



RAG

Retrieval-Augmented *Generation*

Para Leigos

Sandeco Macedo

*Tudo o que você
precisa saber para
ampliar o conhecimento
da sua LLM*

Os meus alunos do Instituto Federal de Goiás, amo vocês!

Copyright © 2025

Meus Livros sobre Inteligência Artificial



Resultado
das minhas
pesquisas em IA



Sumário

1	Gerando Revoluções	6
1.1	Quando usar RAG	7
1.2	Decodificando a Arquitetura: Os Bastidores do RAG	8
1.3	Um RAG Simples	11
1.4	Ele pode lembrar, RAG com Memória	12
1.5	O RAG Autônomo: Agent RAG	13
1.6	O RAG Corretivo: CRAG	14
1.7	O RAG Adaptativo: Adaptive RAG	15
1.8	O RAG em Grafos: GraphRAG	16
1.9	O RAG Híbrido: Hybrid RAG	17
1.10	O RAG-Fusion: Reciprocal Rank Fusion (RRF)	18
1.11	Hypothetical Document Embedding	18
2	O RAG Clássico	20
2.1	O que é um Corpus de Textos	21
2.2	O Fluxo do RAG	22
2.3	A Fase do Indexador	24
2.4	Lendo e convertendo	25
2.5	Chunking	27
2.6	Criação de Embeddings	29
2.7	Banco de dados Vetorial	32
2.8	Nossa Classe de Encoder	33
2.9	Recuperando conhecimento	36
2.10	Aumento de informação	39
2.11	Gerando a Resposta	41
2.12	Rodando com Streamlit	42
2.13	Exercícios	45
3	Rag com Memória	46
3.1	Criando a Memória com o Redis	49
3.2	Redis no Docker	50
3.3	Instalando o Redis	53

3.4	Criando a Memória	54
3.5	Adicionando a Memória ao RAG	57
3.6	Exercícios	62
4	RAG Autônomo com Agentic RAG	63
4.1	AgenticRAG	64
4.2	Classe de Registro de Datasets	67
4.3	Agente Abstrato	68
4.4	Agente com API da LLM	70
4.5	Agentic RAG com CrewAI	73
4.6	Executando o Agentic RAG	76
4.7	Exercícios	78
5	O RAG em Grafos: GraphRAG	79
5.1	Do Vetor ao Grafo	80
5.2	Entidades e Relações	82
5.3	Construindo o Grafo de Conhecimento	85
5.4	Consultando no Grafo	88
5.5	Classe GraphRetriever	91
5.6	Fluxo do GraphRAG	94
5.7	Executando o GraphRAG	95
5.8	Exercícios	96
6	O RAG Híbrido: Hybrid RAG	97
6.1	A Ideia do Híbrido	98
6.2	Combinação Vetorial e Lexical	99
6.3	Integração com Grafos	100
6.4	Classe HybridRetriever	102
6.5	Estratégias de União e Interseção	104
6.6	Reranking Neural	106
6.7	Executando o Hybrid RAG	108
6.8	Exercícios	109
7	RAG Corretivo: CRAG	110
7.1	Entendendo o CRAG	111
7.2	Pipeline Corretivo	113
7.3	Classe Avaliadora de Qualidade	115
7.4	Integrando Avaliação ao Retriever	117
7.5	Fluxo Completo com CRAG	118
7.6	Executando o CRAG	119

7.7	Exercícios	120
8	O RAG Adaptativo: Adaptive RAG	121
8.1	Visão Geral do Adaptive RAG	122
8.2	Classificador de Consultas	123
8.3	RAG Sem Recuperação	124
8.4	RAG com Recuperação Simples	125
8.5	RAG com Recuperação Iterativa	127
8.6	Implementando o Roteador Adaptativo	129
8.7	Executando o Adaptive RAG	131
8.8	Exercícios	132
9	O RAG-Fusion: Reciprocal Rank Fusion (RRF)	133
9.1	Por que Usar Fusão de Rankings	134
9.2	O Algoritmo RRF	136
9.3	Classe RankFusion	138
9.4	Integração com Hybrid RAG	140
9.5	Pipeline Completo com RRF	142
9.6	Executando o RAG-Fusion	145
9.7	Exercícios	146
10	Hypothetical Document Embedding (HyDE)	147
10.1	O Problema do Desalinhamento	149
10.2	Geração de Documentos Hipotéticos	151
10.3	Classe HyDEGenerator	154
10.4	Integração com Retriever	157
10.5	Fluxo do HyDE	159
10.6	Executando o HyDE	161
10.7	Exercícios	162

CAPÍTULO 1

Gerando Revoluções

Imaginem o seguinte cenário: vocês têm um colega de classe, o Léo, que é simplesmente um gênio. O cara leu todos os livros da biblioteca, assistiu a todas as aulas e tem uma capacidade absurda de conversar sobre qualquer assunto. Ele manja muito de história, física, arte e até daquela matéria de cálculo que todo mundo sofre. A escrita dele é fluida, convincente e ele consegue conectar ideias como ninguém. Só que o Léo tem um pequeno problema: a memória dele, apesar de vasta, às vezes prega peças. Ele pode confundir uma data, atribuir uma citação à pessoa errada ou, na pior das hipóteses, inventar um detalhe que soa verdadeiro só pra manter a conversa fluindo. Esse é o LLM puro: brilhante, mas não totalmente confiável.

Agora, vamos botar esse 'gênio' à prova. Entreguem a ele uma prova final sobre a 'Revolução Francesa', mas com uma pergunta bem específica: 'Qual foi o preço exato do pão em Paris na semana anterior à queda da Bastilha e como isso influenciou os discursos de Camille Desmoulins?'. O Léo, usando apenas sua memória, provavelmente vai desenrolar uma resposta incrível, bem escrita e contextualizada. Ele vai falar sobre a fome, a crise econômica e o papel dos oradores. Mas o preço exato do pão? Ele talvez erre por alguns centavos ou invente um valor que pareça plausível. A essência da resposta estará lá, mas o dado crucial, o fato bruto, pode estar incorreto.

É aqui que a mágica do RAG começa a brilhar. Agora, imaginem que, ao lado do Léo, senta a Ana. A Ana não tem a memória enciclopédica do Léo, mas ela é a rainha da organização e da pesquisa. Ela tem um fichário perfeitamente indexado com resumos de todos os livros da biblioteca. Quando o professor faz a pergunta, antes de o Léo começar a escrever, a Ana entra em ação. Ela não lê o livro inteiro sobre a Revolução Francesa. Ela vai direto na ficha certa, encontra o parágrafo exato sobre a economia pré-revolução e entrega ao Léo um pequeno cartão com a informação: 'Preço do pão: 4 sous. Discursos de Desmoulins mencionaram o 'preço absurdo' como estopim para a revolta popular'. A Ana, nesse caso, é o nosso 'Retriever'.

Com esse cartão em mãos, o Léo se transforma. Ele pega aquela informação precisa e a integra em sua genialidade narrativa. A resposta dele agora não é apenas bem escrita e

Retrieval-Augmented Generation



Figura 1.1: Leo e Ana. Gerador e Retriever

contextualizada, ela é factualmente irrefutável. Ele começa o texto com o dado exato, conecta o preço do pão à inflação da época e tece uma análise brilhante sobre como Desmoulins usou essa informação para inflamar a população. Ele não apenas 'colou' a informação da Ana; ele a 'aumentou' com seu poder de geração de texto. Essa colaboração perfeita entre o pesquisador focado (Ana) e o gerador eloquente (Léo) é a essência do 'Retrieval-Augmented Generation'.

Pensem nessa dinâmica em escala. A 'biblioteca' não é apenas sobre a Revolução Francesa, mas sim sobre todos os documentos da sua empresa, todos os artigos médicos publicados nos últimos 20 anos, toda a base de conhecimento de um produto ou todos os livros do Sandeco. O fichário da Ana é o nosso banco de vetores, e a habilidade dela de encontrar a ficha certa é o algoritmo de busca. Parâmetros como **chunk_size** definem o tamanho de cada 'resumo' no fichário dela, enquanto **top_k** determina quantos 'cartões' de informação ela entrega ao Léo para cada pergunta. Ajustar esses detalhes é o que transforma uma boa resposta em uma resposta perfeita.

1.1 QUANDO USAR RAG

A decisão de usar RAG em vez de outras técnicas, como o Fine-Tuning, não é apenas uma escolha técnica, mas uma decisão estratégica. Ambas as abordagens buscam especializar um LLM em um domínio de conhecimento, mas o fazem de maneiras fundamentalmente diferentes e com implicações radicalmente distintas em custo, agilidade e manutenção. Entender quando o RAG brilha é o primeiro passo para construir sistemas de IA robustos e eficientes.

O cenário ideal para o RAG é aquele onde a **verdade é volátil**. Pensem em qualquer base de conhecimento que não seja estática. A legislação de um país, por exemplo, está em

constante fluxo: novas leis são sancionadas, decretos são publicados e artigos são alterados. Um sistema de IA para advogados que foi treinado ou mesmo 'fine-tuned' em um Vade Mecum de 2023 se tornaria obsoleto e perigosamente impreciso em 2024. O custo para retreinar ou refinar o modelo a cada nova portaria seria proibitivo. O RAG resolve isso com uma elegância impressionante. A 'inteligência' do modelo (sua capacidade de ler e interpretar) é separada do 'conhecimento' (os documentos). Quando uma lei muda, você não toca no cérebro do LLM; você simplesmente atualiza o arquivo de texto correspondente no seu Vector Store. A verdade é atualizada em segundos, a um custo marginal.

Isso nos leva ao segundo ponto crucial: **custo e complexidade**. O Fine-Tuning, por mais poderoso que seja para alterar o 'comportamento' ou o 'estilo' de um modelo, é um processo computacionalmente intensivo. Ele exige a preparação de datasets massivos de exemplos, um poder de processamento considerável (geralmente múltiplas GPUs de ponta rodando por horas ou dias) e um conhecimento técnico aprofundado para ajustar os hiperparâmetros e evitar problemas como o 'catastrophic forgetting'. É como reformar a fundação de um prédio. O RAG, em comparação, é como mobiliar o prédio. A implementação consiste em 'plugar' componentes: um carregador de documentos, um modelo de embedding e um banco de vetores. A atualização do conhecimento é uma simples operação de escrita em um banco de dados, algo trivial em termos de custo computacional.

Portanto, a regra geral é clara: se o seu desafio é injetar conhecimento factual, específico e, principalmente, **mutável** em um LLM, o RAG é quase sempre a resposta correta. Ele oferece um caminho mais barato, rápido e sustentável para manter sua IA aterrada nos fatos e sincronizada com a realidade do seu domínio de negócio.

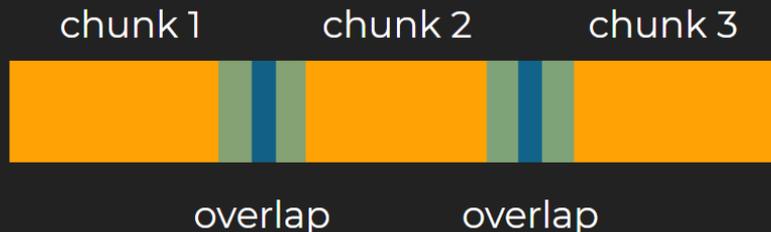
1.2 DECODIFICANDO A ARQUITETURA: OS BASTIDORES DO RAG

A historinha do Léo e da Ana foi legal, né? Ajudou a dar uma clareada nas ideias. Mas agora vamos tirar os apelidos e colocar os nomes técnicos na mesa. O nosso gênio Léo é o que chamamos de **Generator**, geralmente um Modelo de Linguagem (LLM). A nossa pesquisadora supereficiente Ana é o **Retriever**. A colaboração entre eles forma o **sistema RAG**, e esse processo todo tem dois grandes momentos: uma fase de preparação, que acontece antes mesmo de vocês fazerem qualquer pergunta, e uma fase de execução, que rola em tempo real.

Fase 1: A Preparação (Indexação de Conhecimento)

Pensem na Ana montando o fichário dela. Isso não acontece na hora da prova, no desespero. Ela faz antes, com calma e método. No RAG, esse rolê se chama **Indexação**. É aqui que a gente pega uma montanha de informação desestruturada e a organiza de uma forma que o nosso sistema consiga 'pesquisar' de maneira inteligente. A gente basicamente cria o cérebro externo do nosso LLM.

O primeiro passo é **carregar os documentos**. Esqueçam livros físicos; pensem em arquivos PDF, páginas de um site, documentos de texto, transcrições de vídeos, o que for. Depois de carregar tudo, a gente precisa **quebrar esses documentos em pedaços menores**. A gente chama esses pedaços de 'chunks'. Por que fazemos isso? Simples. Mandar um livro de 500 páginas para o LLM analisar e encontrar uma única frase é ineficiente e caro. É muito mais inteligente entregar só o parágrafo certo. Aqui, vocês já encontram dois parâmetros cruciais: o **chunk_size**, que define o tamanho de cada pedaço de texto, e o **chunk_overlap**, que é uma pequena sobreposição entre os 'chunks' pra garantir que a gente não perca o contexto de uma ideia que começa no final de um pedaço e termina no início do outro. Sacaram?



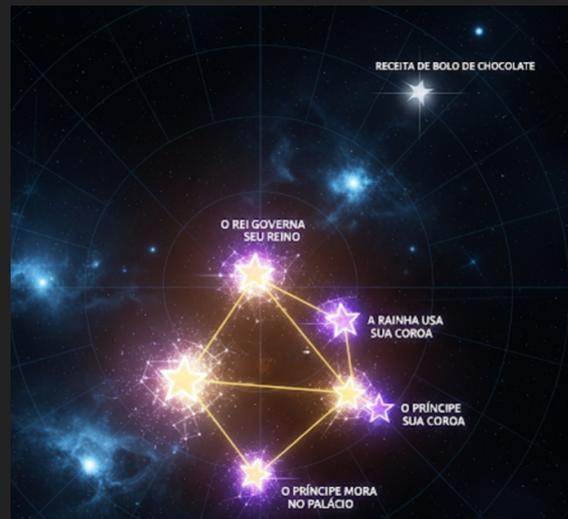
Criando os mapas

Pensem no universo de todas as ideias como uma galáxia gigante e escura. Cada 'chunk' de texto que vocês criaram, cada conceito, é uma estrela brilhando nessa imensidão. Sozinho, é apenas um ponto de luz. O que o modelo de **Embedding** faz é atuar como um astrônomo superavanzado, equipado com um telescópio que enxerga o significado.

Ele não apenas vê as estrelas; ele cria um mapa tridimensional (na verdade, com centenas de dimensões) dessa galáxia. Para cada estrela, ou seja, para cada 'chunk' de texto, ele atribui um conjunto de coordenadas matemáticas únicas e superprecisas. Essa lista de coordenadas é o **vetor de embedding**.

E aqui vem a parte genial, a mágica do mapa: ele não é aleatório. O nosso astrônomo

(o modelo) posiciona as estrelas de forma que aquelas com 'brilho' ou 'energia' conceitual parecida fiquem próximas. A estrela que representa o texto 'O rei governa seu reino' não vai estar perto da estrela 'Receita de bolo de chocolate'. Mas ela vai estar muito próxima das estrelas 'A rainha usa sua coroa' e 'O príncipe mora no palácio'. Elas formam uma **'constelação de significado'**.



Quando vocês fazem uma pergunta, ela também é transformada em uma nova estrela, com suas próprias coordenadas, e colocada nesse mapa. A busca, então, se torna uma tarefa simples: o sistema apenas precisa olhar para a estrela da sua pergunta e identificar quais são as estrelas vizinhas mais próximas. É por isso que a busca é tão poderosa: ela não procura por palavras-chave, ela procura por **vizinhanças de significado** na galáxia das suas informações.

Com os textos todos fatiados, o próximo passo é traduzir esses 'chunks' para uma língua que o computador entende de verdade: números. Esse processo, que é o coração da busca inteligente, se chama **Embedding**. Um modelo de 'embedding' — pensem nele como um tradutor universal de conceitos — pega cada 'chunk' de texto e o converte em um vetor, que é basicamente uma lista gigante de números. A grande magia é que textos com significados parecidos terão vetores (listas de números) matematicamente próximos. É como dar um CEP superpreciso para cada ideia contida nos seus documentos.

E onde a gente guarda todos esses 'CEPs' de ideias? Num lugar especial chamado **Vector Store**, ou banco de dados de vetores. Pensem nele como o armário de fichas da Ana, mas digital, super-rápido e otimizado para buscar por proximidade. Em vez de procurar por uma palavra-chave exata, a gente entrega o 'CEP' (vetor) da nossa pergunta e ele nos devolve os 'CEPs' mais parecidos que ele tem guardado em questão de milissegundos. É a base da nossa busca semântica.

Fase 2: A Execução (Recuperação e Geração)

Certo, o fichário está pronto e organizado. Agora é hora da prova, o momento em que vocês, usuários, entram em cena. Quando vocês mandam uma pergunta, o sistema inicia a fase de **Recuperação e Geração**. A primeira coisa que ele faz é pegar a pergunta de vocês — 'Qual foi o preço do pão em Paris?' — e usar o **mesmo modelo de embedding** da fase anterior para transformá-la em um vetor. É fundamental que seja o mesmo modelo para que a 'linguagem' dos CEPs seja a mesma.

Com o vetor da pergunta em mãos, o sistema vai até o Vector Store e faz a busca por similaridade. Ele basicamente pergunta: 'Ei, me devolva os **top_k** vetores mais parecidos com este aqui que eu te entreguei'. O parâmetro **top_k** é simplesmente o número de 'chunks' que vocês acham relevante recuperar. Se definirem **top_k** como 5, ele vai pegar os 5 pedaços de texto mais relevantes da sua base de dados para responder àquela pergunta específica.

Esses 'chunks' recuperados são o nosso **contexto**. Agora vem o pulo do gato: o sistema monta um novo prompt, muito mais poderoso, para o LLM. Ele junta a pergunta original de vocês com o contexto que acabou de encontrar. A instrução para o LLM (o nosso Léo) é mais ou menos assim: 'Responda à pergunta *'Qual foi o preço do pão em Paris?'* usando **apenas** as informações contidas nos seguintes textos: *[Chunk 1, Chunk 2, Chunk 3...]*'. Ao forçar o LLM a se basear no contexto fornecido, a gente reduz drasticamente a chance de ele inventar coisas, ou como chamamos na área, de 'alucinar'.

E aí está. O LLM, que já é um mestre em interpretar e gerar texto, recebe a pergunta junto com a 'cola' perfeita e factual. A tarefa dele agora não é mais 'lembrar' a resposta, mas sim sintetizar as informações que ele acabou de receber e formular uma resposta coesa, precisa e diretamente baseada nas suas fontes. Entenderam a jogada? O RAG não deixa o LLM mais 'inteligente' no sentido de conhecimento próprio. Ele o torna perfeitamente 'informado' no momento exato da pergunta. E é esse superpoder que vamos aprender a construir e a refinar ao longo deste livro.

1.3 UM RAG SIMPLES

Esta é a forma mais pura e fundamental da nossa arquitetura, o alicerce sobre o qual todas as outras variações são construídas. Pense no RAG Simples, ou 'Vanilla RAG', como o fluxo direto que desenhamos na nossa analogia inicial. É a implementação canônica da colaboração entre o pesquisador focado (o Retriever) e o escritor eloquente (o Generator).

O processo é linear e elegante na sua simplicidade. Tudo começa com a pergunta do usuário. Essa pergunta é transformada em um vetor de embedding e usada para consultar o

nosso banco de vetores. O banco, por sua vez, retorna os 'chunks' de texto mais relevantes que ele possui armazenados. Estes 'chunks' são o nosso contexto. A partir daí, a mágica acontece: a pergunta original e o contexto recuperado são combinados em um novo prompt, que é então enviado ao LLM.

A instrução para o LLM é clara e direta: "Responda a esta pergunta usando estritamente as informações fornecidas neste contexto". O LLM não precisa 'lembrar' de nada do seu treinamento massivo; sua única tarefa é sintetizar, extrair e articular uma resposta a partir do material factual que acabou de receber. Não há loops, não há autocorreção, não há memória de conversas passadas. É um sistema de 'input-processamento-output' direto, focado em uma única coisa: responder a uma pergunta de cada vez com a maior fidelidade possível às fontes fornecidas. Dominar este fluxo é a chave, pois cada técnica avançada que exploraremos a seguir nada mais é do que uma adição inteligente ou uma modificação engenhosa neste processo fundamental.

1.4 ELE PODE LEMBRAR, RAG COM MEMÓRIA

O RAG Simples que acabamos de ver é poderoso, mas tem a memória de um peixinho dourado. Ele trata cada pergunta que você faz como se fosse a primeira vez que vocês conversam. Se você perguntar 'Quem foi Santos Dumont?' e depois 'E onde ele nasceu?', um RAG Simples ficaria perdido. Ele não saberia a quem 'ele' se refere. Para construir chatbots e assistentes verdadeiramente úteis, precisamos superar essa amnésia. É aqui que entra o **RAG com Memória**.

A ideia é exatamente o que o nome sugere: dar ao nosso sistema a capacidade de **lembrar do que foi dito antes**. Em vez de descartar a conversa após cada resposta, o sistema passa a manter um histórico do diálogo, uma espécie de 'buffer de memória'. Esse histórico se torna uma nova fonte de contexto, tão importante quanto os documentos da nossa base de conhecimento.

Na prática, isso geralmente se manifesta de uma forma muito inteligente. Quando uma nova pergunta chega — como 'E onde ele nasceu?' — o sistema não a envia diretamente para o processo de busca. Primeiro, ele olha para a nova pergunta e para o histórico da conversa. Com essas duas informações, um LLM intermediário reescreve a pergunta para que ela se torne autossuficiente. O sistema transforma a pergunta ambígua 'E onde ele nasceu?' na pergunta clara e completa: 'Onde Santos Dumont nasceu?'.

É essa pergunta reescrita, agora cheia de contexto, que é usada para fazer a busca no Vector Store. O resultado é uma mudança transformadora na experiência do usuário. A interação deixa de ser uma série de perguntas e respostas isoladas e se torna um diálogo

fluido e natural. O sistema entende pronomes, resolve ambiguidades e permite que o usuário explore um tópico de forma orgânica. É o primeiro grande passo para transformar nosso sistema de busca em um verdadeiro parceiro conversacional.

1.5 O RAG AUTÔNOMO: AGENT RAG

Se o RAG com Memória deu ao nosso sistema a capacidade de lembrar, o **Agent RAG** o eleva a um novo patamar: o da autonomia e da tomada de decisão. Em vez de seguir uma sequência fixa de passos (buscar e depois gerar), um Agent RAG se comporta como um pequeno 'cérebro' que **decide dinamicamente qual a melhor ação a tomar** para responder a uma pergunta.

Imaginem um investigador. Quando ele recebe uma pergunta complexa, ele não sai buscando em todos os arquivos de uma vez. Ele pensa: "Preciso de mais informações? Qual ferramenta devo usar? Devo consultar o banco de dados interno ou preciso pesquisar na internet? Devo quebrar essa pergunta grande em outras menores?" O Agent RAG simula esse processo de raciocínio.

No coração de um Agent RAG está um LLM que atua como o 'agente' principal, munido de uma lista de **ferramentas**. Essas ferramentas podem ser diversas:

- Uma ferramenta de busca no seu Vector Store (o nosso Retriever padrão).
- Uma ferramenta para fazer buscas na web (como uma API do Google Search).
- Uma ferramenta para executar código ou cálculos.
- Uma ferramenta para consultar uma base de dados estruturada.
- E até mesmo ferramentas para interagir com o usuário e pedir mais esclarecimentos.

Quando o agente recebe uma pergunta, ele analisa o seu pedido e o seu histórico de conversa. Com base nisso, ele **decide qual ferramenta usar, se é que precisa usar alguma**. Ele pode decidir que a pergunta é simples e que ele já tem a resposta. Ou, pode concluir que precisa buscar informações em um documento específico, ou que a busca na web é mais apropriada. Depois de usar uma ferramenta, o resultado é enviado de volta para o agente, que avalia se a informação é suficiente ou se ele precisa usar outra ferramenta, ou talvez refinar a busca.

Essa capacidade de planejar, executar e iterar nas suas próprias ações transforma o RAG de um sistema reativo em um sistema proativo. O Agent RAG pode navegar por problemas complexos, consultando diferentes fontes de informação e até mesmo corrigindo seu próprio curso. É um passo crucial em direção a IAs mais inteligentes e capazes de resolver problemas do mundo real de forma mais independente.

1.6 O RAG CORRETIVO: CRAG

No mundo real, nem toda informação é ouro. Às vezes, o nosso 'Retriever' (o buscador de informações) pode trazer documentos irrelevantes, desatualizados ou até mesmo incorretos. Isso pode levar o LLM a gerar respostas ruins ou, o que é pior, a alucinar com base em um contexto falho. É para combater esse problema crítico que surge o **CRAG (Corrective RAG)**.

Pense no CRAG como um rigoroso editor de jornal que, antes de publicar uma matéria, verifica a qualidade das fontes. Ele não confia cegamente no que o repórter (o Retriever) trouxe. Ele adiciona uma camada inteligente de **autoavaliação e correção** à pipeline do RAG, garantindo que o LLM só receba informações de alta qualidade.

A principal inovação do CRAG reside em sua capacidade de **avaliar a pertinência e a qualidade dos documentos recuperados**. Pesquisas recentes na área de RAG, incluindo o artigo original que propôs o CRAG, destacam que a qualidade da recuperação é o gargalo mais comum para a performance. Um CRAG faz o seguinte:

1. Após o Retriever buscar os documentos iniciais, um módulo de **avaliador de qualidade** (geralmente um LLM menor ou um modelo treinado especificamente para essa tarefa) entra em ação.
2. Este avaliador analisa os 'chunks' recuperados e atribui uma 'pontuação de relevância' ou 'confiança'. Ele pode classificar os documentos como "altamente relevantes", "parcialmente relevantes" ou "irrelevantes".
3. Com base nessa avaliação, o CRAG toma uma decisão:
 - Se os documentos são **altamente relevantes**: Maravilha! Eles são enviados diretamente para o LLM gerador.
 - Se os documentos são **parcialmente relevantes**: O sistema pode decidir que precisa de mais informações. Ele pode, por exemplo, **expandir a busca** para documentos relacionados ou até mesmo fazer uma **nova busca na web** para encontrar informações complementares.
 - Se os documentos são **irrelevantes**: O sistema pode optar por **descartá-los** e, talvez, tentar uma abordagem diferente de busca, ou até mesmo indicar que não conseguiu encontrar informações confiáveis.

Essa capacidade de 'corrigir o curso' da recuperação é o que torna o CRAG tão poderoso. Ele não apenas busca, ele **verifica a validade da busca** e age proativamente para melhorar a qualidade do contexto antes que a geração aconteça. Isso resulta em respostas significativamente mais precisas, com menor taxa de alucinação e uma confiança muito maior na informação final apresentada ao usuário. O CRAG é um passo fundamental para tornar os

sistemas RAG mais robustos e confiáveis em ambientes de dados dinâmicos e heterogêneos.

1.7 O RAG ADAPTATIVO: ADAPTIVE RAG

A realidade é que nem toda pergunta é criada da mesma forma. Algumas são simples e diretas ("Qual a capital da França?"), enquanto outras são complexas e exigem uma investigação profunda ("Quais foram as implicações econômicas do Tratado de Versalhes para a indústria têxtil alemã?"). Um RAG Simples trata ambas da mesma maneira: busca, concatena e gera. Isso pode ser ineficiente e até prejudicial. Para perguntas simples, a busca pode ser um desperdício de tempo e recursos. Para perguntas complexas, uma única busca pode ser insuficiente. O **Adaptive RAG** resolve isso com uma dose de bom senso computacional.

O Adaptive RAG, como o próprio nome diz, **adapta sua estratégia com base na complexidade da pergunta**. Em vez de seguir um caminho único e rígido, ele primeiro analisa a pergunta e decide qual a rota mais eficiente para chegar à melhor resposta. Pense nele como um triagista experiente em um hospital: ele avalia o paciente (a pergunta) e o direciona para o tratamento correto, seja um curativo rápido ou uma cirurgia complexa.

O fluxo de trabalho geralmente envolve um componente inicial, um 'classificador' ou 'roteador', que examina a pergunta do usuário. Com base nessa análise, o sistema pode decidir por um de vários caminhos:

- **Sem Busca (No Retrieval):** Para perguntas de conhecimento geral ou conversas informais ("Como você está hoje?"), o sistema pode concluir que o LLM já sabe a resposta e não precisa de busca externa. Isso economiza tempo e poder computacional.
- **Busca Simples (Single-step Retrieval):** Para a maioria das perguntas factuais, o sistema realiza o nosso processo RAG padrão: uma única busca no banco de vetores para encontrar o contexto relevante.
- **Busca Iterativa ou em Múltiplos Passos (Multi-step Retrieval):** Para perguntas complexas, o sistema pode iniciar um ciclo de busca e raciocínio. Ele pode fazer uma busca inicial, analisar os resultados e decidir que precisa refinar a pergunta e buscar novamente, ou talvez decompor a pergunta original em sub-perguntas e buscar as respostas para cada uma delas antes de sintetizar a resposta final.

Essa capacidade de ajustar o esforço à complexidade da tarefa é o que torna o Adaptive RAG tão eficiente. Ele não desperdiça recursos em problemas fáceis e, ao mesmo tempo, tem a capacidade de aprofundar a investigação quando o desafio exige. É um sistema que não apenas responde, mas primeiro 'pensa' sobre a melhor forma de responder, trazendo um nível de inteligência e otimização muito mais próximo do raciocínio humano.

1.8 O RAG EM GRAFOS: GRAPHRAG

Até agora, nossa visão do RAG tem sido centrada em documentos: artigos, relatórios, páginas da web. Nós quebramos esses textos em pedaços e buscamos os mais relevantes. Mas e se a informação mais valiosa não estiver contida em um parágrafo, mas sim na **relação** entre diferentes pedaços de informação? É para desvendar esse conhecimento conectado que surge o **GraphRAG**.

O GraphRAG troca a nossa tradicional base de dados de vetores por uma estrutura muito mais rica: um **Grafo de Conhecimento (Knowledge Graph)**. Pense em um grafo como um mapa de relacionamentos. Em vez de 'chunks' de texto isolados, temos:

- **Nós (Nodes):** Que representam entidades como pessoas, empresas, lugares ou conceitos (ex: "Santos Dumont", "Paris", "Avião 14-Bis").
- **Arestas (Edges):** Que representam a relação entre essas entidades (ex: "Santos Dumont" *nasceu em* "Palmira", "Santos Dumont" *inventou o* "Avião 14-Bis").

Quando uma pergunta chega a um sistema GraphRAG, a abordagem é completamente diferente. Em vez de simplesmente buscar por similaridade semântica em textos, o sistema primeiro tenta entender as entidades e as relações na pergunta. Ele então traduz essa pergunta em uma consulta que **navega pelo grafo**.

Imagine a pergunta: "Quais inventores brasileiros moraram na mesma cidade que o criador do 14-Bis?". Um RAG tradicional teria muita dificuldade. Ele poderia encontrar documentos sobre inventores e sobre o 14-Bis, mas conectar os pontos seria um desafio. O GraphRAG, por outro lado, executaria uma sequência de passos lógicos:

1. Identifica que o "criador do 14-Bis" é o nó "Santos Dumont".
2. Segue a aresta *morou em* a partir do nó "Santos Dumont" para encontrar o nó "Paris".
3. Busca por outros nós com a propriedade *é um* "Inventor Brasileiro" que também tenham uma aresta *morou em* apontando para "Paris".

O resultado é uma resposta precisa, inferida a partir das conexões explícitas no conhecimento. O GraphRAG é uma técnica poderosa para domínios onde as relações são a chave, como em investigações financeiras (seguindo o dinheiro), descobertas científicas (conectando pesquisadores a artigos e a descobertas) e sistemas de recomendação inteligentes. Ele nos permite fazer um tipo de pergunta totalmente novo: não apenas "o quê?", mas principalmente "como se conecta?".

1.9 O RAG HÍBRIDO: HYBRID RAG

O **Hybrid RAG** surge como uma evolução natural quando percebemos que nenhum único mecanismo de recuperação é suficiente para lidar com a diversidade das perguntas do mundo real. Até agora, exploramos diferentes sabores de RAG: o baseado em vetores, o baseado em memória, RAG por Agentes e o em grafos. Mas a pergunta é direta: por que escolher apenas um, se podemos combinar vários e tirar o melhor de cada abordagem?

Pensem no Hybrid RAG como um time multidisciplinar. Em vez de depender só da Ana (nossa Retriever clássica de vetores), trazemos também o Pedro, que é craque em grafos, e a Júlia, especialista em busca lexical. Cada um tem uma lente diferente para enxergar a informação. O segredo do Hybrid RAG está em orquestrar essas lentes para que o sistema selecione, combine ou até mesmo faça um **reranking** inteligente dos resultados.

Na prática, o Hybrid RAG combina múltiplas estratégias de recuperação, como:

- **Busca Vetorial:** Localiza os **chunks** semanticamente mais próximos da pergunta, usando embeddings.
- **Busca Lexical (BM25, TF-IDF):** Garante que palavras-chave exatas não passem despercebidas, algo crucial em domínios jurídicos e médicos.
- **Grafos de Conhecimento:** Permitem navegar em relações explícitas, conectando conceitos que não estão no mesmo parágrafo, mas fazem parte da mesma rede de significado.

O desafio, claro, é decidir como combinar esses resultados. Alguns pipelines adotam a estratégia de **união**: trazem todos os resultados de todas as buscas e deixam o LLM fazer a síntese. Outros preferem a **interseção**: só os documentos que aparecem em mais de um método são considerados confiáveis. A abordagem mais avançada é o **reranking neural**, onde um modelo adicional atribui pesos de relevância e reorganiza os documentos, privilegiando os mais consistentes.

Um detalhe importante é o papel dos hiperparâmetros. Ajustar **top_k** em um cenário híbrido não é trivial. Pode-se definir um **top_k** específico para cada recuperador (por exemplo, 10 vetoriais, 5 lexicais e 3 de grafo) ou estabelecer um limite global após o reranking. É nessa engenharia fina que o Hybrid RAG mostra sua força: a flexibilidade de moldar a recuperação ao contexto da pergunta.

Em resumo, o Hybrid RAG é como montar um supertime em vez de apostar todas as fichas em um único jogador. Ele reduz a chance de documentos cruciais ficarem de fora, equilibra precisão e abrangência, e garante maior robustez em domínios onde os dados são heterogêneos, ambíguos ou fortemente conectados. É um passo essencial para levar o RAG

de uma solução pontual para uma arquitetura verdadeiramente universal e resiliente.

1.10 O RAG-FUSION: RECIPROCAL RANK FUSION (RRF)

O **RAG-Fusion**, também chamado de **Reciprocal Rank Fusion (RRF)**, é uma técnica refinada de combinação de resultados em buscas híbridas. Até agora, vimos que o Hybrid RAG mistura diferentes recuperadores (vetorial, lexical, grafos). Mas a grande pergunta é: como juntar essas listas de documentos de forma justa e eficiente? É aqui que o RRF brilha.

O truque do RRF é simples, mas genial. Imaginem duas filas de classificação: uma saída da busca vetorial e outra saída da busca lexical. Cada documento aparece em uma posição específica em cada fila (primeiro, segundo, quinto, etc.). O algoritmo RRF pega essas posições e calcula uma pontuação de relevância combinada. Quanto mais alto o documento estiver em qualquer uma das filas, maior será a sua chance de aparecer no resultado final. O efeito é equilibrar os dois mundos: semântica e palavras-chave.

O resultado é que documentos bem classificados em **qualquer um** dos métodos de busca recebem destaque. Isso é fundamental porque, em certos casos, a busca lexical pode capturar um detalhe exato (como o nome de uma lei ou artigo), enquanto a busca vetorial entende melhor o contexto semântico. O RRF garante que nenhum desses sinais seja perdido.

Na prática, o RAG-Fusion melhora muito a robustez do Hybrid RAG. Ele reduz a dependência em um único tipo de busca e cria uma fusão matemática mais estável, garantindo que documentos relevantes tenham chance real de aparecer no contexto entregue ao LLM. Além disso, como é uma técnica leve e independente do modelo, pode ser aplicada em escala, sem precisar de re-treinamento complexo ou modelos adicionais de reranking.

Então é isso: o RAG-Fusion é como aquele juiz justo que pega notas de diferentes jurados e monta um ranking final equilibrado. Ele não deixa que a opinião de um único jurado domine o resultado, mas também não ignora quando alguém dá uma nota muito alta. Para sistemas híbridos de recuperação, o RRF é hoje uma das formas mais práticas e eficazes de combinar evidências.

1.11 HYPOTHETICAL DOCUMENT EMBEDDING

Um dos desafios mais sutis do RAG é o que os pesquisadores chamam de "desalinhamento" entre a pergunta e o documento. Uma pergunta do usuário costuma ser curta, direta e usa palavras-chave específicas ("sintomas de dengue hemorrágica"). Um bom documento que responde a essa pergunta, no entanto, é geralmente longo, detalhado, usa uma linguagem

mais formal e pode nem mesmo conter a frase exata da pergunta (ele pode falar de "manifestações clínicas da febre hemorrágica por dengue"). Essa diferença de estilo e conteúdo pode confundir o nosso buscador. É para resolver esse problema que existe uma técnica brilhante e um tanto contraintuitiva: o **Hypothetical Document Embedding (HyDE)**.

O HyDE parte de uma premissa genial: em vez de usar a pergunta para encontrar uma resposta, que tal se a gente usasse uma **resposta ideal (mesmo que falsa)** para encontrar uma resposta real? É um truque mental que funciona surpreendentemente bem. O processo é o seguinte:

1. O sistema recebe a pergunta do usuário ("sintomas de dengue hemorrágica").
2. Em vez de enviar essa pergunta direto para o banco de vetores, ele primeiro a envia para um LLM com uma instrução simples: "Escreva um documento que responda a esta pergunta".
3. O LLM, então, gera uma **resposta hipotética**. Ele pode 'alucinar' detalhes, mas o documento gerado será semanticamente rico, estruturado como uma boa resposta e conterá o vocabulário e os conceitos relevantes (ex: "O paciente pode apresentar febre alta, dor de cabeça intensa, dor retro-orbital, além de manifestações hemorrágicas como petéquias e sangramento gengival...").
4. É este documento hipotético, e não a pergunta original, que é transformado em um vetor de embedding.
5. Finalmente, este vetor, que representa a 'essência' de uma resposta perfeita, é usado para buscar no nosso banco de vetores. A busca agora não é mais "pergunta-documento", mas sim "documento-documento", o que tende a produzir resultados muito mais relevantes.

O HyDE funciona como uma "ponte semântica". Ele traduz a intenção concisa do usuário em um exemplo detalhado do que ele espera encontrar. Mesmo que o documento hipotético contenha imprecisões, seu embedding captura o padrão geral de uma resposta relevante, tornando-o uma "isca" muito mais eficaz para "pescar" os documentos corretos na nossa base de conhecimento. É uma técnica poderosa para melhorar a precisão da recuperação, especialmente para perguntas complexas ou de nicho.