

▣ Documento Técnico Detalhado: Sistema RAG Jurídico (Versão Atualizada)

1. Visão Geral e Arquitetura do Sistema

O sistema é um **Chatbot de Recuperação Aumentada por Geração (RAG)**, projetado para fornecer respostas precisas, fundamentadas e rastreáveis sobre a legislação do Estado de Santa Catarina¹. A arquitetura divide o fluxo de trabalho em duas fases críticas: **Ingestão (Offline)** e **Consulta (Online)**, priorizando a **busca híbrida** para otimizar latência e custo².

2. Fase de Ingestão e Indexação – `src/data_prep.py` (Antigo `processar_leis.py`)

Esta fase é um processo ETL (Extract, Transform, Load) que prepara os dados para busca de alta velocidade, executado uma única vez ou a cada atualização do corpus legal³.

2.1. Processamento e Normalização (Limpeza de Dados)

O script utiliza o **BeautifulSoup** para extrair o texto de arquivos HTML, seguido por uma rotina de normalização⁴:

- **Extração de Metadados:** Extração do **Ano de Publicação** e **Nome do Arquivo**⁵.
- **Remoção de Conteúdo Taxado:** Implementa uma lógica **Regex (Expressão Regular)** para identificar e remover blocos de texto que contenham indicativos de revogação (ex: tags `<s>` ou notações explícitas de revogação), garantindo que apenas o conteúdo potencialmente vigente seja indexado⁶.
- **Chunking Estratégico:** O texto é dividido em *chunks* (pedaços) de aproximadamente **1500 caracteres** (`CHUNK_SIZE_LIMIT`), com sobreposição de tokens para garantir continuidade semântica entre os pedaços⁷.

2.2. Geração de Embeddings e Gerenciamento da API

O script utiliza o SDK do Google para gerar os vetores⁸. O processo é isolado na função `generate_embeddings_batch` para gerenciamento de erros e limites de taxa⁹.

Detalhe Técnico	Implementação	Finalidade
Endpoint/Método	<code>client.models.embed_content()</code>	Chamada para gerar vetores. ¹⁰
Modelo	<code>"text-embedding-004"</code>	Modelo otimizado para tarefas de recuperação, de alta precisão

Detalhe Técnico	Implementação	Finalidade
		do Gemini. ¹¹
Parâmetro contents	Lista de <i>strings</i> de texto (lote de 100).	Otimiza a latência ao enviar múltiplos <i>chunks</i> em uma única chamada em lote. ¹²
Resiliência	Biblioteca Tenacity (@retry)	Aplica uma estratégia de <i>retry</i> exponencial. Garante que o processo continue em caso de Rate Limit ou falha temporária da API. ¹³

2.3. Indexação e Mapeamento

- FAISS Indexing:** Os vetores gerados são armazenados no índice FAISS (`faiss_index.bin`)¹⁴. O FAISS otimiza consultas de Similaridade de Cosseno (ou Distância Euclidiana) para buscar k Nearest Neighbors (kNN)¹⁵.
 - Mapeamento de ID (documents_map.json):** Estrutura de busca exata¹⁶. Associa a chave Regex (`doc_ID:Art_Número`) ao índice de *chunk* no array de vetores, permitindo acesso direto e rápido¹⁷.
-

3. Fase de Consulta (Online) – `src/rag_service.py` (Antigo `chatbot.py`)

O `src/rag_service.py`¹⁸ contém o *back-end* lógico do RAG, centralizado na função `find_document_chunks` (Busca) e `generate_response` (Geração)¹⁹.

3.1. Lógica de Busca Híbrida (`find_document_chunks`)

Rota de Busca	Lógica de Recuperação	Vantagem/Desvantagem
Rota A: ID Match (Regex)	Busca Determinística: Utiliza Regex sofisticado para extrair o ID da Lei/Artigo. Se houver sucesso, ele ignora o FAISS e recupera o <i>chunk</i> pelo <code>documents_map.json</code> ²⁰ .	Vantagem: Latência e custo mínimos, precisão de 100% no artigo exato. Desvantagem: Não responde a perguntas conceituais. ²¹
Rota B: Fallback FAISS	Busca Semântica: Ativada por falha na Rota A. A <i>query</i> é vetorizada. O vetor é usado para consultar o <code>faiss_index.bin</code> para os <i>top-K chunks</i> mais similares semanticamente (recuperação por similaridade de cosseno) ²² .	Vantagem: Responde a perguntas conceituais. Desvantagem: Maior latência (duas chamadas de API), custo maior. ²³

3.2. Chamadas de API de Geração de Conteúdo

A geração da resposta é a chamada final e de maior custo/latência²⁴.

Detalhe da Chamada	<code>embed_content (Rota B)</code>	<code>generate_content (Rotas A e B)</code>
Função	<code>client.models.embed_content(model, contents=[user_query])</code> ²⁵	<code>client.models.generate_content(model, contents=[prompt])</code> ²⁶
Modelo Usado	"text-embedding-004" ²⁷	"gemini-2.5-flash" (Melhor trade-off entre velocidade e qualidade para tarefas RAG). ²⁸
Payload de Entrada	A <i>string</i> exata da pergunta do usuário. ²⁹	Prompt Augmentation: É a concatenação da <code>SYSTEM_INSTRUCTION</code> + Contexto Recuperado + Pergunta do Usuário. ³⁰
Parâmetros	Não se aplica (apenas a <i>query</i>). ³¹	<code>temperature=0.0</code> : Configuração crucial para respostas factuais, minimiza a criatividade e a alucinação. ³²

4. Implementação da Lógica de Negócio (Regra 4)

A **Regra 4 (Tratamento de Alterações)** é a principal diferenciação do sistema, garantindo a validade jurídica das informações³³. Ela é imposta pela `SYSTEM_INSTRUCTION`³⁴.

Requisito Jurídico	Mecanismo de Implementação	Função da <code>SYSTEM_INSTRUCTION</code> (Regra de Grounding)
Artigos Alterados	O <i>chunk</i> contém o texto atualizado e o metadado da Lei Alteradora (ex: "Redação dada pela LC 789").	O LLM é obrigado a citar a Lei Alteradora como prova da vigência. ³⁵
Artigos Não Alterados (Original)	O <code>src/data_prep.py</code> injetou o Ano de Publicação no cabeçalho do <i>chunk</i> ³⁶ .	O LLM é instruído a usar o Ano de Publicação (no cabeçalho) para declarar que "a redação apresentada é a vigente desde [Ano]", superando o problema de artigos que não possuem notas de alteração. ³⁷

5. Estrutura de Arquivos e Deploy

O deploy é otimizado para **Streamlit Community Cloud**, com foco em segurança de credenciais³⁸.

Arquivo	Propósito	Dependências Críticas
<code>src/data_prep.py</code>	Ingestão, ETL e Indexação.	google-genai, faiss-cpu, beautifulsoup4, tenacity. ³⁹
<code>src/rag_service.py</code>	Core Logic (Coração do RAG). Contém a SYSTEM_INSTRUCTION, a lógica de Busca Híbrida e a chamada de Geração.	google-genai, faiss-cpu, numpy. ⁴⁰
<code>app_web.py</code>	Interface de Usuário (Streamlit) e Gerenciamento de Sessão/UI.	streamlit, <code>src/rag_service.py</code> . ⁴¹
<code>requirements.txt</code>	Lista todas as bibliotecas necessárias para o ambiente de deployment.	Usado pelo Streamlit Cloud/Ambiente de Nuvem. ⁴²
<code>.streamlit/secrets.toml</code>	CRÍTICO PARA DEPLOY. Arquivo que armazena a GEMINI_API_KEY de forma segura.	Acesso via <code>st.secrets["GEMINI_API_KEY"]</code> no <code>app_web.py</code> . ⁴³