

MODELAGEM DE RASTREADOR PARA PORTAIS DA ÁREA DO DIREITO

**Eduardo Pavani Palharini¹; Cristiano Becker Isaia²; Rosane Leal da Silva³;
Robertson Ebling dos Santos⁴; Alexandre de Oliveira Zamberlan⁵**

RESUMO

Este artigo busca apresentar a ideia geral e as modelagens funcionais e estruturais de um sistema de rastreamento de conteúdo em portais da área do Direito. O sistema de rastreamento tem como objetivo qualificar as pesquisas de temas e/ou assuntos relacionados à área do Direito, em sites (fontes) pré-definidas pelo usuário, trazendo resultados relevantes e mais significativos. E todos os resultados encontrados devem ser armazenados no sistema proposto, sob a gestão do usuário que realizou a pesquisa. Para a modelagem, estão sendo utilizados metodologia e ferramentas, como a metodologia SCRUM, associada a técnica *Kanban* de gestão de atividades, a ferramenta ASTAH de diagramação *Unified Modeling Language* (UML). Os resultados deste trabalho são a especificação dos principais aspectos estruturais e de todos os aspectos funcionais do sistema em construção. A pesquisa e o desenvolvimento estão sob a gestão do Laboratório de Práticas da Computação UFN.

Palavras-chave: Inteligência Artificial; Sistema de Informação; Web Crawler.

Eixo Temático: Tecnologia, Inovação e Desenvolvimento Sustentável (TIDS).

1. INTRODUÇÃO

Um rastreador da Web (*Web crawler*) é um software robô denominado *bot* que navega sistematicamente em servidores Web, com a finalidade de recuperar e indexar conteúdos (páginas), a partir de critérios de pesquisa informados em mecanismos de busca. Esses mecanismos de pesquisa que são associados a rastreadores, copiam (baixam) as páginas de interesse para serem processadas localmente com maior relevância e eficiência (MASANES, 2006). Todo esse processo ocorre de forma autônoma e quando um rastreador entra em processamento, ele acessa páginas de servidores com critérios de busca, baixando e indexando dados de inúmeras páginas, como informações de:

¹ Autor/Apresentador – Sistemas de Informação - UFN; e.pavani@ufn.edu.br

² Professor Colaborador – Direito - UFN; cbisaia@ufn.edu.br

³ Professora Colaboradora – Direito - UFN; rosanelealdasilva@ufn.edu.br

⁴ Consultor Externo – ER Sistemas; robertson@ersistema.info

⁵ Professor Orientador – Sistemas de Informação - UFN; alexz@ufn.edu.br

1. acórdãos jurídicos ou sentenças proferidas em Tribunais;
2. leis presentes na Constituição;
3. resumos e palavras-chave de artigos em portais de revistas eletrônicas.

Dessa forma, a pesquisa ou a busca passa ter um caráter de maior relevância, uma vez que pode elencar páginas mais precisas de conteúdo da pesquisa. Um rastreador da Web trabalha com uma lista de URL (endereços Web em navegadores) a serem visitados (neste contexto, sites do Planalto Federal, de Tribunais Federais e Revistas Eletrônicas). Esses URL são chamados de sementes e à medida que o rastreador visita esses URL, ele identifica todos os *hiperlinks* nas páginas da Web recuperadas e os adiciona à lista de URL a serem visitadas. Essa lista de URL comporta-se de acordo com um conjunto de políticas e se o rastreador estiver realizando o arquivamento, ele copia e salva as informações à medida que avança. Os arquivos são geralmente armazenados de forma que possam ser visualizados, lidos e navegados localmente (MASANES, 2006).

A área do Direito da instituição vem trabalhando com questões que envolvem Tecnologia da Informação, em especial a Inteligência Artificial. A saber:

- Técnicas de Inteligência Artificial Aplicadas ao Direito: Representação de Conhecimento e Raciocínio Automatizado;
- Direito à Saúde e à Educação de Crianças e Adolescentes em Tempos De Pandemia: A Atuação dos Entes Públicos Brasileiros na Efetivação de Direitos Fundamentais.

Esses projetos possuem fases que necessitam de buscas relevantes em sites específicos, como Tribunais de Justiça e portais de revistas científicas. Até então, toda a busca é realizada manualmente, site a site, palavra-chave por palavra-chave. Assim, automatizar o processo é importante, além de qualificar o resultado da busca, justificando, portanto, este trabalho. Além disso, o sistema proposto vai obedecer à estrutura e à dinâmica presentes em sistemas já construídos (SIGGEP-COMIC, SIGGEP-SADEPI, TFGOnline, Residencia On-line), pelo Laboratório de Práticas da Computação UFN e registrados formalmente junto ao Instituto Nacional de

Propriedade Intelectual (INPI). Destaca-se, também, que a empresa parceira, ER Sistemas, como nos outros sistemas, está auxiliando no desenvolvimento, principalmente nas questões de disponibilizar na Internet com total segurança os dados da base, além de fornecer boas práticas de desenvolvimento de software.

Dessa forma, o objetivo da pesquisa é projetar, desenvolver e implantar um sistema Web para realizar buscas relevantes em sites e/ou portais jurídicos. Para que o objetivo geral seja alcançado, identificaram-se alguns objetivos específicos:

- entender, testar e aplicar *Web crawlers*;
- modelar aspectos estruturais e funcionais;
- estudar e avaliar bibliotecas do universo Python que promovem buscas relevantes;
- mapear e compilar trabalhos relacionados que usaram técnicas de rastreamento em portais;
- estudar e testar mecanismos de CAPTCHA; estudar e testar API de proteção de portais que bloqueiam sistemas *bots*.

Assim, este artigo apresenta os resultados necessários para o objetivo específico “modelar aspectos estruturais e funcionais” do sistema proposto.

2. METODOLOGIA

Este trabalho é baseado em pesquisa exploratória com modelagem e desenvolvimento de produto computacional, tendo como referência os portais Planalto Federal (constituição federal), Tribunais (acórdãos jurídicos ou sentenças) e Revistas Eletrônicas do Direito (artigos). Nas fases de projeto, modelagem e desenvolvimento do *Web crawler*, são utilizados a metodologia SCRUM (SUTHERLAND, 2016) com a técnica *Kanban* para gestão de atividades, prazos e equipe. As ferramentas utilizadas neste artigo são:

- Trello - técnica kanban;
- ASTAH - diagramação UML do sistema;
- Github - versionamento de código.

3. RESULTADOS E DISCUSSÕES

Os resultados iniciais do trabalho, referem-se a modelagem dos aspectos funcionais e estruturais do sistema proposto. Os aspectos funcionais podem ser visualizados em Diagramas de Atividades (Figura 1) e de Casos de Uso (Figura 2). Quando questões funcionais são modeladas ou mapeadas, destacam-se os atores que fazem relação com o sistema, as funcionalidades ou serviços que o sistema deve atender e o fluxo de funcionamento ou de interação do sistema com esses atores.

Na Figura 1, é possível acompanhar a dinâmica de funcionamento do sistema *Web crawler*, em que um usuário do sistema, depois do *login*, dentro da gestão básica, pode realizar diferentes consultas como áreas do Direito e Fontes de Consulta. Em seguida, esse usuário, no ambiente de pesquisa, pode filtrar áreas e executar o rastreador ou *crawler*. Uma vez disparado o rastreador, o sistema busca nas fontes cadastradas se o site ou o repositório possui algum tipo de segurança, como o CAPTCHA. Caso o site não possua, os resultados são devolvidos ao sistema, em forma de arquivos baixados.

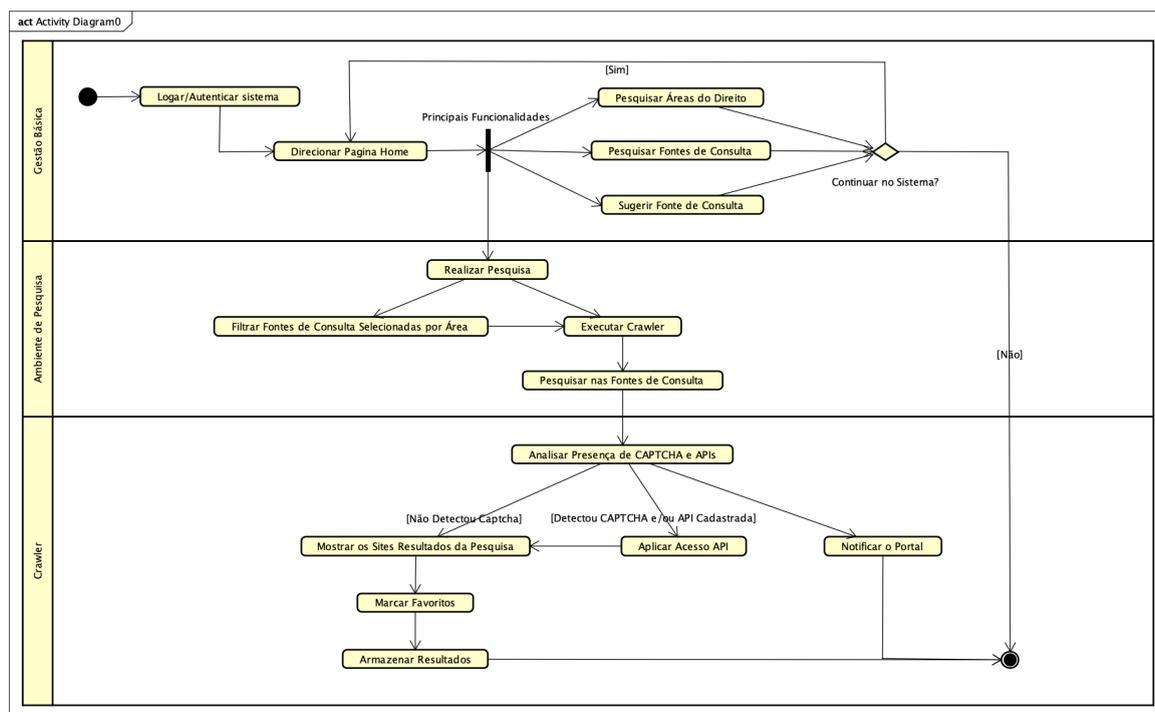


Figura 1: Diagrama de Atividades.

e/ou armazenados. Por fim, o pacote Crawler, que é o principal foco desta pesquisa, uma vez que é nele que serão usadas e implementadas técnicas de rastreamento digital.

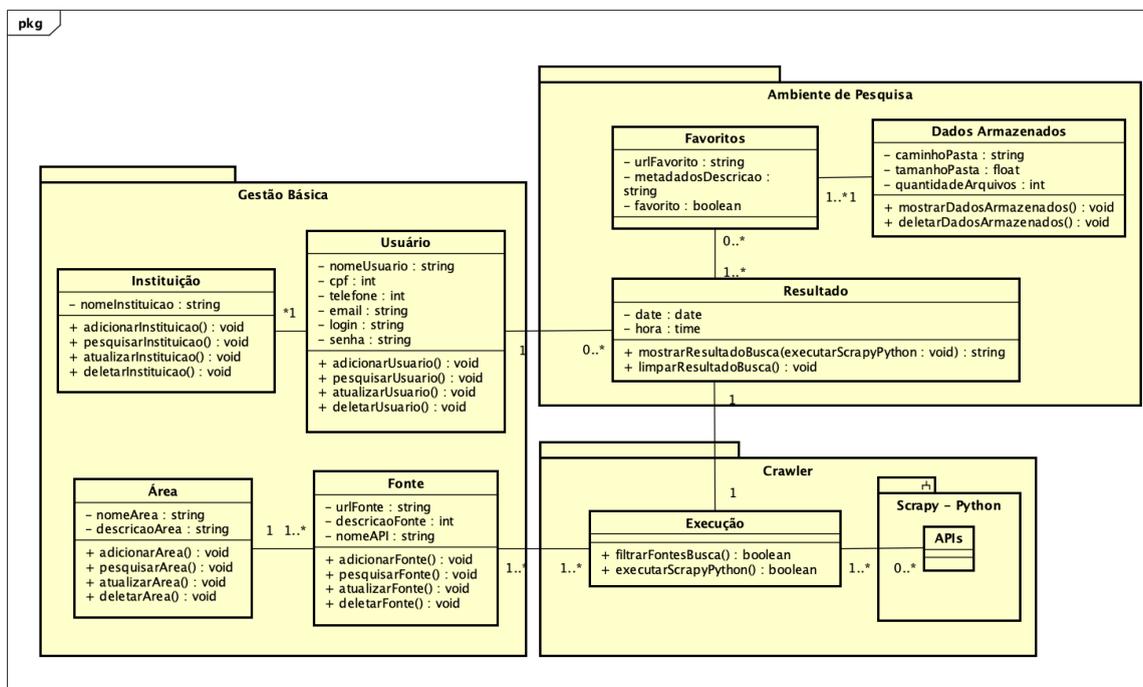


Figura 3: Diagrama de Pacotes.

Como já mencionado, as tabelas detalham cada Caso de Uso mapeado. Esse detalhamento é importante na fase de programação, uma vez que o programador tem acesso a peculiaridades da funcionalidade, como atores envolvidos, condições, regras de funcionamento (fluxo principal) e exceções (fluxo alternativo).

Tabela 1: Realizar Pesquisa.

CASO DE USO	Realizar Pesquisa (UC009)
DESCRIÇÃO	Esses atores poderão realizar a pesquisa (crawler)
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007);
PÓS-CONDIÇÕES	Inclusão: Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Web crawler (UC011);
FLUXO PRINCIPAL	Usuário do sistema selecionará os filtros da consulta (UC010); Sistema irá executar algoritmo de pesquisa (UC011) nos sites cadastrados (fontes de consulta - UC006) que atendam as requisições seleccionadas previamente pelo usuário;
FLUXO ALTERNATIVO	-

Tabela 2: Filtrar Fontes de Consulta Seleccionadas por Área.

CASO DE USO	Filtrar Fontes de Consulta Seleccionadas por Área (UC010)
DESCRIÇÃO	Esses atores poderão selecionar itens (filtros) de pesquisa para direcionar busca
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Área do direito não cadastrada;

Tabela 3: Executar Web crawler.

CASO DE USO	Executar Web crawler (UC011)
DESCRIÇÃO	Esses atores ativarão o algoritmo buscador com base nos critérios selecionados no filtro
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010);
PÓS-CONDIÇÕES	Inclusão: Pesquisar nas Fontes de Consulta (UC016); Baixar os sites resultado da pesquisa (UC013); Marcar Favoritos(UC012);
FLUXO PRINCIPAL	Analisar presença de CAPTCHA e aplicar APIs que irão permitir busca
FLUXO ALTERNATIVO	-

Tabela 4: Marcar Favoritos.

CASO DE USO	Marcar Favoritos (UC012)
DESCRIÇÃO	Esses atores selecionaram sites da busca que julgarem relevantes para armazenar no banco de dados; Os dados favoritados ficarão prontamente disponíveis no sistema para os usuários pesquisarem futuramente
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Webcrawler(UC011);
PÓS-CONDIÇÕES	Exclusão: Mostrar os sites resultados da pesquisa (UC013); Inclusão: descartar não favoritados (UC014); Armazenar resultados favoritados(UC015);
FLUXO PRINCIPAL	Serão mostrados os sites resultado da pesquisa (UC013); Na sequência o usuário irá selecionar os favoritos e consequentemente, não selecionar os sites desnecessários (UC014); Por fim, os sites favoritados serão armazenados no sistema (UC015);
FLUXO ALTERNATIVO	Usuário não identificou sites relevantes na busca para favoritar

Tabela 5: Mostrar os Sites Resultados da Pesquisa.

CASO DE USO	Mostrar os Sites Resultado da Pesquisa (UC013)
DESCRIÇÃO	Os sites mostrados serão aqueles que o crawler identificou na web à partir dos filtros selecionados
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Selecionadas por Área (UC010); Executar Webcrawler(UC011); Marcar Favoritos (UC012);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Não encontrar resultado

Tabela 6: Descartar os Não Favoritos.

CASO DE USO	Descartar os Não Favoritos (UC014)
DESCRIÇÃO	Os sites da pesquisa que os atores julgaram ser irrelevantes para armazenar
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Selecionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa(UC013); Marcar Favoritos (UC012);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	Deixar item de seleção em branco
FLUXO ALTERNATIVO	Todos os itens da pesquisa serem selecionados

Tabela 7: Armazenar Resultados Favoritos.

CASO DE USO	Armazenar Resultados Favoritos (UC015)
DESCRIÇÃO	Os sites armazenados serão aqueles que o crawler identificou na web e o usuário do sistema obrigatoriamente favoritou
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Selecionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa (UC013); Marcar Favoritos (UC012);
PÓS-CONDIÇÕES	Inclusão: Limpar e Classificar Resultados (UC017);
FLUXO PRINCIPAL	Pegar sites mostrados que foram favoritados; Armazenar metadados da busca, como URL da fonte, no banco de dados do sistema
FLUXO ALTERNATIVO	Não ter dados para armazenar

Tabela 8: Pesquisar nas Fontes de Consulta.

CASO DE USO	Pesquisar nas Fontes de Consulta (UC016)
DESCRIÇÃO	Crawler irá pesquisar dados nas fontes cadastradas, previamente, pelo Especialista após execução do crawler
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Webcrawler(UC011);
PÓS-CONDIÇÕES	Inclusão: Avaliar Presença de CAPTCHA e APIs(UC025); Planalto - Legislação (UC018); Tribunais com Jurisprudência (UC019); Revistas Eletrônicas do Direito (UC020);
FLUXO PRINCIPAL	Adentrar nas fontes cadastradas e realizar busca com base nos filtros selecionados
FLUXO ALTERNATIVO	API não funciona; Falta de API; Notificar sistema sobre problema na busca envolvendo mecanismos de segurança

Tabela 9: Limpar e Classificar Resultados.

CASO DE USO	Limpar e Classificar Resultados (UC017)
DESCRIÇÃO	Aplicar IA nos dados armazenados para ajudar o usuário do sistema (Advogado) qual o melhor caminho à seguir
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa (UC013); Marcar Favoritos (UC012); Armazenar Resultados Favoritos(UC015);
PÓS-CONDIÇÕES	Inclusão: Analisar Acórdãos(UC022); Descobrir Embasamento Jurídico(UC021); Quantificar(UC023);
FLUXO PRINCIPAL	Analisar Acórdãos; Descobrir Embasamento Jurídico das Sentenças; Aplicar mineração de dados estimando métricas, valores e/ou probabilidades dos julgamentos
FLUXO ALTERNATIVO	-

Tabela 10: Planalto - Legislação.

CASO DE USO	Planalto - Legislação (UC018)
DESCRIÇÃO	Uma das Fontes de pesquisa
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Webcrawler(UC011); Pesquisar nas Fontes de Consulta(UC016);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	-

Tabela 11: Tribunais com Jurisprudência.

CASO DE USO	Tribunais Com Jurisprudência (UC019)
DESCRIÇÃO	Fontes de pesquisa
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Webcrawler(UC011); Pesquisar nas Fontes de Consulta (UC016);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	Aplicar API; Buscar dados com base nos filtros seleccionados
FLUXO ALTERNATIVO	Problema de validação da API

Tabela 12: Revistas Eletrônicas do Direito.

CASO DE USO	Revistas Eletrônicas do Direito (UC020)
DESCRIÇÃO	Fontes de pesquisa
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Webcrawler(UC011); Pesquisar nas Fontes de Consulta (UC016);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	Aplicar API; Buscar dados com base nos filtros seleccionados;
FLUXO ALTERNATIVO	Problema de validação da API

Tabela 13: Descobrir o Embasamento Jurídico.

CASO DE USO	Descobrir Embasamento Jurídico (UC021)
DESCRIÇÃO	Identificar validade(leis) do que foi armazenado
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa (UC013); Marcar Favoritos (UC012); Armazenar Resultados Favoritos (UC015); Limpar e Classificar Resultados (UC017);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Não ter dados suficientes para análise

Tabela 14: Analisar Acórdãos.

CASO DE USO	Analisar Acórdãos (UC022)
DESCRIÇÃO	Analisar sentenças proferidas um tribunal
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa (UC013); Marcar Favoritos (UC012); Armazenar Resultados Favoritos (UC015); Limpar e Classificar Resultados(UC017);
PÓS-CONDIÇÕES	-
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Não ter dados suficientes para análise

Tabela 15: Quantificar.

CASO DE USO	Quantificar (UC023)
DESCRIÇÃO	Estabelecer métricas nos resultados armazenados
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Seleccionadas por Área (UC010); Executar Web crawler (UC011); Mostrar os Sites Resultado da Pesquisa (UC013); Marcar Favoritos (UC012); Armazenar Resultados Favoritos (UC015); Limpar e Classificar Resultados (UC017);
PÓS-CONDIÇÕES	Mostrar para o usuário do sistema probabilidades de ganho e/ou leis/jurisprudência que validam hipóteses
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Não ter dados suficientes para análise

Tabela 16: Sugerir Fontes de Consulta.

CASO DE USO	Sugerir Fonte de Consulta (UC024)
DESCRIÇÃO	Ator poderá suggestionar fonte de seu agrado para administrador e especialista analisarem viabilidade de inserir como fonte de consulta
ATORES	Usuário
PRÉ-CONDIÇÕES	Estar Logado;
PÓS-CONDIÇÕES	Extensão: Gestão de Fontes de Consulta do Direito do Direito (UC003);
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Fonte não ser aceita pela presença de CAPTCHA e/ou não possibilidade de acesso para busca

Tabela 17: Avaliar Presença de CAPTCHA e APIs.

CASO DE USO	Avaliar Presença de CAPTCHA e APIs (UC025)
DESCRIÇÃO	Identificar mecanismos de segurança nas fontes selecionadas;
ATORES	Usuário; Especialista
PRÉ-CONDIÇÕES	Estar Logado; Fonte(s) de consulta(s) cadastradas no sistema (UC003); Área(s) do direito cadastrada(s) no sistema (UC007); Realizar Pesquisa (UC009); Filtrar Fontes de Consulta Selecionadas por Área (UC010); Executar Webcrawler(UC011); Pesquisar nas Fontes de Consulta(UC016);
PÓS-CONDIÇÕES	Pesquisar e trazer os dados para o usuário do sistema
FLUXO PRINCIPAL	-
FLUXO ALTERNATIVO	Não conseguir realizar a pesquisa em determinada fonte

4. CONCLUSÃO

O texto, inicialmente, apresenta a ideia do trabalho, justificativa, objetivos e alguns conceitos e dinâmica sobre *Web crawlers*, ou o processo de rastreamento. Em seguida, dá ênfase no processo de modelagem, apresentando diagramas e tabelas que especificam ou modelam aspectos funcionais e estruturais. Como as boas práticas da metodologia SCRUM estão sendo usadas, mapear e detalhar as funcionalidades é o foco deste processo.

Finalmente, com a modelagem finalizada, é possível iniciar o processo de implementação das funcionalidades mapeadas.

AGRADECIMENTOS

A Universidade Franciscana pela bolsa PROBIT-UFN 2022-2023 e ao Laboratório de Práticas da Computação UFN.

REFERÊNCIAS

- MASANES, J. Web archiving: issues and methods. In: Web archiving, p.1-53. Springer, 2006.
- SUTHERLAND, J. Scrum: a arte de fazer o dobro do trabalho na metade do tempo. Leya, 2016.