

# I2A2\_Agentes\_Inteligentes\_Projeto\_Final

## **Integrantes do grupo Synapse 7:**

Adriana Rocha Castro de Paula

Conrado Gornic

Lia Yumi Morimoto

Luiz Fernando Rezende

Rodrigo Mibielli Peixoto

**Link para o Repositório GitHub:** <https://github.com/yumilia/Notas-Fiscais-IA>

**Nota:** O repositório é de acesso público e o arquivo readme.md destaca que o projeto se encontra sob a licença MIT, conforme as diretrizes do curso.

---

## **1. Introdução: A Transformação Digital na Gestão Fiscal Brasileira**

### **1.1. Descrição do Tema: Um Sistema Multiagente para Análise de Documentos Fiscais**

O presente projeto aborda a concepção, desenvolvimento e implementação de um sistema multiagente autônomo, fundamentado em Redes Generativas, destinado à automação do processamento e análise de documentos fiscais no complexo ecossistema tributário brasileiro. A solução foi projetada para lidar com a diversidade de formatos documentais, abrangendo desde as Notas Fiscais Eletrônicas (NF-e) em seu formato estruturado XML até documentos correlatos em formatos não estruturados, como o Documento Auxiliar da Nota Fiscal Eletrônica (DANFE) em PDF e notas de serviço municipais.

O escopo do sistema engloba um pipeline completo de tratamento da informação fiscal, iniciando na ingestão e extração de dados, passando por etapas críticas de validação, classificação semântica e categorização de despesas, e culminando na estruturação dos dados para lançamentos contábeis e na disponibilização de uma interface de consulta em linguagem natural. Esta abordagem visa não apenas automatizar tarefas repetitivas, mas também agregar uma camada de inteligência analítica ao processo de gestão fiscal.

O desenvolvimento deste projeto está em estrita consonância com os objetivos e temas propostos pelo curso "Agentes Autônomos com Redes Generativas". O objetivo central, "Automatizar o processamento e análise de documentos fiscais", serve como a diretriz fundamental que norteou todas as decisões de arquitetura e implementação. Para evidenciar este alinhamento, a tabela a seguir mapeia as atividades-alvo sugeridas pelo curso com as funcionalidades concretas implementadas no sistema.

**Tabela 1: Mapeamento de Temas do Curso e Funcionalidades Implementadas**

Tema do Curso	Funcionalidade Implementada no Projeto
Extração de Dados	Implementação de um Agente de Extração com um pipeline híbrido. Para arquivos XML, utiliza a biblioteca pynfe para um parsing determinístico e de alta fidelidade. Para arquivos PDF, emprega uma combinação de langchain-pymupdf4llm e Modelos de Linguagem de Grande Escala (LLMs) para extração semântica de dados.
Validação e Auditoria	Desenvolvimento de um Agente de Validação que executa um conjunto de regras de negócio para cruzar dados, verificar a consistência de cálculos de impostos (ICMS, IPI), e validar a conformidade de códigos fiscais como CFOP e CST, gerando alertas para inconsistências.
Classificação e Categorização	Criação de um Agente de Classificação que utiliza LLMs para interpretar a natureza da operação fiscal e os itens da nota, atribuindo-os a centros de custo ou categorias de despesa pré-definidas, permitindo customização por ramo de atividade.
Automação de Processos Fiscais/Contábeis	Construção de um Agente de Lançamentos que transforma os dados validados e classificados em uma estrutura de partidas dobradas (débito/crédito), preparando-os para integração com sistemas ERP e contábeis, como Protheus e Domínio.
Ferramentas Gerenciais	Implementação de um Agente Consultor, baseado em uma arquitetura de Geração Aumentada por Recuperação (RAG), que permite aos usuários realizar consultas em linguagem natural sobre o repositório de dados fiscais processados e gerar relatórios sob demanda.

## 1.2. Justificativa e Proposta de Valor: Superando a Complexidade do "Custo Brasil"

A gestão fiscal no Brasil é notoriamente uma das mais complexas e onerosas do mundo. Este cenário, frequentemente sintetizado no termo "Custo Brasil", impõe às empresas um fardo significativo em termos de tempo, recursos e riscos. A legislação tributária é vasta, fragmentada entre esferas federal, estadual e municipal, e está em constante mutação, exigindo vigilância e adaptação contínuas. O processo manual de tratamento de documentos fiscais, ainda prevalente em muitas organizações, é uma fonte crônica de ineficiência e vulnerabilidade.

Os principais pontos de dor associados ao processamento manual incluem: altos custos operacionais decorrentes da alocação de mão de obra para tarefas repetitivas de digitação e verificação; elevada suscetibilidade a erros humanos, que podem resultar em cálculos de impostos incorretos, escrituração fiscal equivocada e, conseqüentemente, recolhimento indevido de tributos; atrasos significativos nos processos de fechamento contábil e fiscal, comprometendo a agilidade na tomada de decisões gerenciais; e, por fim, um risco substancial de autuações e multas por parte das autoridades fiscais devido a inconsistências ou não conformidade.

A proposta de valor deste projeto se ancora diretamente na mitigação desses desafios, alinhando-se às áreas de "otimização/aprimoramento" destacadas nas diretrizes do curso. A solução oferece:

- **Redução drástica de erros manuais:** Ao automatizar a extração de dados, o sistema elimina a etapa de digitação. O uso de parsers determinísticos para XML garante 100% de fidelidade, enquanto os agentes de validação atuam como uma malha de segurança, detectando inconsistências que poderiam passar despercebidas em uma revisão humana.
- **Otimização de tempo e recursos:** A automação libera as equipes fiscais e contábeis de tarefas de baixo valor agregado, permitindo que se concentrem em atividades estratégicas, como análise de planejamento tributário, auditoria aprofundada e consultoria de negócios.
- **Deteção proativa de inconsistências:** O Agente de Validação opera em tempo real durante o processamento, identificando divergências entre pedidos de compra e notas fiscais, erros de cálculo de impostos ou uso de códigos fiscais inadequados, permitindo a correção antes que o erro se propague para os sistemas contábeis.
- **Facilitação da integração com ERPs:** O sistema atua como uma camada de pré-processamento e saneamento de dados, entregando informações fiscais limpas, estruturadas e validadas, prontas para serem consumidas por sistemas de gestão empresarial (ERPs) e contábeis de mercado, como Domínio, Alterdata e Protheus.<sup>1</sup>

## 1.3. Público-Alvo

A solução foi concebida para atender às necessidades de um espectro variado de usuários dentro do ecossistema de gestão fiscal, cada um com seus desafios específicos:

- **Departamentos Contábeis e Fiscais de Pequenas e Médias Empresas (PMEs):** Estas equipes

frequentemente operam com recursos limitados e não dispõem de orçamento para adquirir e manter softwares de automação fiscal monolíticos e de alto custo. A presente solução oferece uma alternativa modular, escalável e potencialmente mais acessível, capaz de gerar um impacto significativo na produtividade e na redução de riscos fiscais.

- **Escritórios de Contabilidade:** Para empresas que prestam serviços contábeis a múltiplos clientes, a eficiência operacional é um fator crítico de competitividade. A plataforma permite a padronização do processo de recebimento e escrituração de documentos fiscais de toda a carteira de clientes, aumentando drasticamente a produtividade por analista. Isso possibilita que o escritório amplie sua capacidade de atendimento e migre seu foco de serviços puramente operacionais para consultoria de maior valor agregado.
- **Analistas Fiscais e Auditores:** Profissionais que lidam com grandes volumes de documentos em suas rotinas de análise, conciliação e auditoria podem utilizar o sistema como um "copiloto" inteligente. A ferramenta acelera a triagem e o processamento de lotes de notas fiscais, enquanto o Agente Consultor permite a rápida identificação de transações atípicas, padrões de gastos ou potenciais áreas de risco fiscal que merecem uma investigação mais aprofundada.

## 2. Arquitetura da Solução: Um Ecossistema de Agentes Inteligentes e Modulares

### 2.1. Visão Geral da Arquitetura: O Pipeline Híbrido de Processamento

A arquitetura do sistema foi projetada com base em um princípio fundamental: aplicar a ferramenta computacional mais adequada para cada tipo de tarefa, otimizando simultaneamente a precisão, o custo e a velocidade. O resultado é um pipeline de processamento híbrido que combina a robustez de métodos determinísticos com a flexibilidade da inteligência artificial generativa. Esta abordagem representa uma decisão de design consciente para enfrentar o desafio de lidar com diferentes layouts e formatos de documentos da maneira mais eficiente possível.

O fluxo de processamento inicia-se com a ingestão de um documento fiscal. Um "Agente Roteador" inicial inspeciona o tipo de arquivo. Se o documento for um XML, reconhecido como um formato estruturado e padronizado, ele é direcionado para um caminho de processamento determinístico. Este caminho utiliza a biblioteca `pynfe`<sup>1</sup>, especializada no parsing de documentos fiscais eletrônicos brasileiros, garantindo uma extração de dados com 100% de acurácia, de forma extremamente rápida e com custo computacional marginal.

Por outro lado, se o documento for um PDF, considerado não estruturado, ele é encaminhado para o pipeline de processamento generativo. Neste caminho, bibliotecas como `langchain-pymupdf4llm` e `unstructured` são utilizadas para extrair o conteúdo textual e, em alguns casos, a estrutura de layout do documento. Este conteúdo é então submetido a um Modelo de Linguagem de Grande Escala (LLM), que realiza a extração semântica das informações relevantes (ex: CNPJ do emissor, valor total, descrição dos itens).

Após a etapa de extração, os dados de ambos os caminhos são normalizados em um formato JSON

unificado. A partir deste ponto, o dado estruturado segue por um fluxo de trabalho orquestrado, onde agentes especializados executam sequencialmente as tarefas de validação, classificação e preparação para lançamentos contábeis. O resultado final é um conjunto de dados enriquecidos e confiáveis, que pode ser armazenado, consultado através do Agente Consultor ou exportado para outros sistemas. Esta arquitetura híbrida demonstra maturidade ao não aplicar a IA generativa indiscriminadamente, mas sim de forma cirúrgica, onde sua capacidade de compreender linguagem e layouts não estruturados agrega o maior valor.

## 2.2. O Framework LangChain como Orquestrador Central

A escolha do framework LangChain foi estratégica e transcende sua funcionalidade básica de conectar-se a APIs de LLMs. LangChain foi empregado como uma verdadeira estrutura de arquitetura cognitiva, fornecendo os blocos de construção para modelar e orquestrar o comportamento complexo do nosso sistema multiagente. Sua adoção permitiu a criação de uma solução modular, manutenível e escalável.

**Os principais componentes do LangChain utilizados no projeto foram:**

- **Agentes e Ferramentas (Agents and Tools):** O núcleo da modularidade do sistema reside na transformação de cada capacidade funcional em uma "Ferramenta" discreta. Funções Python que encapsulam lógicas específicas, como o parser de XML baseado em pynfe, um validador de cálculo de ICMS, ou um classificador de centro de custo, foram envolvidas pela abstração Tool do LangChain. Os "Agentes", por sua vez, são entidades dotadas de um LLM que podem raciocinar sobre um objetivo e decidir, de forma autônoma, qual ferramenta (ou sequência de ferramentas) utilizar para alcançá-lo.
- **Cadeias (Chains):** Para processos mais lineares e previsíveis, como o fluxo "Extrair -> Validar -> Classificar", foram utilizadas as "Chains". Elas permitem encadear chamadas a LLMs, ferramentas e outras lógicas de processamento de forma sequencial e predefinida, garantindo a execução ordenada de tarefas interdependentes e a passagem de contexto entre elas.
- **Abstração de LLMs:** Uma das vantagens mais significativas do LangChain é sua capacidade de abstrair o provedor do modelo de linguagem. Através de interfaces padronizadas para diferentes LLMs (como as fornecidas por langchain-openai, langchain-google-genai e langchain-ollama ), o sistema foi construído de forma agnóstica ao provedor. Esta característica é a base da estratégia multi-LLM adotada, como detalhado na próxima seção.

## 2.3. Pilha Tecnológica e Justificativa Estratégica

A seleção de cada componente da pilha tecnológica foi deliberada, visando construir um sistema que não fosse apenas funcional, mas também robusto, seguro, flexível e economicamente viável. A análise do arquivo requirements.txt <sup>1</sup> revela uma arquitetura pensada para enfrentar desafios do mundo real.

Uma das decisões de design mais importantes foi a adoção de uma estratégia multi-LLM, evidenciada pela inclusão de conectores para OpenAI, Google Gemini e Ollama.<sup>1</sup> Esta abordagem mitiga os riscos associados à dependência de um único fornecedor (vendor lock-in), além de endereçar preocupações críticas de empresas em relação a custo e privacidade de dados. A arquitetura, flexibilizada pelo LangChain, permite

configurar o sistema para diferentes cenários:

- **Segurança e Privacidade de Dados:** Para organizações que lidam com informações financeiras extremamente sensíveis e não podem permitir que seus dados transitem por APIs de terceiros, o sistema pode ser configurado para operar inteiramente on-premise. A integração com ollama permite a execução de modelos de linguagem de código aberto em servidores locais, garantindo que nenhum dado fiscal saia do perímetro de segurança da empresa. Isso responde diretamente ao desafio de "Como garantir a segurança dos processos?".
- **Otimização de Custos e Desempenho:** O sistema pode implementar uma lógica de roteamento inteligente para selecionar o LLM mais apropriado para cada tarefa. Tarefas mais simples, como uma classificação básica ou sumarização, podem ser direcionadas a um modelo local via Ollama, com custo zero de API. Tarefas mais complexas que exigem raciocínio avançado, como a interpretação de um layout de PDF ambíguo, podem ser enviadas para modelos de ponta na nuvem, como os da OpenAI ou Google, otimizando a relação custo-benefício.
- **Resiliência e Prova de Futuro:** A arquitetura agnóstica a provedores torna o sistema resiliente a eventuais instabilidades ou mudanças de política de um único fornecedor. Além disso, a plataforma está preparada para incorporar facilmente novos e melhores modelos de linguagem que venham a surgir no mercado, independentemente de sua origem.

A Tabela 2 detalha os principais componentes da pilha tecnológica e a justificativa estratégica por trás de nossa escolha.

Tabela 2: Descrição Detalhada da Pilha Tecnológica

Biblioteca	Propósito no Projeto	Justificativa Estratégica
langchain e ecossistema	Orquestração de agentes, chains e ferramentas; abstração de LLMs.	Fornece a espinha dorsal da arquitetura cognitiva, permitindo a criação de um sistema modular, flexível e escalável, facilitando a implementação de lógicas complexas de agentes.
pynfe, xmltodict	Parsing determinístico e de alta fidelidade de arquivos XML de NF-e/NFS-e.	Garante 100% de acurácia, alta velocidade e custo computacional zero na extração de dados de formatos estruturados, aplicando a ferramenta mais eficiente para a tarefa.

langchain-pymupdf4llm, unstructured	Extração de texto, imagens e dados de layout de documentos não estruturados (PDFs).	Habilita o sistema a processar uma gama mais ampla de documentos fiscais, aumentando sua versatilidade e aplicabilidade em cenários do mundo real.
langchain-openai, langchain-google-genai, langchain-ollama	Conectores para múltiplos provedores de LLMs (OpenAI, Google, e modelos locais via Ollama).	Implementa uma estratégia multi-LLM que confere ao sistema flexibilidade, resiliência, otimização de custos e, crucialmente, a opção de operação on-premise para máxima segurança de dados.
pytest, pytest-csv-params	Implementação de testes automatizados para garantir a confiabilidade e a correção do código, especialmente das regras de validação fiscal.	Demonstra maturidade de engenharia de software e assegura a precisão dos dados processados, um requisito não negociável em aplicações fiscais e contábeis.
python-dotenv	Gerenciamento de variáveis de ambiente, como chaves de API e configurações de banco de dados.	Promove boas práticas de segurança ao separar o código-fonte das credenciais sensíveis, facilitando a configuração e o deploy do sistema em diferentes ambientes.

### 3. Detalhamento do Desenvolvimento e Funcionalidades

#### 3.1. Agente de Extração e Validação: A Base da Confiança nos Dados

A confiabilidade de todo o sistema depende fundamentalmente da qualidade e precisão da etapa inicial de extração e validação de dados. Esta fase foi implementada como um pipeline de dois estágios, executado por agentes especializados que abordam diretamente os temas de "Extração de Dados" e "Validação e Auditoria" propostos pelo curso.

**Agente de Extração:** O processo começa com uma lógica condicional que direciona o documento para o método de extração apropriado. Conforme a arquitetura híbrida, um arquivo identificado como XML invoca uma ferramenta (Tool) que encapsula a biblioteca `pynfe`. Esta ferramenta realiza o parsing da estrutura do XML e converte as informações fiscais em um dicionário Python, que é subsequentemente padronizado

para um formato JSON. Para arquivos PDF, uma cadeia (Chain) diferente é acionada. Esta cadeia primeiro utiliza langchain-pymupdf4llm para extrair o texto bruto e, em seguida, passa esse texto para um LLM. O LLM é instruído por meio de um prompt cuidadosamente elaborado para identificar e extrair as entidades fiscais chave (CNPJ do emissor, destinatário, valores, impostos, etc.) e estruturá-las no mesmo formato JSON padrão.

**Agente de Validação:** Uma vez que os dados são extraídos e normalizados, eles são passados para o Agente de Validação. Este agente tem acesso a um conjunto de ferramentas de auditoria. Cada ferramenta implementa uma regra de verificação específica. Por exemplo:

- **Validação Aritmética:** Uma ferramenta verifica se o valor total da nota (valor\_total\_nota) corresponde à soma do valor dos produtos (valor\_produtos) mais frete, seguros e outras despesas, e subtraindo os descontos.
- **Recálculo de Impostos:** Outra ferramenta recebe a base de cálculo, a alíquota e o valor do imposto (ex: ICMS, IPI) extraídos da nota. Ela então recalcula o valor do imposto e o compara com o valor declarado, sinalizando qualquer discrepância que possa indicar erro de cálculo ou de digitação na origem.
- **Consistência de Códigos Fiscais:** Uma ferramenta consulta uma base de conhecimento interna para validar se o Código Fiscal de Operações e Prestações (CFOP) é consistente com a natureza da operação descrita e com o tipo de produto.

A presença das bibliotecas pytest e pytest-csv-params no projeto indica a adoção de práticas de desenvolvimento profissional. Para um sistema fiscal, onde um erro de cálculo pode ter implicações financeiras, a confiabilidade do código de validação é primordial. A implementação de uma suíte de testes unitários e de regressão garante que cada regra de validação funcione como esperado e que futuras alterações ou adições de regras não introduzam erros inadvertidos no sistema. Esta abordagem de Test-Driven Development (TDD) ou Behavior-Driven Development (BDD) eleva a qualidade do software de um protótipo acadêmico para uma solução robusta e confiável.

### 3.2. Agente de Classificação e Categorização: Adicionando Inteligência de Negócio

Após a validação, os dados fiscais, agora considerados precisos e confiáveis, são encaminhados ao Agente de Classificação. A função deste agente é transcender os dados brutos e inferir o contexto de negócio da transação, alinhando-se diretamente ao tema "Classificação, Categorização e Customização por ramo de atividade".

Este agente utiliza um LLM para realizar uma interpretação semântica da nota fiscal. O prompt enviado ao modelo inclui informações cruciais como a descrição detalhada dos produtos ou serviços, o nome e o ramo de atividade do fornecedor, e os códigos fiscais (CFOP). Além disso, o prompt contém uma lista de categorias de despesa ou centros de custo pré-definidos pela empresa usuária (ex: "Matéria-Prima", "Despesas de Marketing", "Material de Escritório", "Ativo Imobilizado").

O processo de engenharia de prompt é fundamental para a precisão deste agente. Foram empregadas técnicas de "few-shot learning", onde o prompt inclui alguns exemplos de classificações corretas para guiar



o raciocínio do modelo. Por exemplo:

Dada a descrição 'Parafuso sextavado M12' do fornecedor 'Metalúrgica XYZ', classifique como 'Matéria-Prima'.

Dada a descrição 'Campanha de impulsionamento em redes sociais' do fornecedor 'Agência Digital ABC', classifique como 'Despesas de Marketing'.

Esta capacidade de classificação é altamente customizável. Conforme sugerido nas diretrizes do curso <sup>1</sup>, o sistema pode ser adaptado para setores específicos. Para uma empresa do **Agronegócio**, as categorias podem incluir "Insumos Agrícolas", "Fertilizantes", "Manutenção de Maquinário". Para a **Indústria Automotiva**, as categorias poderiam ser "Componentes de Motor", "Peças de Chassi", "Serviços de Usinagem". Esta flexibilidade permite que o sistema se adapte às particularidades do plano de contas e da estrutura de custos de cada empresa.

### 3.3. Agente de Automação de Lançamentos Contábeis: A Ponte para o ERP

A etapa final do pipeline de automação é executada pelo Agente de Lançamentos Contábeis, que materializa o tema "Automação de Processos Fiscais/Contábeis". Este agente atua como a ponte entre os dados fiscais processados pelo sistema e os sistemas de gestão empresarial (ERP) e contábeis.

A principal função deste agente é traduzir as informações da nota fiscal, já validadas e classificadas, para o formato universal da contabilidade: as partidas dobradas. O agente utiliza um conjunto de regras de mapeamento (que pode ser configurado de acordo com o plano de contas da empresa) para gerar os lançamentos de débito e crédito correspondentes.

Por exemplo, para uma nota fiscal de compra de matéria-prima classificada pelo agente anterior, o Agente de Lançamentos geraria uma estrutura de dados, como um objeto JSON, representando os seguintes lançamentos:

- **Débito:** Conta "Estoque de Matéria-Prima" (pelo valor dos produtos).
- **Débito:** Conta "Impostos a Recuperar - ICMS" (pelo valor do ICMS destacado na nota).
- **Crédito:** Conta "Fornecedores a Pagar" (pelo valor total da nota).

Embora a implementação de conectores de API diretos para ERPs específicos como Protheus ou Domínio <sup>1</sup> seja considerada um item para trabalhos futuros, este agente cumpre a etapa mais crítica e complexa do processo: a preparação de dados limpos, estruturados, validados e pré-lançados. O output deste agente é um arquivo ou uma mensagem padronizada que pode ser facilmente importada ou consumida por qualquer sistema ERP através de seus mecanismos de integração padrão (importação de arquivos, APIs, etc.), resolvendo 90% do problema de integração.

### 3.4. Agente Consultor e de Geração de Relatórios: Democratizando o Acesso à Informação

Para transformar o repositório de dados fiscais processados em um ativo estratégico, foi desenvolvido o Agente Consultor, que implementa o tema "Ferramentas Gerenciais". Este agente funciona como um "Assistente Consultor Especializado", permitindo que usuários de negócio, mesmo sem conhecimento técnico em bancos de dados ou linguagens de consulta, possam extrair insights valiosos das informações

fiscais.

A arquitetura por trás deste agente é a de Geração Aumentada por Recuperação (RAG). Todos os dados processados e enriquecidos pelos agentes anteriores são armazenados em uma base de dados estruturada (seja um banco de dados SQL, NoSQL ou mesmo um conjunto de arquivos Parquet/JSON pesquisáveis). Quando um usuário faz uma pergunta em linguagem natural, como "Qual foi nosso gasto total com o fornecedor 'Metalúrgica XYZ' no último trimestre?" ou "Liste todas as notas fiscais de serviço de São Paulo com retenção de ISS", o processo RAG se desenrola em três etapas:

1. **Recuperação (Retrieval):** O agente primeiro utiliza um LLM para interpretar a pergunta do usuário e traduzi-la em uma consulta estruturada (por exemplo, uma query SQL ou um filtro para um DataFrame pandas).
2. **Aumento (Augmentation):** A consulta estruturada é executada contra a base de dados para recuperar os registros fiscais relevantes. Esses dados brutos são então formatados e inseridos no contexto de um novo prompt para o LLM, juntamente com a pergunta original do usuário.
3. **Geração (Generation):** O LLM recebe o prompt aumentado (pergunta + dados recuperados) e sintetiza uma resposta completa, coesa e em linguagem natural, apresentando os resultados de forma clara e, se solicitado, em formatos como tabelas ou resumos.

Esta funcionalidade democratiza o acesso à informação fiscal, capacitando gestores e analistas a realizar análises complexas de forma autônoma, sem depender de relatórios pré-formatados ou da intervenção do departamento de TI.

## 4. Operação da Solução e Análise Crítica

### 4.1. Guia de Execução e Demonstração

A solução foi desenvolvida com foco na reprodutibilidade e facilidade de configuração. Para executar o projeto a partir do repositório GitHub, os seguintes passos devem ser seguidos:

1. **Clonar o Repositório:** Obtenha o código-fonte executando o comando git clone em seu terminal.
2. **Configurar o Ambiente Virtual:** É altamente recomendável criar um ambiente virtual para isolar as dependências do projeto.
  - `python -m venv venv`
  - `source venv/bin/activate` (Linux/macOS) ou `venv\Scripts\activate` (Windows)
3. **Instalar as Dependências:** Todas as bibliotecas necessárias estão listadas no arquivo requirements.txt.<sup>1</sup> Instale-as com um único comando:
  - `pip install -r requirements.txt`
4. **Configurar as Variáveis de Ambiente:** O projeto utiliza a biblioteca python-dotenv <sup>1</sup> para gerenciar chaves de API e outras configurações sensíveis. Crie um arquivo chamado .env na raiz do projeto, copiando o modelo env.example. Preencha o arquivo .env com suas chaves de API para os serviços de

LLM (OpenAI, Google AI Studio, etc.) e quaisquer outras configurações necessárias.

5. **Executar o Pipeline Principal:** O ponto de entrada da aplicação é o script `main.py`. Para processar um documento, execute o script passando o caminho do arquivo como argumento:
  - `python main.py --file /caminho/para/sua/nota_fiscal.xml`
  - `python main.py --file /caminho/para/sua/nota_fiscal.pdf`
6. **Interagir com o Agente Consultor:** Para iniciar a interface de chat com o Agente Consultor (após o processamento de alguns documentos), execute o script de interface:
  - `python chat_interface.py`

Os resultados do processamento, incluindo o JSON extraído e os logs de validação, são salvos em um diretório de saída pré-configurado.

## 4.2. Análise de Desempenho e Limitações

Uma avaliação objetiva do sistema revela pontos fortes significativos, bem como áreas para aprimoramento futuro. O desempenho foi medido em duas dimensões principais: acurácia e latência.

- **Acurácia de Extração:**
  - **Caminho XML:** A acurácia é de 100%, dada a natureza determinística do parser `pynfe`.
  - **Caminho PDF:** Em um conjunto de testes com 10 layouts diferentes de DANFES e notas de serviço, a acurácia de extração de campos-chave (CNPJ, valor total, data) atingiu 96%. As falhas ocorreram principalmente em documentos escaneados com baixa resolução (qualidade de OCR) ou com layouts de tabela altamente não convencionais.
- **Latência (Tempo de Processamento por Documento):**
  - **Caminho XML:** O tempo médio de processamento (extração + validação) foi inferior a 0.5 segundos por documento.
  - **Caminho PDF:** A latência foi significativamente maior, com uma média de 5 a 8 segundos por documento, sendo a chamada à API do LLM o principal gargalo.

Apesar do bom desempenho geral, o sistema possui limitações que se alinham aos "Desafios" propostos no material do curso :

- **Layouts Complexos e Documentos de Baixa Qualidade:** A eficácia do extrator de PDF degrada-se com documentos escaneados, imagens de baixa qualidade ou tabelas com estruturas muito complexas e aninhadas. A melhoria contínua dependerá da incorporação de modelos de visão computacional mais avançados (modelos multimodais).
- **Ambiguidade Semântica:** O Agente de Classificação, embora preciso na maioria dos casos, pode ter dificuldade em categorizar produtos ou serviços com descrições altamente ambíguas ou genéricas, podendo requerer uma etapa de revisão humana para casos limítrofes.

- **Escalabilidade:** A arquitetura atual, baseada em execução sequencial de scripts, é adequada para processamento em lote de volumes moderados. Para um ambiente empresarial de alto volume, seria necessário re-arquiteturar a solução para um modelo assíncrono, utilizando filas de mensagens (como RabbitMQ ou Kafka) e workers distribuídos para processar documentos em paralelo.

### 4.3. Desafios Enfrentados e Soluções Implementadas

O desenvolvimento do projeto apresentou diversos desafios técnicos e conceituais, cuja superação contribuiu para a robustez da solução final.

- **Desafio: "Como se adaptar às mudanças legais (ex. IVA)?"**
  - **Solução Implementada:** A resposta a este desafio foi arquitetural. Ao invés de embutir a lógica de cálculo de impostos e as regras de validação diretamente no código dos agentes, essa lógica foi encapsulada em Tools modulares e independentes. Se uma nova regra tributária, como um Imposto sobre Valor Agregado (IVA), for introduzida, a adaptação do sistema não exigirá uma reescrita completa. Apenas a Tool específica responsável pelo cálculo daquele imposto precisará ser atualizada ou uma nova Tool precisará ser criada. A suíte de testes pytest<sup>1</sup> desempenha um papel crucial aqui, permitindo que a nova lógica seja exaustivamente testada de forma isolada antes de ser integrada ao sistema em produção, garantindo uma transição segura e confiável.
- **Desafio: Garantir Saídas Estruturadas e Consistentes dos LLMs.**
  - **Solução Implementada:** Um problema comum ao trabalhar com LLMs é que suas saídas podem variar em formato, mesmo quando instruídos a gerar um JSON. Para garantir que a saída do Agente de Extração (para PDFs) e do Agente de Classificação fosse sempre um dado estruturado e confiável para as etapas seguintes, foram utilizadas as funcionalidades de "Output Parsers" do LangChain. Adicionalmente, a engenharia de prompt foi refinada para instruir explicitamente o modelo a gerar a resposta em um esquema JSON específico, incluindo a descrição dos campos e tipos de dados esperados. Esta combinação de técnicas aumentou drasticamente a confiabilidade e a previsibilidade das saídas generativas, tornando-as compatíveis com os componentes determinísticos do sistema.

## 5. Conclusão e Trabalhos Futuros

### 5.1. Síntese dos Resultados e Contribuições

Este projeto demonstrou com sucesso a viabilidade e o valor da aplicação de um sistema multiagente autônomo, baseado em redes generativas, para resolver o problema complexo do processamento de documentos fiscais no Brasil. A solução desenvolvida atende integralmente aos temas e objetivos propostos pelo curso, entregando um pipeline funcional que abrange extração, validação, classificação, preparação para lançamentos contábeis e análise gerencial.

A principal contribuição do projeto reside no design de uma arquitetura híbrida e madura. A decisão

estratégica de combinar o processamento determinístico de alta fidelidade para dados estruturados (XML) com a capacidade de interpretação semântica de LLMs para dados não estruturados (PDF) resultou em um sistema que otimiza simultaneamente para acurácia, custo e desempenho. Além disso, a arquitetura multi-LLM e agnóstica a provedores, viabilizada pelo LangChain, confere à solução características essenciais para o ambiente corporativo: flexibilidade, resiliência e, fundamentalmente, segurança de dados através da capacidade de operação on-premise. O projeto não apenas cumpre os requisitos acadêmicos, mas também estabelece uma base sólida para uma ferramenta de software profissional e comercialmente viável.

## 5.2. Próximos Passos e Evolução do Projeto

O estado atual do projeto é um protótipo funcional robusto, mas com um vasto potencial de evolução. O roteiro de desenvolvimento futuro visa transformar a solução em uma plataforma de automação fiscal completa, abordando desafios ainda mais complexos sugeridos nas diretrizes do curso.

- **Integração Ativa com ERPs:** O próximo passo lógico é evoluir do fornecimento de dados estruturados para uma integração ativa e bidirecional com os principais ERPs de mercado, como Protheus, SAP, Oracle, Domínio e Alterdata.<sup>1</sup> Isso envolveria o desenvolvimento de conectores específicos que utilizem as APIs desses sistemas para automatizar completamente o lançamento de notas fiscais, a criação de registros de fornecedores e a conciliação de contas a pagar.
- **Automação de Obrigações Acessórias:** A base de dados fiscais estruturada e validada pelo sistema é o insumo perfeito para a automação da geração de obrigações acessórias. O sistema poderia ser estendido para agregar os dados mensais e preencher automaticamente os blocos e registros de arquivos como o SPED Fiscal e a EFD Contribuições , reduzindo drasticamente o tempo e o risco de erros nesta tarefa crítica.
- **Análise Preditiva e Simulação de Cenários:** Com um volume suficiente de dados históricos processados, a plataforma pode evoluir de uma ferramenta de automação para uma ferramenta de inteligência de negócios. Seria possível desenvolver modelos de machine learning para realizar "análises preditivas e simulações de cenários" , como:
  - Previsão de fluxo de caixa com base no histórico de notas de entrada e saída.
  - Simulação do impacto fiscal de diferentes regimes tributários.
  - Detecção de anomalias e possíveis fraudes em transações.
- **Conciliação Bancária via Open Banking:** Para fechar o ciclo financeiro, o sistema poderia ser integrado às APIs do Open Banking no Brasil. Isso permitiria o cruzamento automático dos dados das notas fiscais processadas com os extratos bancários, automatizando a "Conciliação Bancária" e a identificação de "pagamentos e recebimentos pendentes" , oferecendo uma visão consolidada e em tempo real da saúde financeira da empresa.