinado com as informações solicitadas ou por conta de limitações em sua configuração, resultando em sequências prováveis, mas não precisas necessariamente.

A RAG começa onde a IA generativa termina, fornecendo informações que um modelo de LLM não possui para responder com precisão. A RAG otimiza tarefas de recuperação de informações e adiciona os dados recuperados durante a entrada (seja consulta do usuário ou um prompt automatizado), gerando uma saída melhorada e mais amigável ao usuário. O funcionamento geral do RAG pode ser resumido na figura a seguir:

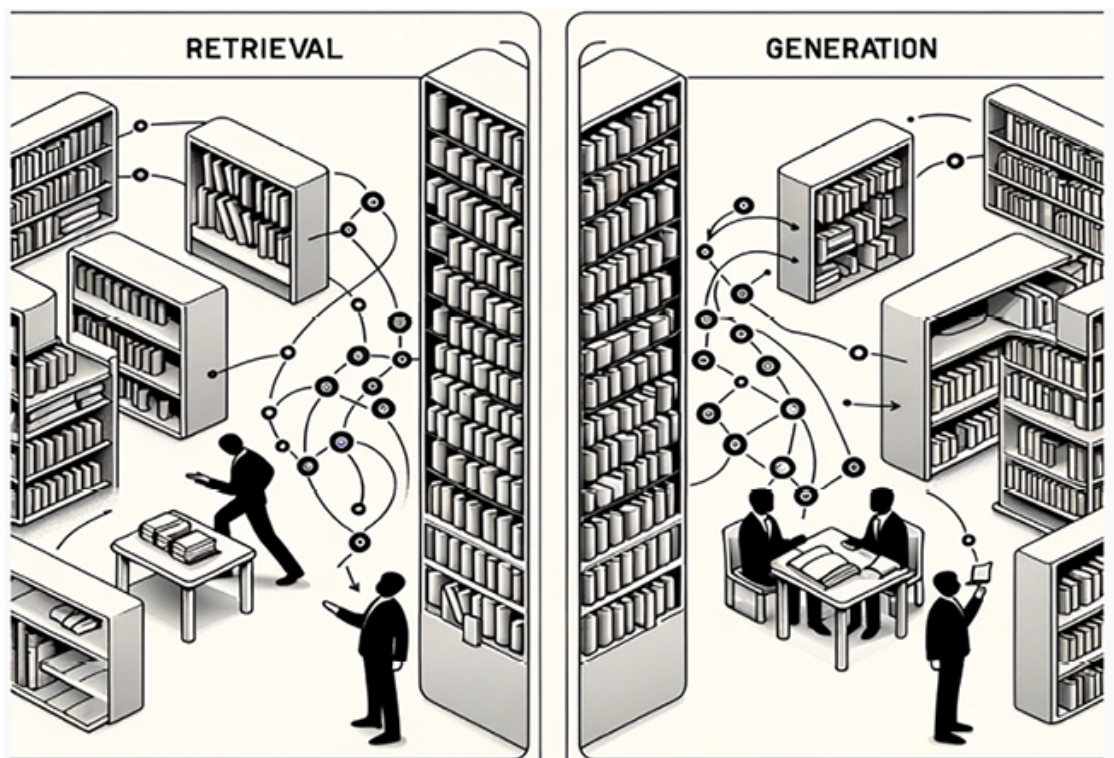


Figura 2: Funcionamento geral de uma estrutura RAG, [Rothman, 2024]

Imagine um estudante em uma biblioteca, com a tarefa de escrever uma dissertação sobre RAG. Assim como o ChatGPT ou outros copilotos de IA, o estudante sabe ler e escrever. Como qualquer LLM, ele é treinado para compreender informações avançadas, resumir e criar conteúdo. No entanto, como qualquer IA avançada, seja do Hugging Face, Vertex AI ou OpenAI, há muitas informações que este estudante ainda desconhece.

Na fase de recuperação, ele busca por livros sobre o tema necessário (lado esquerdo da figura 1) na biblioteca. Em seguida, ele retorna ao seu lugar, realiza a tarefa de recuperação sozinho ou com a ajuda de um colega, extraíndo as informações relevantes dos livros adquiridos. Na fase de geração (lado direito da figura 1), o estudante começa a escrever sua dissertação. Assim, funciona como um agente humano guiado por RAG, de maneira semelhante a um framework de IA generativa baseado em RAG.

Enquanto escreve o seu ensaio sobre RAG, o estudante encontra tópicos difíceis com os quais não tem tempo para consultar todas as informações disponíveis. Como um agente humano generativo, ele fica travado, assim como um modelo de IA generativa. Ele até pode tentar escrever algo, mas, como a IA, não saberá se o conteúdo está correto até que alguém corrija a dissertação e lhe avalie de alguma maneira.

Neste ponto, ele já atingiu seu limite e decide recorrer a um copiloto de IA generativa com RAG para obter respostas corretas. No entanto, existe uma variedade tão grande de modelos LLM e configurações RAG disponíveis que o estudante acaba ficando confuso. Antes de prosseguir, é necessário entender os recursos disponíveis e como o RAG está organizado.

2.5 Ecossistema RAG

A IA generativa baseada em RAG é um framework que pode ser implementado com diversas configurações, funcionando dentro de um ecossistema amplo (Figura 2). Independentemente da quantidade de frameworks de recuperação e geração disponíveis, tudo se resume a quatro domínios principais e suas respectivas questões:

- **Dados:** De onde vêm os dados? São confiáveis e suficientes? Há questões de direitos autorais, privacidade ou segurança?
- **Armazenamento:** Como os dados serão armazenados antes ou depois do processamento? Qual será o volume armazenado?
- **Recuperação:** Como os dados corretos serão recuperados para complementar o input do usuário? Qual tipo de framework RAG será mais adequado ao projeto?
- **Geração:** Qual modelo de IA generativa melhor se adapta ao framework RAG escolhido?

Esses domínios dependem do tipo de framework RAG utilizado. Antes de escolher, é essencial avaliar a proporção de conhecimento paramétrico e não paramétrico no ecossistema implementado. A Figura 1.3 ilustra os principais componentes do framework RAG, independentemente do tipo implementado.

- **Retriever (D):** Responsável pela coleta, processamento, armazenamento e recuperação de dados.
- **Generator (G):** Cuida da complementação do input, engenharia de prompts e geração de respostas.

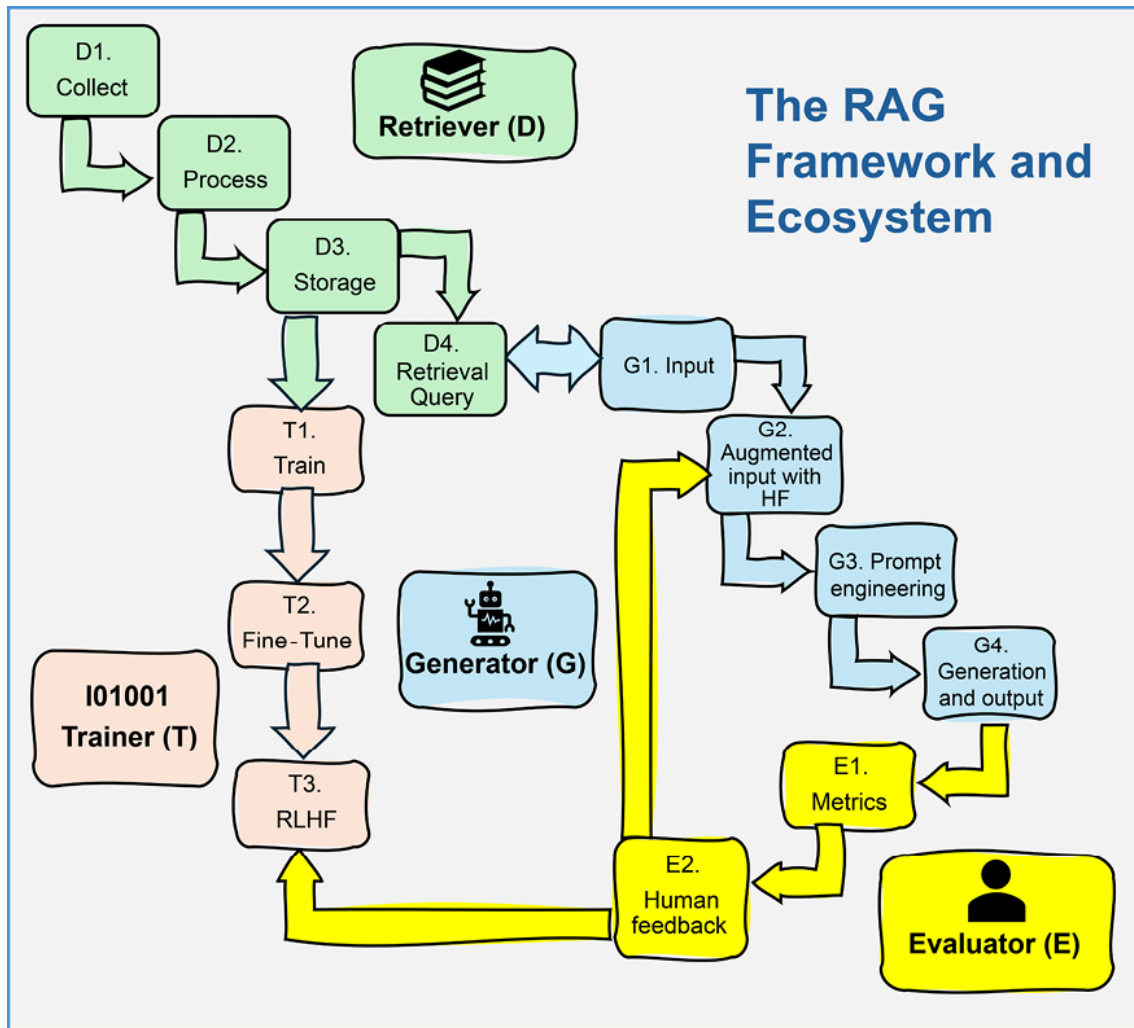


Figura 3: Arquitetura da estrutura RAG, [Rothman, 2024]

- **Evaluator (E):** Avalia o desempenho usando métricas matemáticas, feedback humano e outras formas de validação.
- **Trainer (T):** Gerencia o modelo pré-treinado inicial e sua posterior fine-tuning.

Esses quatro componentes dependem de seus respectivos ecossistemas, formando o pipeline de IA generativa baseada em RAG. Nas seções a seguir, usaremos as siglas D, G, E e T para representar respectivamente Retriever, Generator, Evaluator e Trainer. Começando pelo Retriever (D).

2.5.1 Retriever (D)

O componente retriever de um ecossistema RAG coleta, processa, armazena e recupera dados. O ponto de partida de um ecossistema RAG é, portanto, um processo de ingestão

de dados, cujo primeiro passo é a coleta de dados.

1. **Coleta de Dados (D1):** Atualmente dados são extremamente diversos, podendo ser textos, arquivos de mídia (como músicas ou vídeos em mp4) ou arquivos estruturados e não estruturados (PDFs, JSONs e páginas web). Além disso, grande parte desses dados é não estruturada e pode ser encontrada de maneiras imprevisíveis e complexas. Felizmente, várias plataformas, como Pinecone, OpenAI, Chroma e Activerloop, oferecem ferramentas prontas para processar e armazenar essa vasta quantidade de dados.
2. **Processamento de Dados (D2):** Na fase de coleta de dados (D1) no processamento de dados multimodais, diferentes tipos de dados, como texto, imagens e vídeos, podem ser extraídos de websites utilizando técnicas de web scraping ou outras fontes de informação. Esses objetos de dados são então transformados para criar representações uniformes. Alguns exemplos dessas transformações incluem: chunking, embedding e indexação. Essas técnicas serão discutidas mais adiante.
3. **Armazenamento de Dados (D3):** Neste estágio do pipeline, já coletamos e começamos a processar uma grande quantidade de dados diversos. Mas para fazermos com que esses dados sejam úteis, vamos fazer uso de vetor stores (armazenamento de vetores), como Elastic Search. Ele não apenas armazena os dados, mas os convertem em entidades matemáticas, representadas como vetores, permitindo realizar cálculos poderosos. Esses sistemas também utilizam técnicas de indexação e outras abordagens para garantir acesso rápido e eficiente aos dados. Em vez de manter os dados em arquivos estáticos, transformamos tudo em um sistema dinâmico e pesquisável, pronto para alimentar chatbots, motores de busca e outras aplicações.
4. **Consulta de Recuperação (D4):** O processo de recuperação é acionado pelo input do usuário ou input automatizado (G1). Para recuperar dados rapidamente, carregamos os dados nos vetor stores e datasets após transformá-los para um formato adequado. Em seguida, utilizamos uma combinação de pesquisas por palavras-chave, embeddings inteligentes e indexação para recuperar os dados de forma eficiente. A similaridade cosseno, por exemplo, encontra itens que estão

intimamente relacionados, garantindo que os resultados da busca não sejam apenas rápidos, mas também altamente relevantes. Após a recuperação dos dados, o próximo passo é aumentar o input, ou seja, adicionar as informações recuperadas para enriquecer a resposta gerada.

2.5.2 Generator (G)

No ecossistema RAG, as linhas entre a entrada e a recuperação não são tão nítidas, como mostrado na figura 2, que representa o framework e ecossistema RAG. O input do usuário (G1), seja automatizado ou humano, interage com a consulta de recuperação (D4) para complementar o input antes de enviá-lo ao modelo generativo. O fluxo gerativo começa com o input, que é aprimorado com dados recuperados antes de ser processado pelo modelo de IA generativa.

1. **Entrada (G1):** O input pode ser uma série de tarefas automatizadas (como o processamento de e-mails, por exemplo) ou prompts humanos por meio de uma Interface de Usuário (UI). Essa flexibilidade permite integrar a IA de forma fluida em diversos ambientes profissionais, aprimorando a produtividade em diferentes setores.
2. **Entrada Aumentada com Feedback Humano (G2):** O feedback humano (HF, de Human Feedback) pode ser adicionado ao input, conforme descrito na seção Feedback Humano (E2), sob o componente Evaluator (E). O feedback humano torna o ecossistema RAG consideravelmente mais adaptável, permitindo total controle sobre a recuperação de dados e os inputs para a IA generativa.
3. **Engenharia de Prompts (G3):** Tanto o retriever (D) quanto o generator (G) dependem fortemente da engenharia de prompts para preparar a mensagem padrão e aumentada que o modelo de IA generativa deverá processar. A engenharia de prompts combina a saída do retriever com o input do usuário, garantindo que o modelo receba uma entrada bem estruturada e relevante para gerar a resposta desejada.
4. **Geração e Saída (G4):** A escolha de um modelo de IA generativa depende dos

objetivos do projeto. Modelos como Llama¹, Gemini², GPT³ e outros podem atender a diferentes requisitos. No entanto, o prompt precisa estar alinhado com as especificações de cada modelo.

2.5.3 Evaluator (E)

Frequentemente, dependemos de métricas matemáticas para avaliar o desempenho de um modelo de IA generativa. No entanto, essas métricas fornecem apenas uma parte do todo. É importante lembrar que o teste final da eficácia de uma IA depende da avaliação humana, que garante uma compreensão mais completa da qualidade e relevância dos resultados gerados.

1. Métricas (E1):

Um modelo não pode ser avaliado sem métricas matemáticas, como a similaridade cosseno, assim como em qualquer sistema de IA. Essas métricas garantem que os dados recuperados sejam relevantes e precisos. Ao quantificar as relações e a relevância dos pontos de dados, elas fornecem uma base sólida para avaliar o desempenho e a confiabilidade do modelo.

2. Feedback Humano (E2):

Em um sistema de IA generativa, seja ele baseado em RAG ou não, e independentemente de as métricas matemáticas parecerem suficientes, o feedback humano é essencial. A avaliação humana é o fator decisivo que determina se um sistema projetado para usuários humanos será aceito ou rejeitado, elogiado ou criticado. O RAG adaptativo introduz o feedback humano, pragmático e da vida real, o que melhora o ecossistema de IA generativa baseado em RAG.

2.5.4 Trainer (T)

Um modelo de IA generativa padrão é pré-treinado com uma grande quantidade de dados de uso geral. Em seguida, podemos ajustar (T2) o modelo com dados específicos de um domínio.

1. Modelos disponíveis em [Meta, 2024]
2. Modelos disponíveis em [Google, 2024]
3. Modelos disponíveis em [OpenAI, 2024]

3 Modelagem

Com os conceitos de um ecossistema RAG em mente, é possível explicar de forma objetiva a modelagem da solução proposta para este trabalho. Antes de abordar as tecnologias específicas, tema reservado para um capítulo posterior, será apresentada uma visão abstrata da modelagem da solução.

De maneira objetiva, o projeto pode ser descrito como um sistema composto por três artefatos, cada um responsável por uma parte do ecossistema proposto. Esses artefatos implementam sistemas ou serviços que se comunicam entre si ao longo da solução, podendo implementar uma LLM, uma API ou uma instância de banco vetorial. Apesar de estarem conectados, cada artefato é independente e funciona de forma autônoma.

Juntos, esses três artefatos formam a solução proposta. Proporcionando uma abordagem modular e integrada para alcançar os objetivos propostos neste trabalho. Dentre os artefatos temos respectivamente um pipeline de processamento, uma interface de programação de aplicações (API) e uma interface gráfica de usuário.

3.0.1 Pipeline de Processamento

Artefato responsável por gerenciar fluxos de análise e transformação de documentos para a solução. Neste contexto, ele desempenha a função de lidar com bibliotecas de documentos, analisando cada documento e garantindo que cada um seja processado corretamente. Esse processo envolve etapas como leitura, análise e transformação do conteúdo dos documentos, com o objetivo de prepará-los para demais etapas do processo.

Ao término do processamento de cada documento na biblioteca, os dados processados são inseridos em um banco vetorial. Assim, o pipeline converte os dados para um formato propício para utilização no ecossistema em questão. Essa abordagem garante escalabilidade e precisão no gerenciamento e utilização dos dados processados.

3.0.2 Interface de Programação de Aplicações (API)

Artefato central no ecossistema, sendo responsável por mediar a comunicação entre os diferentes componentes do sistema. Sua principal função é receber e processar as re-

quisições feitas pelos usuário na GUI, garantindo que as operações sejam realizadas de forma eficiente e segura.

Além disso, a API desempenha um papel crítico na segurança do sistema ao restringir o acesso direto ao banco vetorial por parte dos outros artefatos. Essa camada de abstração impede interações não autorizadas ou inadequadas com o banco de dados, garantindo que apenas as operações de usuários autenticados sejam realizadas. Com isso, a API promove uma integração controlada e otimizada entre os artefatos, além de assegurar a integridade dos dados e a consistência do sistema como um todo.

3.0.3 Interface Gráfica de Usuário (GUI)

Artefato responsável por receber e guiar o usuário durante suas interações com o sistema, assim como também é responsável em garantir a comunicação com o artefato da API.

Artefato dedicado a mediar a interação entre o sistema e usuário de forma intuitiva. Sua principal função é receber as entradas do usuário e orientá-lo durante a navegação pelo sistema, oferecendo elementos visuais para simplificar o uso das funcionalidades disponíveis.

Além de interagir diretamente com o usuário, a GUI desempenha um papel importante na comunicação com a API, garantindo que as solicitações realizadas pelo usuário sejam traduzidas em requisições apropriadas para o sistema. Essa integração possibilita que somente ações como consultas de dados sejam realizadas pelo usuário, de forma transparente e sincronizada com o restante dos artefatos.

3.1 Modelo Conceitual

Analisar uma biblioteca de documentos pode ser complexo, pois cada biblioteca pode comportar diversos tipos de documentos, em formatos distintos. Este projeto foi construído visando a análise de qualquer biblioteca de documentos, ou seja, o processamento dos dados da biblioteca será realizado mesmo que ela tenha arquivos em formatos diferentes. O projeto foi contruído com base nos três artefatos previamente explicados e funciona de acordo com a seguinte visão arquitetural:

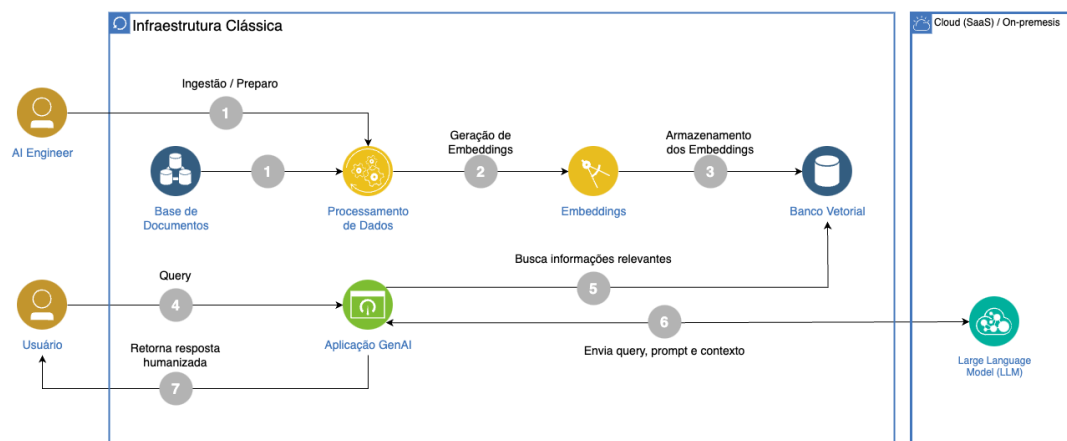


Figura 4: Visão arquitetural da solução RAG proposta para o trabalho

3.1.1 Ingestão e Preparo

Etapa onde a biblioteca é analisada, seus documentos são processados e transformados de maneira sequencial para obtermos um determinado modelo com as informações necessárias de cada documento. Esta etapa é dividida em duas partes, uma de ingestão e outra de preparo.

- **Ingestão:** Nesta parte os documentos são lidos e sua extensão é analisada, determinando a forma como o documento terá seus dados extraídos. Uma vez analisados, o documento é dividido em chunks, um bloco de texto com os dados do documento.
- **Preparo:** Para cada chunk obtido na ingestão, é instanciado um modelo contendo os dados gerais do documento, assim como os dados do próprio chunk. Este modelo, em específico, visa uma inserção no banco vetorial que será feita na etapa de armazenamento dos embeddings.

3.1.2 Geração de Embeddings

Antes de armazenar os modelos obtidos na etapa anterior, a informação de cada chunk precisa estar em um formato específico para o banco vetorial, este formato é obtido quando geramos os embeddings, que neste contexto, são uma representação numérica de baixa dimensão dos dados de um chunk específico.

3.1.3 Armazenamento dos Embeddings

Uma vez gerados, os embeddings de cada chunk estão presentes no modelo previamente adquirido na etapa de Ingestão e Preparo. Nesta etapa ocorre o armazenamento dos modelos na instância do banco vetorial do projeto. O armazenamento é feito de forma sequencial, porém em batches ou grupos com tamanho baseado na quantidade de chunks obtidos. Isso ocorre pois reduz o tempo da etapa de inserção ao banco, tendo em vista que uma vez que um documento é particionado em chunks, dependendo do tamanho do documento, podemos obter mais de 100 chunks, ou seja, 100 inserções.

3.1.4 Query

Aqui ocorre a interação do usuário com o sistema, nesta etapa o usuário digita uma pergunta dentro do espaço de input ou entrada da interface gráfica, que neste contexto é composta por um chat.

3.1.5 Busca por informações relevantes

No momento em que o usuário realiza a query, o sistema realiza uma busca por contexto ou informações relevantes como a biblioteca que o usuário está consultando (caso hajam múltiplas bibliotecas), uma determinada data, autor ou demais filtros que o usuário possa vir a fazer. Com as informações relevantes analisadas, é gerado um objeto de contexto que contém todas as informações relevantes

3.1.6 Envio de Query, Prompt e Contexto

Com o contexto e a query em mãos, o sistema prossegue com o fluxo principal e envia uma solicitação a API contendo a query que o usuário realizou, o prompt a ser seguido pela LLM e o contexto.

3.1.7 Retorno de Resposta Humanizada

Ao final da requisição feita pelo usuário, a API responde com uma lista contendo os objetos presentes no banco vetorial cujo embedding contém informações relevantes para

a query do usuário. Esta lista contém os modelos presentes no banco vetorial, definidos na etapa de Preparação, e a partir dessa resposta passamos a lista para uma LLM que irá analisar primeiramente um prompt focado em entregar respostas humanizadas para o usuário. Neste prompt, há todo um cuidado para que qualquer desenvolvedor que deseje implementar este projeto possa utilizar do prompt engineering a fim de realizar as afinações necessárias para seu próprio contexto.

3.2 Tecnologias

3.2.1 Python e Ambientes Virtuais

O Python 3.11 é uma das versões mais recentes da linguagem de programação Python, lançada com melhorias significativas em desempenho, novas funcionalidades e correções de bugs. Uma das principais inovações dessa versão é a otimização do tempo de execução, com um aumento significativo na velocidade em relação às versões anteriores, devido a melhorias no interpretador e na execução de código. Além disso, a versão 3.11 introduziu aprimoramentos na sintaxe, como a simplificação de anotações de tipos e novos recursos que facilitam o desenvolvimento de aplicativos mais robustos e eficientes. A linguagem continua a ser uma das mais populares, especialmente em áreas como desenvolvimento web, ciência de dados e inteligência artificial.

Ambientes virtuais são espaços isolados onde é possível instalar e gerenciar dependências específicas de um projeto, sem afetar o sistema global ou outros projetos. Em Python, esses ambientes são criados com o módulo `venv`, permitindo que cada projeto tenha suas próprias versões de pacotes e bibliotecas. Isso evita conflitos entre dependências e facilita o gerenciamento de diferentes versões de pacotes para cada projeto.

As principais vantagens dos ambientes virtuais incluem a capacidade de isolar dependências, garantindo que cada projeto utilize a versão correta de pacotes, a facilidade de replicar e compartilhar ambientes de desenvolvimento e a redução de problemas de compatibilidade entre pacotes ou versões do Python. Eles oferecem um controle mais preciso sobre o ambiente de execução, tornando o desenvolvimento mais seguro e eficiente.

3.2.2 Docling

Uma biblioteca Python projetada para facilitar a extração e organização de documentos. Ela permite que desenvolvedores processem automaticamente documentos de diferentes formatos, como PDFs, textos e arquivos de Word, para extrair informações estruturadas e relevantes. A biblioteca é útil para automatizar tarefas de leitura e processamento de grandes volumes de dados, organizando-os de forma a facilitar a análise posterior.

A principal vantagem do Docling é a sua capacidade de simplificar o processo de extração de dados, tornando-o mais rápido e eficiente, sem a necessidade de escrever código complexo para lidar com diferentes tipos de documentos. Isso a torna uma ferramenta valiosa para projetos que requerem análise de texto ou a criação de bancos de dados a partir de documentos não estruturados.

3.2.3 Llama Index

Uma biblioteca Python que oferece ferramentas para criar índices eficientes, que podem ser usados em conjunto com modelos de linguagem, como o GPT, para realizar tarefas de busca e processamento de texto com alta performance.

A principal vantagem do LlamaIndex é sua capacidade de organizar dados complexos de forma que seja fácil e rápido realizar consultas sobre eles, integrando-se bem com outras ferramentas de IA e processamento de linguagem natural. Isso a torna útil para criar soluções que necessitam de pesquisa e análise de grandes quantidades de informações, como em assistentes virtuais, sistemas de recomendação e análise de documentos.

3.2.4 Elastic Search

O Elasticsearch é um serviço de banco de dados baseado em busca de texto completo e análise de grandes volumes de dados. Ele é amplamente utilizado para armazenar, pesquisar e analisar dados de forma rápida e escalável. Baseado no motor de busca Apache Lucene, o Elasticsearch permite que dados sejam indexados e consultados com alta performance, suportando operações como pesquisa por palavras-chave, agregações e análises em tempo real.

Em um banco vetorial, como o Elasticsearch, os dados são representados como vetores, o que permite consultas de similaridade, ou seja, encontrar documentos que sejam semanticamente semelhantes aos dados consultados. Isso é especialmente útil em casos de busca por texto, recomendação de conteúdo ou sistemas de recuperação de informações avançados. O Elasticsearch organiza e indexa esses vetores para proporcionar consultas eficientes e rápidas, sendo um componente essencial em sistemas que exigem grandes volumes de dados e análises complexas.

3.2.5 Ollama

A LLM do Ollama é uma ferramenta de inteligência artificial focada em fornecer modelos de linguagem avançados para tarefas como processamento de texto, geração de conteúdo, tradução, resumo e respostas a perguntas. Desenvolvido para ser altamente eficiente e acessível, o Ollama oferece modelos de linguagem que são otimizados para uma variedade de aplicações, desde automação de processos até assistência personalizada em diferentes domínios.

A principal vantagem da LLM do Ollama é sua facilidade de integração em soluções de software, além de seu desempenho rápido e preciso. A plataforma é projetada para ser simples de usar, com uma interface intuitiva e suporte a diferentes linguagens de programação, tornando a tecnologia de IA mais acessível para desenvolvedores e empresas. Com um foco em flexibilidade e escalabilidade, o Ollama facilita a implementação de modelos de linguagem em larga escala sem comprometer a eficiência.

3.2.6 Fast API

Um framework moderno e de alto desempenho para construção de APIs em Python. Ele é baseado em padrões como Python type hints e pydantic, permitindo que desenvolvedores escrevam código de forma rápida e eficiente, com uma validação de dados robusta. O FastAPI se destaca pela sua velocidade, sendo capaz de processar requisições de maneira extremamente rápida, o que o torna ideal para aplicações que exigem alta performance.

As principais vantagens do FastAPI incluem sua facilidade de uso, a validação automática de dados, a documentação interativa gerada com o Swagger UI, e sua perfor-

mance comparável a frameworks como Node.js e Go. Além disso, o FastAPI aproveita os recursos de tipagem do Python, o que resulta em um código mais claro e com menos erros, além de ser altamente escalável e fácil de integrar com outras ferramentas.

3.2.7 Streamlit

Uma biblioteca em Python usada para criar interfaces web interativas de forma rápida e simples, sem a necessidade de conhecimentos avançados em desenvolvimento front-end. Ele permite que desenvolvedores criem dashboards, visualizações de dados e aplicativos interativos com poucos comandos, utilizando apenas Python. O Streamlit é ideal para protótipos rápidos e aplicativos que precisam integrar visualizações com análises de dados.

As principais vantagens do Streamlit incluem sua simplicidade de uso, já que permite criar interfaces com poucas linhas de código, e sua integração direta com bibliotecas como Pandas, Matplotlib e Plotly, facilitando a visualização e análise de dados. Além disso, ele permite a criação de aplicativos interativos com atualizações em tempo real, sem a necessidade de configurar complexos sistemas de front-end. Isso torna o Streamlit uma opção excelente para cientistas de dados e desenvolvedores que desejam criar interfaces intuitivas sem precisar aprender HTML, CSS ou JavaScript.

4 Solução Desenvolvida

Com uma visão geral de como funciona um ecossistema RAG e quais artefatos são necessários para implementá-lo, prosseguiremos para o desenvolvimento da solução deste projeto. Levando em consideração as tecnologias apresentadas anteriormente, os próximos tópicos irão abordar o desenvolvimento concreto da solução buscada para este projeto.

Para a implementação da solução descrita, a biblioteca utilizada foi de todos os trabalhos de conclusão de curso do ano de 2023 publicados pela secretaria do Bacharelado em Sistemas de Informação da Universidade Federal do Estado do Rio de Janeiro, UNIRIO.

4.1 Pipeline de Ingestão

Tendo em vista que uma biblioteca de documentos contém diversos documentos e cada documento possui sua própria extensão, é necessário desenvolver uma solução capaz de processar diversos tipos de documentos, podendo atender diversas bibliotecas distintas.

Este processamento foi implementado seguindo um Design Pattern, mais especificamente o factory design pattern. O factory design pattern implementa uma interface para a criação de objetos que herdam de uma superclasse, permitindo que as classes filhas possam alterar o tipo de objeto criado. Neste contexto, foi implementada uma fábrica de extratores, todo extrator possui uma instancia dos clientes provedores de serviços do LlamaIndex e Docling, assim como variáveis pré definidas e os métodos de extração.

No entanto, cada extrator implementa o método de extração a sua própria maneira, assim é possível implementar a extração de um arquivo .pdf ou .docx mudando somente a implementação de um único método. Isso contribui para uma melhor manutenção do código e atribui a responsabilidade de extração para um único extrator específico de cada documento.

Após extraídos os dados dos documentos, é necessário separar os dados relevantes de cada documento e particionar o mesmo em chunks. Para isso, antes de particionar os documentos é realizada uma busca por dados relevantes do documento, neste caso, por se tratar de uma biblioteca de TCCs os dados relevantes são o resumo (abstract),

nome, biblioteca a qual ele pertence, caminho (filepath), extensão e mídias do arquivo.

As mídias são separadas em dois tipos, imagens e tabelas, durante a primeira leitura do documento a biblioteca Docling identifica as mídias dos documentos e as salva em um diretório específico dentro do código do artefato do pipeline, posteriormente este diretório será injetado no código do artefato da API.

As informações relevantes devem então ser extraídas no começo da extração do documento, uma vez que a busca pelo resumo do trabalho pode tomar tempo de processamento, seria imprudente realizar uma busca após a partição dos documentos. Uma vez extraídas, é gerado um breve resumo do abstract por uma LLM, neste caso foi utilizado o Ollama com um prompt que enfatiza o resumo de no máximo 2 paragrafos contendo 250 caracteres. Essas informações nos geram o seguinte modelo:

```
{
  "abstract": "resumo extraído do documento",
  "summary": "um breve resumo do abstract",
  "filename": "nome do arquivo sem extensão",
  "database": "biblioteca a qual o arquivo
    pertence",
  "filepath": "caminho onde o arquivo se encontra",
  "medias": {
    "images": "lista com o caminho das imagens
      achadas no arquivo",
    "tables": "lista com o caminhos das tabelas
      achadas no arquivo"
  },
  "extension": "extensão do arquivo"
}
```

Com as informações gerais acima, podemos dar continuidade a partição dos documentos em chunks. Neste projeto, cada chunk possui 2048 tokens. Cada chunk gerado será lido e terá o embedding gerado para a inserção no banco vetorial, neste projeto o modelo de embedding é parametrizado mas optei por seguir esta implementação utili-

zando o modelo `impira/layoutlm-document-qa` pois é um modelo específico para RAG e já treinado para melhor atender nosso caso de uso. Assim, obtemos o seguinte modelo para os chunks:

```
{
  "_id": "chave única do chunk",
  "chunk": "objeto imutável obtido durante
    partição",
  "sequence": "sequência do chunk em questão",
  "pages": "lista de páginas com a informação do
    chunk em questão",
  "tables": "lista de caminhos das tabelas
    achadas no arquivo",
  "figures": "lista de caminhos das imagens
    achadas no arquivo",
  "vector": "embeddings gerado para o chunk"
}
```

Este processo de leitura e processamento irá se repetir diversas vezes até que o último documento seja devidamente processado. Uma vez processados, os documentos estarão no formato de chunk seguindo o modelo acima e estarão prontos para inserção no banco vetorial. A inserção no banco irá selecionar cem chunks e realizar uma inserção em batch dos mesmos, devido ao tamanho da biblioteca em questão cem foi um número razoável de chunks, mas dependendo da biblioteca será necessário aumentar este número.

4.2 Interface de Programação de Aplicações

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna.

Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

4.3 Interface Gráfica de Usuário

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

5 Conclusão

5.1 Considerações Finais

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

5.2 Limitações

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum

pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

5.3 Trabalhos Futuros

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Referências

Referências

- Alan Hevner, Salvatore T. March e Salvatore T., Alan R. Março de 2004. “Design Science in Information Systems Research” (). https://www.researchgate.net/publication/201168946_Design_Science_in_Information_Systems_Research.
- Aline Dresch, Daniel Pacheco Lacerda e José Antonio Valle Antunes Júnior. 2020. *Design Science Research: Método de Pesquisa para Avanço da Ciência e Tecnologia*. 1ª edição. Bookman Editora.
- Ashish Vaswani, et al., Noam Shazeer. Junho de 2017. “Attention Is All You Need” (). <https://arxiv.org/abs/1706.03762>.
- Company, The Search AI. 2024. *Elasticsearch*. <https://www.elastic.co/>.
- David Massey, et al., Mihai Criveti. Janeiro de 2024. “Retrieval Augmented Generation” (). <https://pages.github.ibm.com/solution-architectures/generative-ai-ref-arch/patterns/rag/>.
- David Reinsel, John Grantz e John Rydning. 2018. *The Digitization of the World From Edge to Core*. <https://www.seagate.com/files/www-content/our-story/trends/files/dataage-idc-report-final.pdf>. Último acesso em: 19 de dezembro de 2024.
- Duarte, Fabio. 2024. *Amount of Data Created Daily*. <https://explodingtopics.com/blog/data-generated-per-day>.
- Foundation, Python Software. 2024. *Python*. <https://www.python.org/doc/>.
- Google. 2024. *Gemini Models*. https://ai.google.dev/?utm_source=website&utm_medium=referral&utm_campaign=geminichat&utm_content.
- IBM. Novembro de 2023. “What are LLMs?” (). <https://www.ibm.com/think/topics/language-models>.
- . 2024. *Docling*. <https://ds4sd.github.io/docling/>.
- Inc., Snowflake. 2024. *Streamlit*. <https://streamlit.io/>.

- Jan Vom Brocke, Alan Hevner e Alexander Maedche. Setembro de 2020. “Introduction to Design Science Research” (). https://www.researchgate.net/publication/345430098_Introduction_to_Design_Science_Research.
- Mark Scapicchio, Cole Stryker e. Março de 2024. “What is generative AI?” (). <https://www.ibm.com/think/topics/generative-ai>.
- Meta. 2024. *Llama Models*. <https://www.llama.com/llama-downloads/>.
- OpenAI. 2024. *GPT Models*. <https://platform.openai.com/docs/models>.
- Rothman, Denis. 2024. *RAG-Driven Generative AI*. 1ª edição. Packt Publishing.
- Simon, Herbert A. 2019. *The Sciences of the Artificial*. Reissue of the third edition with a new introduction by John Laird. MIT Press.
- UNIRIO, BSI. 2024. *Publicações de TCCs do Bacharelado em Sistemas de Informação*. <https://bsi.uniriotec.br/publicacoes-de-tcc/>. Último acesso em: 19 de dezembro de 2024.
- Yran Bartolomeu Dias, Pepe Cafferata e. Julho de 2024. “65 por cento das empresas usam Gen AI no mundo; líderes apontam caminhos para obter retorno financeiro” (). <https://www.mckinsey.com.br/our-insights/all-insights/65-das-empresas-usam-gen-ai-no-mundo>.