



Desafio Técnico: Inteligência Artificial

Objetivo:

Desenvolver um pipeline funcional que:

- Indexa documentos heterogêneos (PDFs com e sem texto, imagens)
- Aplique OCR quando necessário
- Gere uma base vetorial
- Implemente recuperação de informações
- Utilize uma LLM para responder perguntas com base nos documentos recuperados

Desafio:

Você deve construir uma solução composta pelos seguintes blocos:

1. Extração de conteúdo (OCR quando necessário)

A pipeline deve ser capaz de ler:

- PDFs com texto nativo
- PDFs digitalizados (ex: imagens de documentos)
- Imagens (ex: JPEG/PNG com texto)

2. Indexação

Os textos extraídos devem ser:

- Chunkados
- Convertidos em vetores semânticos
- Armazenados em um banco vetorial

3. Recuperação e geração

O sistema deve:

- Receber uma pergunta em linguagem natural
- Buscar os trechos mais relevantes na base vetorial



- Gerar uma resposta usando uma LLM, baseada nos trechos encontrados

Recursos fornecidos para indexação e recuperação -

📄 Documento Desafio SISAMB

Um pequeno conjunto de documentos:

- Código de Obras- PDF com texto
- Tabela de custos de materiais de construção - imagem

Entrega Esperada:

Repositório com:

- Código-fonte (pode ser .py ou Jupyter Notebook)
- Arquitetura do pipeline (em Markdown ou diagrama)
- README explicando execução e escolhas técnicas
- Scripts de:
 - Pré-processamento e OCR
 - Indexação
 - Consulta e resposta

Critérios de Avaliação

Critério	Peso
Extração e tratamento dos documentos	20
Indexação e uso de embeddings	20
Recuperação	20
Integração com LLM	20
Organização do código e clareza	10
Documentação e justificativa técnica	10



Prazos e Duração:

- Tempo total sugerido para resolução: até 72 horas
- Entrega via repositório Git (público ou link ZIP)

Observações:

- Avaliaremos clareza, estrutura, raciocínio e funcionalidade geral
- O foco não é uma solução de produção, mas sim viabilidade e organização técnica