



TESE DE DOUTORADO

**Serviço de Segurança Baseado em IoT  
para a Cadeia de Custódia Digital**

**Daniel Alves da Silva**

**Brasília, 14 de julho de 2021**

**UNIVERSIDADE DE BRASÍLIA**

FACULDADE DE TECNOLOGIA

**UNIVERSIDADE DE BRASÍLIA  
FACULDADE DE TECNOLOGIA  
DEPARTAMENTO DE ENGENHARIA ELÉTRICA**

**SERVIÇO DE SEGURANÇA BASEADO EM IOT  
PARA A CADEIA DE CUSTÓDIA DIGITAL**

**DANIEL ALVES DA SILVA**

**Orientador: PROF. DR. RAFAEL TIMÓTEO DE SOUSA JR., ENE/UNB**

**TESE DE DOUTORADO EM ENGENHARIA ELÉTRICA**

**PUBLICAÇÃO PPGEE.TD - 179/21  
BRASÍLIA-DF, DE 14 DE JULHO DE 2021.**

**UNIVERSIDADE DE BRASÍLIA  
FACULDADE DE TECNOLOGIA  
DEPARTAMENTO DE ENGENHARIA ELÉTRICA**

**SERVIÇO DE SEGURANÇA BASEADO EM IOT  
PARA A CADEIA DE CUSTÓDIA DIGITAL**

**DANIEL ALVES DA SILVA**

TESE DE DOUTORADO SUBMETIDA AO DEPARTAMENTO DE ENGENHARIA ELÉTRICA DA FACULDADE DE TECNOLOGIA DA UNIVERSIDADE DE BRASÍLIA COMO PARTE DOS REQUISITOS NECESSÁRIOS PARA A OBTENÇÃO DO GRAU DE DOUTOR EM ENGENHARIA ELÉTRICA.

APROVADA POR:

Prof. Dr. Rafael Timóteo de Sousa Jr., ENE/UnB  
Orientador

Prof. Dr. William Ferreira Giozza, ENE/UnB  
Examinador interno

Prof. Dr. Mário Antônio Ribeiro Dantas, DCC/UFJF  
Examinador externo

Prof. Dr. Clarimar José Coelho, PUC Goiás  
Examinador externo

**BRASÍLIA, 14 DE JULHO DE 2021.**

## FICHA CATALOGRÁFICA

SILVA, DANIEL ALVES

**Serviço de Segurança Baseado em IoT para a Cadeia de Custódia Digital** [Distrito Federal] **2021**.  
xvi, 110 p., 210 x 297 mm (ENE/FT/UnB, Doutor, Engenharia Elétrica, 2021).

Tese de Doutorado - Universidade de Brasília, Faculdade de Tecnologia.

Departamento de Engenharia Elétrica

- |  |  |
|--|--|
| 1. Digitalização de documentos e processos | 2. Cadeia de custódia documental       |
| 3. Confiança computacional                 | 4. Segurança cibernética e privacidade |
| 5. Internet das coisas industrial (IIoT)   | 6. Cidades inteligentes e sustentáveis |
| I. ENE/FT/UnB                              |  |

## REFERÊNCIA BIBLIOGRÁFICA

SILVA, D. A. (2021). **Serviço de Segurança Baseado em IoT para a Cadeia de Custódia Digital**. Tese de Doutorado em Engenharia Elétrica, Publicação PPGENE.DM-179/21, Departamento de Engenharia Elétrica, Universidade de Brasília, Brasília, DF, 110 p.

## CESSÃO DE DIREITOS

AUTOR: Daniel Alves da Silva

TÍTULO: **Serviço de Segurança Baseado em IoT para a Cadeia de Custódia Digital**.

GRAU: Doutor em Engenharia Elétrica ANO: 2021

É concedida à Universidade de Brasília permissão para reproduzir cópias desta Tese de Doutorado e para emprestar ou vender tais cópias somente para propósitos acadêmicos e científicos. Os autores reservam outros direitos de publicação e nenhuma parte desta Tese de Doutorado pode ser reproduzida sem autorização por escrito dos autores.

---

Daniel Alves da Silva

Departamento de Engenharia Elétrica (ENE) - FT

Universidade de Brasília (UnB)

Campus Darcy Ribeiro

CEP 70919-970 - Brasília - DF - Brasil



*Dedico este trabalho a minha esposa Stela  
e aos meus filhos Davi e Nicolas.*

# Agradecimentos

Obrigado a meu orientador e amigo, Prof. Dr. Rafael Timóteo de Sousa Júnior, por nunca ter desistido de mim, do nosso país, da nossa cidade, da nossa educação e da nossa Universidade. Obrigado pela empatia e pela porta sempre aberta, muitas vezes parando suas atividades para atender seus alunos e colegas de trabalho em sua concorrida sala. Obrigado por ser o primeiro a chegar e o último a sair, por exigir sempre a excelência acadêmica, ética e o respeito com os recursos públicos. Muito obrigado, pela lição de humanidade que nos deu no enfrentamento à pandemia em que trabalhamos sem nenhum investimento e entregamos soluções funcionais de suporte ao diagnóstico médico e apoio a comunidade. Saiba que repassarei tudo o que foi aprendido contigo a meus filhos.

Agradeço aos professores do Departamento de Engenharia Elétrica da UnB, Robson Oliveira Albuquerque, Fábio Lúcio Lopes Mendonça, Georges Daniel Amvame Nze, Ricardo Staciardini Puttinin, William Ferreira Giozza, Flávio Elias Gomes de Deus e Ugo Silva Dias, pelos ensinamentos e cortesia. Não poderia deixar de citar a coordenação e secretaria do programa que foram essenciais nesta jornada, principalmente nos tempos de pandemia em que todos estavam muito ansiosos e preocupados.

Meu muito obrigado aos professores João Paulo Javidi da Costa (UnB / EFS / THI), Gordon Elger (THI), Alessandro Zimmer (UFPR / THI) e Christian Facchi (THI) pela oportunidade da visita técnica a Technische Hochschule Ingolstadt (THI) e seus centros de pesquisa; ao Professor Luis Javier García Villalba (UCM) pela grande colaboração na publicação do journal deste trabalho; aos professores Mário Antônio Ribeiro Dantas (DCC/UFJ) e Claudio Palasciano (PoliMi) pela oportunidade de fazer a primeira apresentação relativa a este trabalho no Brasil-Italy webinar do projeto FASTEN; e ao Professor Anderson Clayton A. Nascimento (UW Tacoma) pela competência de ministrar conteúdos extremamente complexos de forma clara e acessível em suas aulas de criptografia.

Agradecimento sinceros a Carlos Eduardo L. Veiga, Maria Fernanda Bittencourt, Lucimar Rizzo L. dos Santos, Guilherme Fay Vergara, Matheus S. Fonseca, Dario Santos, Bruno J. G. Praciano, Gabriel Albanese D. de Araújo, Ariovaldo D. Furtado, Francisco L. de Caldas Filho, José Alberto Torres, Ludmila B. Silva, Francisco Vitor Lopes, Vinícius Coutinho, Paulo L. Machado, Alessandro S. Mendes, Rodrigo F. Vergara, Tiago Ianuk, Rodrigo M. dos Santos e a todos os membros do Laboratório LATITUDE/UnB e IEEE VTS Centro-Norte Brasil Chapter, pelo incentivo, pelas valiosas sugestões e discussões construtivas sobre este trabalho.

Agradeço o apoio técnico e computacional do Laboratório de Tecnologias da Tomada de Decisão - LATITUDE, da Universidade de Brasília, que conta com apoio do CNPq - Conselho Nacional de Pesquisa (Outorgas 312180/2019-5 PQ-2, BRICS2017-591 LargEWiN e 465741/2014-2 INCT em Cibersegurança), da CAPES - Coordenação de Aperfeiçoamento do Pessoal de Ní-

vel Superior (Outorgas 23038.007604/2014-69 FORTE e 88887.144009/2017-00 PROBRAL), da FAP-DF - Fundação de Amparo à Pesquisa do Distrito Federal (Outorgas 0193.001366/2016 UIoT e 0193.001365/2016 SSDDC), do Ministério da Economia (Outorgas 005/2016 DIPLA e 083/2016 ENAP), da Secretaria de Segurança Institucional da Presidência da República do Brasil (Outorga ABIN 002/2017), do Conselho Administrativo de Defesa Econômica (Outorga CADE 08700.000047/2019-14), da Advocacia Geral da União (Outorga AGU 697.935/2019), do Ministério das Cidades, (Outorga MC 01/2019), do Ministério da Justiça e Segurança Pública, (Outorga MJSP 01/2019) e dos Decanatos de Pesquisa e Inovação e de Pós-Graduação da Universidade de Brasília (DPI/DPG/UnB).

# SERVIÇO DE SEGURANÇA BASEADO EM IoT PARA A CADEIA DE CUSTÓDIA DIGITAL

**Autor: Daniel Alves da Silva**

**Orientador: Rafael Timóteo de Sousa Júnior**

**Programa de Pós-graduação em Engenharia Elétrica - PPGEE**

**Brasília, 14 de julho de 2021**

A digitalização de documentos e processos é considerada um fator essencial para cidades sustentáveis, pois facilita a interação socioeconômica e contribui para ambientes urbanos e sociedades favoráveis ao clima. Como os documentos digitais são vulneráveis a modificações ilegítimas e falsificações, há o risco de esses documentos serem rejeitados como não confiáveis. Portanto, para realmente se beneficiar dos documentos e processos digitalizados em ambientes urbanos e sociais sustentáveis, é imperativo criar meios para gerar confiança em relação a esses documentos e processos. Como a confiança computacional e a segurança da informação contribuem mutuamente para construir uma à outra, este trabalho se dedica a construir a confiança em documentos digitais, garantindo a cadeia de custódia (CoC) dos documentos produzidos e preservados. O CoC compreende um conjunto de recursos e meios para garantir as propriedades de segurança exigidas aos documentos digitalizados e é um conceito abrangente que permite integrar as medidas de segurança em todo o fluxo de informação produtor-consumidor, sendo explicitamente informado ao utilizador para que este possa confiar que as proteções disponíveis foram implementadas para seus documentos digitais. Como o CoC precisa ser operativo em diversas cadeias de produção de documentos, nosso projeto considera o paradigma Industry Internet of Things (IIoT) para conceber um serviço de segurança baseado em IIoT para a cadeia de custódia documental - IoTSec2CoC, com um modelo de serviço para identificar e autorizar todos os seus componentes, fornecendo microsserviços para a integração das medidas de segurança pró-ativas e reativas totalmente distribuídas, buscando garantir a confidencialidade, disponibilidade, e integridade das unidades de documentos em movimento. Esta tese constitui uma cadeia produtiva que apresenta a plasticidade adequada para cidades inteligentes e sustentáveis. A validação da proposta foi realizada com um protótipo implementado cujos testes de resiliência mostram que o novo modelo proposto aumenta a segurança em todas as etapas do CoC de processamento digital de documentos.

**Palavras-chave:** Digitalização de documentos e processos. Cadeia de custódia documental. Confiança computacional. Segurança cibernética e privacidade. Internet das Coisas Industrial (IIoT). Cidades inteligentes e sustentáveis.

# **IOT-BASED SECURITY SERVICE FOR THE DIGITAL CHAIN OF CUSTODY**

**Author: Daniel Alves da Silva**

**Supervisor: Rafael Timóteo de Sousa Júnior**

**Programa de Pós-graduação em Engenharia Elétrica - PPGEE**

**Brasília, Jul 14, 2021**

Document and process digitalization is considered an essential factor for sustainable cities as it facilitates socio-economic interaction and contributes to climate-friendly urban environments and societies. Since digital documents are vulnerable to illegitimate modifications and forgery, there is a risk of these documents being rejected as untrustful. Therefore, to really benefit from digitalized documents and processes in sustainable urban environments and society, it is imperative to put in place means to generate trust towards these documents and processes. As computational trust and information security mutually contribute to building each other, this thesis is devoted to building trust in digital documents by ensuring the digital chain of custody (CoC) of the produced and preserved documents. The CoC comprises a set of resources and means to ensure the properties of security required for digitalized documents and is a comprehensive concept allowing to integrate the security measures in the whole producer-to-consumer information flow, being explicitly informed to the user so that she/he can trust that available protections were put in place for their digital documents. As the CoC needs to be operative in diverse document production chains, our design considers the Industry Internet of Things (IIoT) paradigm to conceive an IIoT-Based Security Service for the documentary chain of custody – IoTSec2CoC, with a service model to identify and authorize all of its components by providing microservices for the integration of the fully distributed proactive and reactive security measures, to ensure the confidentiality, availability, and integrity of the flowing document units. This thesis constitutes a production chain that presents the plasticity adequate for smart sustainable cities. The proposal validation was carried out with an implemented prototype whose resilience tests show that the proposed new model increases security at all stages of the digital document processing CoC.

**Keywords:** Document and process digitalization. Documentary chain of custody. Computational trust. Cybernetic security and privacy. Industrial Internet of Things (IIoT). Smart sustainable cities.

# SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO</b>	<b>1</b>
1.1	OBJETIVO	4
1.2	OBJETIVOS ESPECÍFICOS	5
1.3	PUBLICAÇÕES VINCULADAS A ESTE TRABALHO	6
1.4	METODOLOGIA	7
1.5	ORGANIZAÇÃO DO TRABALHO	8
<b>2</b>	<b>ESTADO DA ARTE E TRABALHOS CORRELATOS</b>	<b>9</b>
2.1	DOCUMENTO ARQUIVÍSTICO DIGITAL	9
2.2	PRESERVAÇÃO E MANUTENÇÃO DA CADEIA DE CUSTÓDIA DOCUMENTAL	11
2.2.1	MODELO OAIS	13
2.2.2	REPOSITÓRIO ARQUIVÍSTICO DIGITAL	14
2.3	SEGURANÇA COMPUTACIONAL	15
2.3.1	ATAQUES	16
2.3.2	CRIPTOGRAFIA SIMÉTRICA E ASSIMÉTRICA	18
2.3.3	ASSINATURA DIGITAL	19
2.3.4	GERENCIAMENTO DE EVENTOS E INFORMAÇÕES DE SEGURANÇA	20
2.4	COMPUTAÇÃO DISTRIBUÍDA	21
2.4.1	COMPUTAÇÃO EM NUVEM	21
2.4.2	<i>Edge Computing</i>	22
2.4.3	MICROSSERVIÇOS	23
2.4.4	INTERNET DAS COISAS E INDÚSTRIA 4.0	24
2.4.5	MIDDLEWARE IOT	25
2.5	INTELIGÊNCIA ARTIFICIAL E APRENDIZADO DE MÁQUINA	30
2.5.1	REDES NEURAIS	30
2.6	TRABALHOS RELACIONADOS	32
2.6.1	CONFIANÇA E CADEIA DE CUSTÓDIA	32
2.6.2	INDÚSTRIA 4.0	33
2.6.3	SEGURANÇA EM IOT	34
<b>3</b>	<b>ESTUDO DO PROBLEMA E HIPÓTESES DA TESE</b>	<b>35</b>
3.1	CADEIA DE CUSTÓDIA DOCUMENTAL CONVENCIONAL	36
3.2	HIPÓTESES	37
<b>4</b>	<b>PROPOSTA DA ARQUITETURA IoTSec2DCoC PARA GARANTIR A CADEIA DE CUSTÓDIA DOCUMENTAL</b>	<b>39</b>
4.1	VISÃO GERAL DA ARQUITETURA PROPOSTA	39

4.2	ESTRUTURA DA ARQUITETURA DE SOFTWARE IOTSec2DCoC .....	41
4.3	DETALHAMENTO DOS PROCESSOS E FLUXOS DO MODELO .....	42
<b>5</b>	<b>PROVA DE CONCEITO .....</b>	<b>50</b>
5.1	DESCRIÇÃO DO AMBIENTE DE TESTES .....	50
5.1.1	DESCRIÇÃO DO HARDWARE UTILIZADO .....	51
5.1.2	DESCRIÇÃO DO AMBIENTE DE SOFTWARE .....	51
5.1.3	PROTÓTIPO DO DISPOSITIVO INTELIGENTE COM COMPUTADOR DE PLACA ÚNICA .....	52
5.1.4	AMOSTRA DE DADOS .....	53
5.2	DISPOSITIVOS INTELIGENTES DE ADMISSÃO .....	54
5.3	IOT E SUPORTE DE MICROSERVIÇOS PARA COMUNICAÇÕES E INTEGRAÇÃO .....	60
5.4	SERVIÇOS DE SEGURANÇA INTEGRADOS .....	63
5.4.1	SERVIÇO DE AUTORIDADE CERTIFICADORA .....	63
5.4.2	SERVIÇOS DE ASSINATURA DIGITAL E VALIDAÇÃO DE DOCUMENTOS (CHECKOUT E CHECK-IN) .....	65
5.4.3	SERVIÇO DE MONITORAMENTO .....	66
5.4.4	SERVIÇO DE AUTENTICAÇÃO (JWT) .....	66
5.5	MÓDULOS DE APLICAÇÃO .....	67
5.5.1	EMPACOTADOR .....	67
5.5.2	REPOSITÓRIO ARQUIVÍSTICO DIGITAL CONFIÁVEL - RDC-ARQ .....	69
5.5.3	PLATAFORMA DE ACESSO .....	72
5.5.4	COMPONENTE DE CLASSIFICAÇÃO E EXTRAÇÃO DE METADADOS .....	73
<b>6</b>	<b>VALIDAÇÃO DA RESILIÊNCIA DO MODELO .....</b>	<b>79</b>
6.1	RESILIÊNCIA A ATAQUES DE FALSIFICAÇÃO .....	79
6.1.1	TENTATIVA DE INTERCEPTAR O TRÁFEGO ENTRE OS DISPOSITIVOS INTELIGENTES E O MIDDLEWARE EDGE .....	79
6.1.2	TENTATIVA DE TROCAR DADOS COM O MIDDLEWARE EDGE SEM AUTENTICAÇÃO VÁLIDA .....	80
6.1.3	TENTATIVA DE ENVIAR PACOTES DE DOCUMENTOS COM ASSINATURAS INVÁLIDAS PARA O MIDDLEWARE EDGE .....	83
6.1.4	CONTRAMEDIDAS CONTRA VAZAMENTO DE CREDENCIAIS E CHAVES DE CRIPTOGRAFIA .....	86
6.2	RESILIÊNCIA DO MODELO A ATAQUES CONTRA SENHAS .....	86
6.3	CAPACIDADE EFETIVA DE MONITORAMENTO E RASTREABILIDADE .....	88
<b>7</b>	<b>DISCUSSÃO DOS RESULTADOS .....</b>	<b>90</b>
7.1	ASPECTOS DO MODELO E SEU RELACIONAMENTO COM OS COMPONENTES DA SOLUÇÃO .....	90

7.1.1	COMPONENTES E O REFORÇO ÀS PROPRIEDADES DE SEGURANÇA DA INFORMAÇÃO .....	90
7.1.2	COMPONENTES E O COMBATE ÀS AMEAÇAS DE SEGURANÇA DA INFORMAÇÃO .....	91
7.1.3	COMPONENTES E O REFORÇO ÀS PROPRIEDADES DA SOLUÇÃO E DE INTEGRAÇÃO .....	91
7.2	RESULTADOS DA CLASSIFICAÇÃO AUTOMÁTICA DOS DOCUMENTOS E EXTRAÇÃO DE METADADOS .....	92
7.3	RESULTADOS DA RESILIÊNCIA DO MODELO A PERTURBAÇÃO DO FUNCIONAMENTO NORMAL .....	95
7.4	RESULTADOS DA RESILIÊNCIA A ATAQUES DE FALSIFICAÇÃO.....	96
7.5	RESULTADOS DAS CONTRAMEDIDAS CONTRA VAZAMENTO DE CREDENCIAIS E CHAVES DE CRIPTOGRAFIA .....	97
7.6	RESULTADOS DA RESILIÊNCIA DO MODELO A ATAQUES CONTRA SENHAS	97
7.7	RESULTADOS DA RESILIÊNCIA DO MODELO A ATAQUES DOS.....	97
7.8	RESULTADOS DO MONITORAMENTO E RASTREABILIDADE.....	98
<b>8</b>	<b>CONCLUSÃO.....</b>	<b>100</b>
8.1	TRABALHOS FUTUROS .....	102
	<b>REFERÊNCIAS BIBLIOGRÁFICAS.....</b>	<b>104</b>



# LISTA DE FIGURAS

2.1	Cadeia de Custódia dos documentos arquivísticos tradicionais. Adaptado de: [21]..	12
2.2	Diagrama do modelo de referência OAIS. Fonte: Adaptado de: [23].....	13
2.3	Processo de criptografia simétrica. Adaptado de: [33]. .....	19
2.4	Processo de criptografia assimétrica. Adaptado de: [33].....	20
2.5	Modelo NIST de definição de Computação em Nuvem. Fonte: [39]. .....	21
2.6	O paradigma de computação de borda. Fonte: [40]. .....	23
2.7	Comparação entre as arquiteturas Monolítica e Microsserviços. Fonte: [42].....	24
2.8	Estrutura do Middleware IoT. Fonte: Adaptado de: [50].....	26
2.9	Arquitetura do Middleware hierárquico do UIoT. Fonte: Adaptado de: [52]. .....	27
3.1	Arquitetura de uma solução de cadeia de custódia convencional. Fonte: Adaptado de: [23]. .....	36
4.1	Visão geral da estrutura IoTSec2DCoC. ....	40
4.2	Arquitetura de Software IoTSec2DCoC. ....	41
4.3	Arquitetura da solução IoTSec2DCoC e seus fluxos.....	45
4.4	Identificação de Dispositivo Inteligente e outros componentes IoT. ....	46
4.5	Identificação dos componentes da cadeia de custódia convencional. ....	47
4.6	Identificação dos componentes Check-in e Check-out. ....	48
4.7	Identificação dos serviços de segurança.....	49
5.1	Foto do protótipo de dispositivo montado. ....	53
5.2	Fluxo de login do Dispositivo Inteligente na rede IOT. ....	54
5.3	Fluxo de processamento do componente OCR. ....	55
5.4	Log de processamento do Dispositivo Inteligente implantado na PoC. ....	57
5.5	Registro do Dispositivo Inteligente na rede IOT. ....	59
5.6	Transferência de pacote de documentos do Dispositivo Inteligente para o Middleware IoT Edge. ....	60
5.7	Dashboard para visualização de dados do UIoT implantado na PoC. ....	61
5.8	Arquitetura com um servidor Edge. Adaptado de: [84]. ....	62
5.9	Tela de revisão de pacotes do Empacotador implantado na PoC.....	68
5.10	Estrutura do SIP transferida do Empacotador para o RDC-Arq.....	68
5.11	Tela de transferência de pacotes do Archivematica implantado na PoC.....	69
5.12	Tela de administração de pacotes (AIP e DIP) do Archivematica implantado na PoC.....	70
5.13	Registro da Política de Formatos. ....	71
5.14	Tela de visualização de documentos do Atom implantado na PoC. ....	72
5.15	Arquitetura da rede neural com 2 camadas.....	73

5.16	Gráfico do resultado de treinamento e validação de acurácia. ....	74
5.17	Matriz de confusão. ....	75
5.18	Marcação de metadados para treinamento. ....	76
5.19	Modelo preditivo treinado. ....	77
6.1	Interceptação malsucedida do conteúdo da mensagem. ....	80
6.2	Diagrama de sequência com exemplificação do ataque sem autenticação válida. ....	81
6.3	Falha de validação do JWT no Middleware Edge. ....	82
6.4	Diagrama de sequência com exemplificação do ataque sem assinatura digital de documentos válida. ....	83
6.5	Falha de validação de assinatura digital solicitado pelo Componente de Check-in...	85
6.6	Falha de validação do certificado digital solicitado pelo Componente de Check-in..	86
6.7	Número máximo configurado de autenticações falhas atingido e desabilitação da conta. ....	87
6.8	Tentativa de autenticação com a conta desabilitada. ....	88
6.9	Log de auditoria consultada na interface web Kibana. ....	89
7.1	Distribuição de classes ao longo do dataset de treinamento de detecção de objetos.	93
7.2	Análise de precisão ao longo das 300 épocas de treinamento. ....	93
7.3	Análise da revocação ao longo das 300 épocas de treinamento para o sistema de detecção de objetos. ....	94
7.4	Valor médio de precisão com um threshold de 0.5 ao longo das 300 épocas de treinamento. ....	94
7.5	Valor médio de precisão com um threshold de 0.95 ao longo das 300 épocas de treinamento. ....	95
7.6	Tratamentos e transferências sequenciados de componentes IoTSec2DCoC. ....	98

# LISTA DE TABELAS

2.1	Equações de predição do YOLO [62].....	32
4.1	Componentes da IoTSec2DCoC. ....	42
5.1	Máquinas Virtuais na Amazon EC2.....	51
5.2	Sumarização de tecnologias aplicadas na PoC da IoTSec2DCoC. ....	52
5.3	Descrição da amostra de dados.....	53
7.1	Contribuições de cada componente para os serviços de segurança do IoTSec2DCoC.	90

# LISTA DE TERMOS E SIGLAS

AES	Advanced Encryption Standard
API	Application Programming Interface
CNN	Convolutional Neural Network
CoC	Chain of custody
ECDSA	Elliptic Curve Digital Signature Algorithm
EDGE	Edge computing
GNU	GNU Operating System
GPG	GNU Privacy Guard
HTTP	Hypertext Transfer Protocol
ISAD(G)	General International Standard Archival Description
ISO	International Organization for Standardization
IIoT	Industrial Internet Of Things
IoT	Internet Of Things
JSON	JavaScript Object Notation
M2M	Machine-to-Machine
NLP	Natural Language Processing
PREMIS	Preservation Metadata International Standard
RSA	Public key cryptography algorithm
UIoT	UnB Internet Of Things
SHA	Secure Hash Algorith
TLS	Transport Layer Security
XML	eXtensible Markup Language
URI	Uniform Resource Identifier
XSS	Cross-Site Scripting
XXE	XML External Entities
NIST	National Institute of Standards and Technology
IaaS	Infrastructure-as-a-Service
PaaS	Platform-as-a-Service
SaaS	Software-as-a-Service
UIT	União Internacional de Telecomunicações
SPA	Single Page Applications
DIMS	Data Interface Management System
UIMS	User Interface Management System
DoS	Ataques de negação de serviço
DDoS	Ataques Distribuídos
ICMP	Internet Control Message Protocol
TCP	Transmission Control Protocol

METS	Metadata Encoding and Transmission Standard
DCMI	Dublin Core
FPR	Registro da Política de Formatos
OAIS	Open Archival Information System
OCR	Optical Character Recognition

# 1 INTRODUÇÃO

A digitalização é considerada um fator essencial para o desenvolvimento sustentável das cidades, tanto como facilitador da interação socioeconômica quanto contribuinte para ambientes urbanos e sociedades favoráveis ao clima [1]. A desmaterialização das coleções de documentos impressos a partir do suporte analógico e a substituição dos impressos pelos digitais tem efeito direto na pegada de carbono, ao permitir linhas de produção documental que apresentem melhor desempenho ambiental em relação às emissões diretas e indiretas, sequestro de carbono nas florestas, valor bioenergético, e emissões evitadas, bem como custo de armazenamento, energia e ocupação do espaço urbano.

Porém, os documentos digitais são vulneráveis a modificações ilegítimas e falsificações, fazendo com que as pessoas questionem tais documentos e os processos usados para produzi-los e preservá-los. Com efeito, a presunção da autenticidade de um documento digital, sem a verificação adequada da sua origem e sem medidas de preservação da sua integridade, pode conduzir ao risco de repúdio desse documento. Essas questões existiam em mídia analógica e se expandiram para o mundo digital, o que possivelmente acelerou a produção indiscriminada de documentos e informações sem valor e sua multiplicação em ambientes não integrados e não controlados, dificultando a busca e recuperação de informações, levando até mesmo ao consentimento da produção, armazenamento e preservação de documentos feitos ou modificados de forma maliciosa.

Portanto, para realmente se beneficiar de documentos e processos digitalizados em ambientes urbanos e sociais sustentáveis, é imperativo estabelecer meios para gerar confiança em relação a esses documentos e processos [1]. O trabalho anterior relacionado [2] discute a sinergia e dualidade entre confiança computacional e segurança da informação, mostrando como a confiança e a segurança contribuem mutuamente para a construção uma da outra. Assim, este trabalho é dedicado a construir confiança em documentos digitais, propondo medidas de segurança inteligentes e distribuídas baseadas em IoT para toda a linha de produção de documentos digitais, garantindo a cadeia de custódia (CoC) digital dos documentos produzidos e preservados, fornecendo, ainda, meios para verificar a autenticidade e a integridade de documentos digitais, sejam eles convertidos de impressos ou nascidos digitalmente.

Outra preocupação é que a implantação do processo de digitalização seja obrigada a seguir os padrões úteis aceitos em gestão, tecnologia e segurança da informação como a *International Organization for Standardization* (ISO). Por exemplo, o modelo de referência *Open Archival Information System* (OAIS) ISO 14721 [3], o qual no Brasil foi incorporado ao Modelo de Requisitos para Sistemas Informatizados de Gestão Arquivística de Documentos (e-ARQ Brasil) [4], visa estabelecer os requisitos mínimos e desejáveis para garantir a chamada cadeia de custódia documental (CoC).

A cadeia de custódia é um conceito central para a preservação dos documentos arquivísticos

digitais. Tendo origem jurídica e remetendo à sucessão de fatos cronológicos encadeados que provem algum evento, conforme observado no trabalho de [5], a CoC contribui para manter e documentar a história cronológica da evidência e para rastrear a posse e o manuseio da amostra a partir do preparo do recipiente coletor, o procedimento de coleta, transporte, recebimento, análise e armazenamento. O conceito inclui toda a sequência de posse para a evidência, portanto a CoC deve ser ininterrupta como uma representação da linha contínua de custodiadores dos documentos arquivísticos (desde seus produtores até seus legítimos consumidores). Isso também assegura que os documentos preservados sejam os mesmos desde o início da cadeia e que eles não sofram qualquer processo de mudança, sendo, portanto, autênticos.

A cadeia de custódia (CoC) dos documentos produzidos e preservados é essencial para construir a confiança nos documentos e demais objetos digitais, fornecendo um conjunto de recursos e meios para garantir as propriedades de segurança exigidas aos documentos digitalizados, ou seja, sua integridade, disponibilidade e confidencialidade, de acordo com as necessidades dos usuários. A CoC é um conceito abrangente, pois permite modularizar e integrar as medidas de segurança para documentos digitais convertidos de impressos ou nascidos digitalmente. Essas medidas devem ser distribuídas em todo o fluxo de informação do produtor ao consumidor, sendo explicitamente informadas ao usuário para que este possa confiar que as proteções disponíveis foram implementadas para seus documentos digitais. Como a CoC precisa ser operativa nas cadeias de produção de documentos, nosso projeto considerou o paradigma da Internet das Coisas (IoT) para conceber um Serviço de Segurança Baseado em IoT para a cadeia de custódia digital - IoTSec2DCoC. Esse modelo de serviço visa identificar e autorizar todos os seus componentes dentro de uma Infraestrutura de rede de Industrial IoT (IIoT) para fornecer microsserviços para a integração das medidas de segurança pró-ativas e reativas totalmente distribuídas, a fim de garantir a confidencialidade, disponibilidade e integridade das unidades de informação em fluxo, constituindo, assim, uma cadeia produtiva que apresenta a plasticidade e flexibilidade adequadas para cidades sustentáveis inteligentes

Os documentos arquivísticos em todo o mundo, em meio analógico ou digital, são importantes fontes de informações e meios para que se cumpram as iniciativas de transparência governamental. Nesse sentido, precisam ser armazenados e preservados com critérios que visem à garantia de suas características de confiabilidade, autenticidade e acessibilidade.

Como uma consequência necessária dessa definição, a autenticidade dos documentos arquivísticos digitais deve estar apoiada na evidência de que eles são mantidos utilizando tecnologias, medidas de segurança e procedimentos administrativos que garantam sua identidade e integridade (fatores essenciais da autenticidade) ou, ao menos, que essas medidas minimizem o risco de modificações em um documento, desde sua entrada na CoC e em todas as transferências e acessos subsequentes. Portanto, presume-se que cada elemento de processamento de um documento em uma cadeia de custódia seja capaz de confirmar o processamento e fornecer garantias sobre a integridade do documento processado. Essas confirmações asseguram a ininterruptabilidade da cadeia de custódia, desde a produção inicial do documento até a sua transferência para a instituição arquivística responsável pela sua preservação no longo prazo. Caso essa cadeia de custódia seja

interrompida, o tempo em que os documentos não estiveram sob a proteção do seu produtor ou sucessor pode causar dúvidas sobre a sua autenticidade [4].

Assegurar que a cadeia de produção e custódia documental funcione como planejado traz desafios interessantes relacionados a medidas para mitigar as vulnerabilidades dos documentos arquivísticos digitais, os ciclos de obsolescência tecnológica de mídias e formatos e a dificuldade de se provar a autenticidade dos documentos digitais.

Como há muitas vulnerabilidades na cadeia de custódia, os documentos digitais são suscetíveis a vários tipos de ataques, como interceptação, modificação e fabricação. Portanto, a fim de conter quaisquer ações que comprometam a integridade documental, é necessário adotar medidas de segurança da origem até o destino do documento digital. Essa destinação é comumente um repositório confiável, por exemplo, um repositório arquivístico digital confiável [6], e, para cumprir os requisitos do modelo de referência OAIS [3], as medidas de segurança da cadeia de custódia devem ser projetadas para integrarem-se com as do repositório digital confiável.

Motivado por essa necessidade, o presente trabalho propõe um abrangente conjunto de medidas de segurança computacionais concebido e orquestrado para proteger e conferir confiabilidade à cadeia de custódia de documentos digitais, tanto os natos digitais quanto os digitalizados de seus equivalentes impressos. Nossa proposta está em conformidade com o modelo de referência OAIS [3] e foi concebida considerando um modelo adversarial oriundo da análise de vulnerabilidades e correspondentes possibilidades de ataque à cadeia de custódia documental.

A presente proposta inclui sensores de hardware e de software que são implantados para supervisionar todo o fluxo de documentos em uma linha de produção desde a fonte analógica (scanner) ou digital até o repositório arquivístico. Esses sensores são orquestrados similarmente a uma cadeia industrial de controladores, que segue o paradigma de Internet das Coisas (IoT) para prover serviços de segurança orientados à autenticidade e integridade dos documentos em fluxo, automatizando a cadeia de custódia desse fluxo documental.

Com o nome de Modelo de Segurança com Internet das Coisas para Cadeia de Custódia Digital – IoTSec2DCoC, nossa proposta compreende um conjunto de componentes e protocolos de software que visa identificar e autorizar todos os componentes utilizados dentro de uma infraestrutura de rede IoT hierarquizada sob o controle de um Middleware IoT. Os componentes de software utilizam a arquitetura de microsserviços para fornecer serviços de segurança, utilizando técnicas de criptografia, assinatura digital e verificação sistemática da integridade dos documentos em cada transferência dentro da cadeia de produção documental, constituindo um conjunto abrangente de medidas de confidencialidade, autenticidade e integridade do fluxo de pacotes de informações em relação aos documentos transferidos. O modelo proposto inclui um serviço de gestão de certificados digitais a fim de configurar e validar assinaturas para túneis de comunicação *Transport Layer Security* (TLS) versão 1.3 seguros. A IoTSec2DCoC também inclui serviços de autenticação para usuários e dispositivos, bem como um serviço para monitorar eventos de segurança e detectar tentativas de intrusão contra a cadeia de custódia documental.

A arquitetura IoTSec2DCoC proposta compõe-se de: (1) dispositivos inteligentes de IoT para



captura documental, contendo componentes de captura de documentos, reconhecimento ótico de caracteres (OCR) e serviço de admissão de arquivo nato-digital; (2) módulo de comunicado e integração, contendo Middleware Hierárquico e integração com o gateway da interface de programação de aplicativos (API) e banco de dados; (3) serviços de segurança, contendo certificado e assinatura digital, autenticação e monitoramento de eventos; e (4) módulos de aplicação, contendo componentes de classificação e extração de metadados baseados em aprendizado de máquina, módulo de empacotamento e integração ao repositório arquivístico confiável e plataforma de acesso.

Para validar a proposta, foi escolhido um método experimental para verificar, em uma instalação laboratorial controlada, as funcionalidades propostas usando um protótipo operacional que foi desenvolvido com componentes de hardware e software. Baterias de testes no cenário controlado foram realizadas para demonstrar as funcionalidades e obter registros do funcionamento e da interoperação das instâncias da arquitetura. Essas operações foram analisadas com inputs representativos das ameaças de segurança previstas em nosso modelo adversarial inicial, o qual foi criado para comprometer a integridade dos metadados e documentos digitais, tramitando pelo protótipo implementado até sua custódia definitiva em um repositório confiável. O comportamento observado da arquitetura IoTSec2DCoC proposta apresenta evidências de sua resiliência contra as ameaças e mostra que serviços de segurança ativados por uma estrutura IoT são capazes de garantir a cadeia de custódia para a cadeia de produção documental.

## 1.1 OBJETIVO

O objetivo geral deste trabalho é conceber um conjunto de medidas de segurança computacional para proteger e prover confiança quanto à cadeia de custódia de documentos digitais que utiliza modelo OAIS e Repositório Digital Arquivístico Seguro (RDC-Arq), com suporte de uma arquitetura distribuída de IoT, para se contrapor a um modelo adversarial dessa cadeia de custódia, modelo este também elaborado durante o trabalho. Os serviços de segurança totalmente distribuídos apoiados por uma instância IIoT, sendo explicitamente ativos e visíveis para o usuário dos documentos arquivísticos CoC, constituem um *framework* que contribui para a confiança do usuário em relação aos documentos protegidos e preservados, para que a proposta possa beneficiar cidades sustentáveis, pois facilita as interações socioeconômicas do usuário e leva a ambientes urbanos e sociedades favoráveis ao clima, uma vez que a substituição de documentos impressos por documentos digitais tem um efeito direto na pegada de carbono e na ocupação do espaço urbano.

Para isso, é proposto o Modelo de Segurança com Internet das Coisas para Cadeia de Custódia Documental - IoTSec2DCoC, sendo sua arquitetura composta de: (1) dispositivos físicos e virtuais para a captura e importação documental, contendo os componentes de captura/OCR e serviço de admissão de arquivo nato-digital; (2) IoT e suporte de microsserviços para comunicações e integração, com os componentes Gateway de API, Middleware IoT Edge, Middleware IoT Cloud, banco de dados; (3) serviços de segurança integrados, com os serviços de autoridade cer-

tificadora (AC), assinatura e validação de documentos, de monitoramento e de autenticação; e (4) módulos de aplicação, contendo componentes de classificação e extração de metadados baseados em aprendizado de máquina, módulo de empacotamento e integração ao RDC-Arq e plataforma de acesso.

## 1.2 OBJETIVOS ESPECÍFICOS

De forma a alcançar o objetivo geral, o trabalho propõe um modelo de arquitetura norteada pelas normas e padrões nacionais e internacionais de gestão arquivística, assim como nas melhores práticas de segurança computacional. Para isso, será criado um ambiente experimental para o processo de digitalização dos arquivos permanentes em suporte analógico, seguindo todas as normas nacionais e internacionais que respaldam uma cadeia de custódia documental. Em seguida, em uma abordagem adversarial, serão feitos testes de penetração e propostas de soluções para mitigar as vulnerabilidades.

Além disso, busca-se também a prototipação de sistema para teste com as seguintes características ou objetivos específicos:

- Proposição de um modelo que unifica diversas áreas de conhecimento originalmente não relacionadas, como a cadeia de custódia de documentos arquivísticos digitais, arquitetura de soluções IoT, aprendizado de máquina aplicada à extração de dados de uma fonte não estruturada, e a aplicação de um modelo adversarial almejando a melhoria da segurança de toda a solução proposta.
- Proposta e validação do protótipo de um componente de software que se integre com scanners via uma API padrão de acesso à hardware desses dispositivos para realizar as seguintes atividades de processamento: captura de imagens de documentos arquivísticos; pré-processamento de imagens e sua preparação para facilitar o reconhecimento óptico de caracteres; integração de reconhecimento óptico de caracteres via ferramenta Tesseract [7]; geração paralela de arquivos PDF a partir do OCR e TIFF multi-páginas com as imagens originais, além de gerar metadados técnicos oriundos do scanner e de propriedades dos arquivos gerados, seguindo as recomendações estabelecidas na estrutura normativa pelo Conarq.
- Proposição e validação de um componente de software que aplica algoritmos de aprendizado de máquina para extração de dados do texto de um documento. O componente se baseia em dois modelos de aprendizado: um com treinamento para a aplicação de redes neurais profundas com intuito de classificar o documento, e outro com treinamento para a aplicação de detecção de objetos com o intuito de extrair metadados sobre o documento (título, data de publicação, autor, etc.).
- Proposta e validação do protótipo de um software baseado em microsserviços para uma

arquitetura de aplicação IoT com o objetivo de controlar a cadeia de custódia documental e serviços de segurança relacionados, validando assinaturas digitais dos pacotes trafegando na cadeia de custódia para confirmar a integridade e autenticidade das mensagens trocadas entre os módulos interoperantes. O controlador baseado em IoT também realiza a gestão automatizada de certificados digitais utilizados para comunicação segura com *Transport Layer Security* (TLS) e para a assinatura digital de documentos.

### 1.3 PUBLICAÇÕES VINCULADAS A ESTE TRABALHO

Durante o desenvolvimento deste trabalho, foram realizadas pesquisas e publicações em diversos domínios da engenharia, as quais refletem na característica multidisciplinar desta obra. Na presente lista de produções, excluimos aquelas que são de carácter introdutório ou preparatório dos diversos estudos realizados. Dessa forma, listamos aqui apenas às publicações que tem vínculo direto com a tese.

- Daniel Alves da Silva, Rafael Timóteo de Sousa Jr, Robson de Oliveira Albuquerque, Ana Lucila Sandoval Orozcoa, Luis Javier Garca Villalba. IoT-based security service for the documentary chain of custody, *Sustainable Cities and Society*, Volume 71, 2021,102940, ISSN 2210-6707, <https://doi.org/10.1016/j.scs.2021.102940> [8].
- Daniel Alves da Silva, Rafael Timóteo de Sousa Jr. IoT-Based Security Service for the Documentary Chain of Custody (IoTsec2CoC):how to protect digital documents in the IoT flow from information producer to consumer? Brasil-Italy webinar, *Flexible and Autonomous Manufacturing Systems for Custom-Designed Products (FASTEN)*, July, Rome, 2020 [9].
- Silva, Daniel Alves da; Torres, José Alberto Sousa; Pinheiro, Alexandre; de Caldas Filho, Francisco L.; Mendonça, Fabio L. L.; Praciano, Bruno J. G; Kfourir, Guilherme Oliveira; de Sousa, Jr, Rafael T. Inference of driver behavior using correlated IoT data from the vehicle telemetry and the driver mobile phone In: *2019 Federated Conference on Computer Science and Information Systems*, 2019, Leipzig. ACSIS, Vol. 18, pp. 487–491 [10].
- Silva, D. A., Machado, P. L., Coelho, V. C. G., Barbosa, R. V., Mendonça, F. L. L., Santos, D. P., de Sousa Júnior, R. T. Produção de indicadores de empregabilidade com base em técnicas de mineração de Big Data e Business Intelligence. Os desafios da qualificação profissional no Brasil: as experiências da escola do trabalhador. *Revista Inclusão Social/Instituto Brasileiro de Informação em Ciência e Tecnologia*. 2019, v. 12, n. 2 . e-ISSN 1808-8678 ISSN impresso 1808-8392 [11].

Vale notar que o primeiro trabalho é considerado A1 pela capes *IoT-based security service for the documentary chain of custody* (Revista *Sustainable Cities and Society* da Elsevier; Fator de Impacto: 7.587 / A1).

## 1.4 METODOLOGIA

A metodologia clássica da pesquisa na engenharia consiste em verificar na literatura a possibilidade de existirem soluções para a problemática aqui tratada, ou soluções para problemas semelhantes. No que tange a estrutura específica da cadeia de custódia de documentos arquivísticos digitais, normatizada por organismos internacionais e nacionais, tais como ISO e Arquivo Nacional, ficou clara a necessidade de conceber uma proposta com foco na segurança computacional.

O trabalho de concepção consistiu em descrever quais são os módulos propostos e o funcionamento deles através do desenho de um diagrama da arquitetura de sistema. Do ponto de vista da metodologia aqui usada, conceber essa arquitetura é muito importante, pois, além de ser o ponto central de orquestração dos diversos módulos e componentes, ela cria um ambiente realista, uma vez que incorpora todo o arcabouço legal e normativo acerca da preservação dos documentos digitais. No entanto, pelo cunho científico do trabalho, faz-se necessária a verificação se tais medidas podem de fato garantir a preservação, autenticidade e não repúdio de tais documentos.

A abordagem para verificar se a arquitetura proposta soluciona o problema não pode ser feita pelo ponto de vista analítico com equações, pois, em primeiro lugar, é muito difícil descrever de forma adequada em equações tantos módulos e componentes em interoperação. Contudo, pode-se, no caso, adotar uma abordagem de simulação, designando cada um desses módulos e componentes e simulando situações em que os documentos são digitalizados. Dessa forma, verifica-se, a cada passo, uma situação documental, porém isso envolveria desenvolver softwares de relativa complexidade dentro de um simulador, sendo necessário o emprego massivo de trabalho de pouca função fora do ambiente simulado.

Por essa razão, optou-se pelo uso de uma metodologia experimental que consiste em desenvolver um protótipo e colocá-lo para operar de forma conjunta, a fim de que a avaliação possa ser feita sobre o protótipo, executando testes em cenários de validação que sejam representativos da realidade, sendo essa a proposta da terceira etapa metodológica, cuja a validação experimental é baseada em um protótipo. A quarta etapa é feita com análise de resultado à luz do modelo adversarial, buscando soluções de segurança frente as vulnerabilidades pelo sistema proposto.

Sendo assim, fica organizada da seguinte maneira a metodologia do trabalho: (1) estudo do problema e da literatura; (2) concepção de uma solução; (3) validação experimental; e (4) verificação se os resultados da validação experimental respondem às necessidades do modelo adversarial.

De forma complementar, foram utilizadas as metodologias Crisp-dm no componente da arquitetura, onde foram aplicadas a Inteligência Artificial (IA) e Metodologia Ágil para o desenvolvimento dos módulos e componentes de software e montagem de hardware. Tais metodologias não serão descritas profundamente aqui, pois já são muito conhecidas, assim como a parte de hardware que também tem seus métodos próprios, de dimensionamento e integração eletrônica que são bastante simplificados, haja vista a utilização de componentes de fácil integração, como os componentes da plataforma Raspberry PI.

## 1.5 ORGANIZAÇÃO DO TRABALHO

Para otimizar o entendimento e a organização do trabalho, os capítulos estão dispostos da seguinte forma:

O capítulo 2 apresenta conceitos básicos relevantes para a compreensão da proposta deste trabalho, revisões do estado da arte e trabalhos correlatos.

O capítulo 3 discute os problemas relacionados à proteção de uma cadeia de produção documental e apresenta as hipóteses da tese.

O capítulo 4 descreve a proposta da arquitetura IoT2SecDCoC com seus: serviços; componentes; fluxos e operações.

O capítulo 5 relata todo o experimento desta tese através da implementação do protótipo da solução IoT2Sec2DCoC em formato de prova de conceito.

O capítulo 6 é dedicado à validação da resiliência do modelo proposto, considerando um modelo adversarial.

O capítulo 7 é dedicado à apresentação e análise dos resultados obtidos.

Finalmente, o capítulo 8 apresenta as considerações finais e propostas para trabalhos futuros.

## 2 ESTADO DA ARTE E TRABALHOS CORRELATOS

Neste capítulo serão expostos os conceitos que norteiam a pesquisa aplicada a este trabalho, entre eles, o documento arquivístico com suas características diplomáticas, normas e desafios para sua preservação segura e garantia de legitimidade. Também se faz necessário explicar os conceitos fundamentais de segurança computacional, IoT (*Internet of Things*) e Inteligência Artificial.

### 2.1 DOCUMENTO ARQUIVÍSTICO DIGITAL

Apenas para utilizar uma definição paradigmática, neste trabalho um documento arquivístico digital é considerado em vista de sua utilização em alguma atividade pública oficial. Por exemplo, usando uma definição que é comum a diversos países: no Brasil, de acordo com o regulamento nacional e-ARQ Brasil [12], um documento arquivístico digital “é um documento digital tratado e gerenciado como um documento arquivístico, ou seja, incorporado ao sistema de arquivos”.

Dessa forma, o documento arquivístico digital tem as mesmas funções probatórias, informativas e garantidoras de direitos, como outros documentos, tanto para organizações como para cidadãos. Portanto, deve ser tratado conforme procedimentos de ciência arquivística como os documentos convencionais e, apesar disso, é claramente distinto de um documento analógico (impresso). O planejamento e a execução de ações para tratar documentos digitais envolvem procedimentos diferentes, sendo necessárias provisões para a conversão de documentos da esfera analógica para a digital.

Segundo o Glossário de Documentos Arquivísticos Digitais [13], outro conceito associado aos Documentos Arquivísticos Digitais de extrema importância é o de Trilha de auditoria definido como “conjunto de informações registradas que permite o rastreamento de intervenções ou tentativas de intervenções feitas no documento arquivístico digital ou no sistema computacional”.

Cabe ressaltar também o conceito de Objeto Digital que, segundo o glossário do Conarq, trata-se de um:

Texto Conjunto de uma ou mais cadeias de bits que registram o conteúdo do objeto e de seus Metadados associados. Constitui-se de três níveis: 1. Nível Físico – Refere-se ao Objeto Digital enquanto fenômeno físico que registra as codificações lógicas dos bits nos Suportes. (e.g. No suporte magnético, o Objeto Físico é a sequência do estado de polaridades negativas e positivas; No suporte óptico, é a sequência de estados de translucidez, transparência e opacidade); 2. Nível Lógico – Refere-se ao Objeto Digital enquanto conjunto de sequên-

cias de bits, que constitui a base dos objetos conceituais; 3. Nível Conceitual – Refere-se ao Objeto Digital que se apresenta de maneira compreensível para o usuário [13].

Evidencia-se que, em seu Nível Conceitual, a fim de o Objeto Digital ser compreensível para o usuário consumidor da informação, é necessário que seus níveis Físico (Hardware) e Lógico (Software) estejam preservados e operacionais.

No entanto, a diferença entre os conceitos de Documento Arquivístico Digital e Objeto Digital fica mais clara ao se analisar o conceito de organicidade, proposto por Rousseau *et al.* [14]: “Informação orgânica é a que foi elaborada, expedida ou recebida no âmbito da missão de um organismo”. Essa informação orgânica é citada nos termos do artigo 2º da Lei nº 8.159/91, em que os seguintes documentos necessitam de tratamento arquivístico:

Texto [...] documentos produzidos e recebidos por órgãos públicos, instituições de caráter público e entidades privadas, em decorrência do exercício de atividades específicas, bem como por pessoa física, qualquer que seja o suporte da informação ou a natureza dos documentos [15].

Considerando o momento de transição e as peculiaridades do documento arquivístico digital, o Conarq aprovou diretrizes para a presunção da autenticidade de documentos arquivísticos digitais, conforme consta na Resolução nº 37, de 19 de dezembro de 2012, em que é delimitado o contexto da necessidade de adaptação, como se segue:

Texto A autenticidade dos documentos arquivísticos digitais é ameaçada sempre que eles são transmitidos através do espaço (entre pessoas e sistemas ou aplicativos) ou do tempo (armazenagem contínua ou atualização/substituição de hardware/software usados para armazenar, processar e comunicar os documentos). Como a guarda de documentos arquivísticos digitais é inexoravelmente ameaçada pela obsolescência tecnológica, a presunção da sua autenticidade deve se apoiar na evidência de que eles foram mantidos com uso de tecnologias e procedimentos administrativos que garantiram a sua identidade e integridade (componentes da autenticidade); ou que, pelo menos, minimizaram os riscos de modificações dos documentos a partir do momento em que foram salvos pela primeira vez e em todos os acessos subsequentes [4].

Mais uma vez tomando o Brasil como exemplo, conforme o artigo 4º do Decreto nº 8.539/2015 [16]: “Os órgãos e as entidades da administração pública federal direta, autárquica e fundacional utilizarão sistemas informatizados para a gestão e o trâmite de processos administrativos eletrô-

nicos.” Esse decreto representa um marco na gestão dos documentos digitais, sendo comum a outros países por tratar de regulamentos internacionais como os do modelo OAIS, apresentado daqui em diante. É interessante notar que, embora a adoção de sistemas informatizados pela administração pública federal direta brasileira tenha ocorrido em tempos remotos, só a contar do momento em que o decreto mencionado entrou em vigor que os órgãos passaram a ter a obrigação legal de adotar medidas para produzir seus documentos em meio digital e digitalizar um grande volume de documentos a partir de seus originais impressos.

A evolução legal brasileira teve prosseguimento com a Lei nº 13.874/2019 [17], a qual reconhece em seu artigo 2º a validade jurídica do documento digitalizado. Buscou-se com essa medida fomentar a transformação digital dos serviços públicos, beneficiando os empreendedores e a sociedade como um todo com essa desburocratização. Dessa forma, é reconhecida a validade jurídica do documento digitalizado.

Mais recentemente foi estabelecido o Decreto nº 10.278/2020 [18]. Este Decreto regulamenta a digitalização de documentos, sejam eles públicos ou privados, garantindo a estes, assim como a documentos nato-digitais, a mesma validade jurídica dos documentos físicos originais, na condição de representantes legais dos originais, desde que observados os requisitos estabelecidos no citado Decreto. Este Decreto passa a ser o novo marco dos documentos digitais, pois esclarece a aplicabilidade e outros pontos que causavam dúvidas no arcabouço legal e normativo até então\*, resguardando a preservação de documentos de guarda permanente, em sua forma original, mesmo após a digitalização, conforme já previa a Lei nº 8159/91 [15].

## **2.2 PRESERVAÇÃO E MANUTENÇÃO DA CADEIA DE CUSTÓDIA DOCUMENTAL**

O projeto InterPARES define a preservação do documento como: “o conjunto de princípios, políticas e estratégias que orienta as atividades prestadas para assegurar a estabilidade física e tecnológica, bem como a proteção do conteúdo intelectual dos materiais, dados, documentos e documentos arquivísticos” [19].

Para que os documentos arquivísticos digitais possam ser preservados em longo prazo e com critérios que garantam a manutenção de sua identidade e integridade ao longo do tempo, assim como a conservação de suas características elementares, torna-se essencial a utilização dos chamados repositórios arquivísticos digitais. Esses repositórios são os únicos que implementam a chamada cadeia de custódia de documentos arquivísticos digitais, ou seja, com eles, os documentos são controlados desde a sua gênese até o processo de preservação e disponibilização da informação. Esse procedimento precisa ser ininterrupto, garantindo a autenticidade dos documentos.

Ao trazermos essa abordagem clássica da cadeia de custódia dos documentos analógicos para os documentos digitais, evidencia-se que os documentos não podem ficar armazenados eterna-



mente na fase de gestão (na imagem representada pelo arquivo corrente e intermediário), e sim que precisam ser exportados para um ambiente confiável e idôneo (representado pelo arquivo permanente, cujo o papel seria de repositório arquivístico digital). Cabe ressaltar que esse processo de exportação deve manter a cadeia de custódia, garantindo a autenticidade e identidade dos documentos. Além disso, as ações de preservação digital precisam ser incorporadas desde o início do ciclo de vida do documento [20].

No trabalho [21], os autores afirmam que “para que um documento arquivístico seja considerado íntegro, é necessário que seja inalterado e completo”. Essa integridade está relacionada diretamente aos ambientes de produção e preservação de documentos, ou seja, aos ambientes custodiadores. A cadeia de custódia confiável de documentos arquivísticos tradicionais é mantida através de uma linha ininterrupta, a qual compreende as três idades do arquivo, sendo elas: corrente, intermediária e permanente. Em decorrência disso, a confiabilidade ocorre por conta da própria instituição que faz a produção, gestão e preservação dos documentos. A Figura 2.1 exemplifica a cadeia de custódia dos documentos arquivísticos tradicionais.

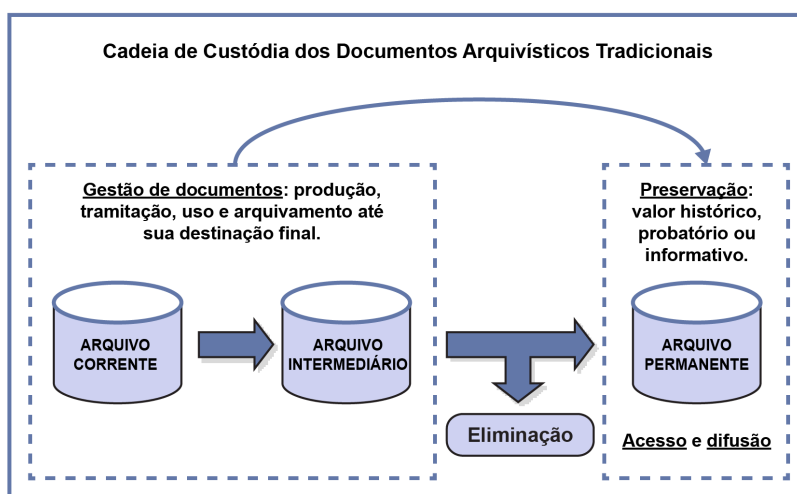


Figura 2.1: Cadeia de Custódia dos documentos arquivísticos tradicionais. Adaptado de: [21].

Essa estruturação de ambientes para os documentos digitais é inclusive recomendada pelos mais recentes estudos conduzidos internacionalmente para a preservação de objetos digitais, tais como: o *Research Library Group*<sup>1</sup> (RLG), o *Online Computer Library Center*<sup>2</sup> (OCLC) e o modelo OAIS. Essas iniciativas e as corretas definições do que vem a ser considerado um repositório arquivístico digital serão elencados nos subtópicos a seguir.

A Resolução nº 39 do Conselho Nacional de Arquivos conceitua o repositório digital como sendo um ambiente de armazenamento e gerenciamento de materiais digitais [22]. Trata-se de uma solução informatizada em que os materiais podem ser capturados, armazenados, preservados e acessados. Um repositório digital é, então, um complexo que apoia o gerenciamento dos materiais digitais pelo tempo necessário, sendo formado por elementos de hardware, software e

<sup>1</sup>Disponível em: <https://www.rlg.org/>.

<sup>2</sup>Disponível em: <https://www.oclc.org/en/home.html>.

metadados, bem como por uma infraestrutura organizacional e procedimentos normativos e técnicos. Tais ambientes são empregados em diversas situações, tais como:

- Arquivo corrente e intermediário (em associação com um SIGAD).
- Arquivo permanente.
- Biblioteca digital.
- Acervo de obras de arte digitais.
- Depósito legal de material digital.
- Curadoria de dados digitais de pesquisa.

### 2.2.1 Modelo OAIS

Motivado pela mesma necessidade de preservar a autenticidade e a integridade de documentos, o modelo OAIS surgiu de uma iniciativa do Comitê Consultivo para Sistemas de Dados Espaciais (CCSDS) da NASA, e seu objetivo inicial era estabelecer padrões capazes de regular o armazenamento em longo prazo de informações digitais produzidas em missões espaciais [23].

A primeira versão do OAIS foi publicada em 1999 [24] e a última em 2012 [23]. Em 2003, tornou-se a norma ISO 14721:2003 [3] e, na última atualização, a ISO 14721:2012 [25]. No Brasil, essa norma foi traduzida e publicada pela Associação Brasileira de Normas Técnicas (ABNT) como NBR 15472:2007 – Modelo de referência para um Sistema Aberto de Arquivamento de Informação - SAAI [26].

O OAIS é um modelo conceitual de referência que visa identificar os componentes funcionais que poderão compor um sistema de informação dedicado à preservação digital, descrevendo as características da interface do sistema e os objetos informacionais que ele deve manter. A Figura 2.2 ilustra o modelo conceitual de OAIS, com funcionalidades, agentes e pacotes de informação.

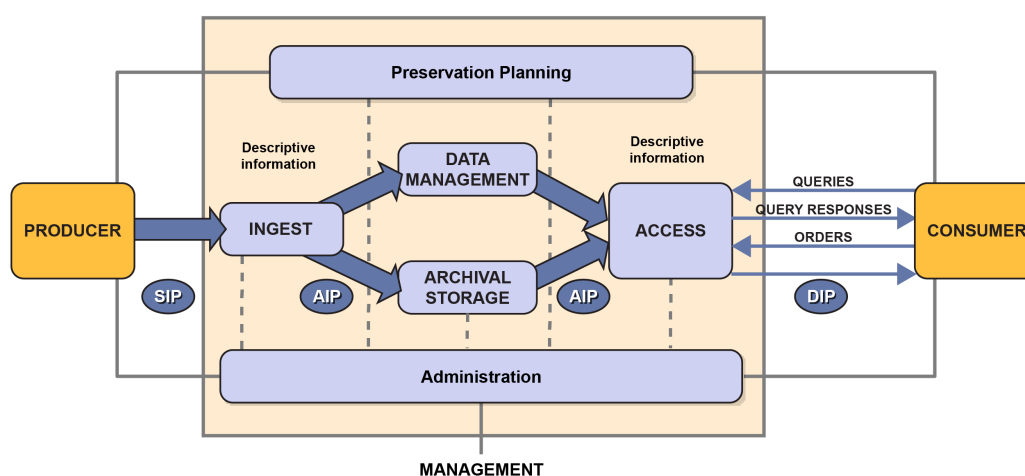


Figura 2.2: Diagrama do modelo de referência OAIS. Fonte: Adaptado de: [23].

Como foi ilustrado na Figura 2.2, o modelo OAIS estabelece seis grupos funcionais no fluxo gerenciado por uma entidade repositório desde o produtor até o consumidor: admissão (ingestão); armazenamento; gestão de dados; planejamento da preservação; administração; e acesso. De acordo com [27], as funcionalidades de cada grupo, conforme a Figura 2.2, são detalhadas no modelo e podem ser acessadas por três tipos de agentes: produtores (pessoas ou sistemas que depositam objetos digitais no repositório); consumidores (pessoas ou sistemas que interagem com as plataformas de disseminação para acessar os objetos digitais); e gestores (responsáveis pelo estabelecimento das políticas e pela gestão dos objetos digitais preservados).

O modelo OAIS prevê a criação de recipientes conceituais chamados de pacotes, os quais armazenam a informação de conteúdo (o documento em si), informação de descrição de preservação (metadados necessários para apoiar a preservação e o acesso ao documento no longo prazo), bem como informações descritivas do pacote (metadados descritivos, possibilitando localizar o pacote no repositório) [27]. As normas de descrição arquivística, como a *General International Standard Archival Description* (ISAD (G)) [28], fornecem esses metadados e a NOBRADE - Norma Brasileira de Descrição Arquivística [29].

Esses pacotes são os meios de garantir a unificação da informação nos ambientes interoperáveis, carregando em seus metadados cada modificação realizada no caminho percorrido pelo documento. Os agentes realizam as seguintes ações com esses pacotes [23]: o produtor realiza a submissão de um Pacote de Submissão de Informação (SIP) contendo documentos e informações de descrição enviados para a entidade de ingestão ou admissão (*ingest*). Após o pacote SIP passar pela etapa da ingestão e ter recebido a informação descritiva, ele se transforma em um Pacote de Arquivamento de Informação (AIP) para armazenamento, o pacote de armazenamento dos documentos de valor permanente. Após o AIP ser armazenado nas entidades de gestão de metadados e repositório de arquivos, é possível gerar o Pacote de Disseminação de Informação (DIP), permitindo aos consumidores pesquisar e acessar os documentos arquivísticos [21].

### **2.2.2 Repositório Arquivístico Digital**

De acordo com o [30], um repositório arquivístico digital é um repositório digital com o objetivo de armazenar e gerenciar um acervo, sejam eles de fase corrente, intermediária ou permanente. Para tanto, esse repositório deve gerenciar os documentos e metadados, com o intuito de proteger as características do documento arquivístico, em especial a autenticidade (identidade e integridade) e a relação orgânica entre os documentos, de acordo com as boas práticas da arquivologia e normas vigentes.

Já um Repositório Arquivístico Digital Confiável (RDC-Arq) é um repositório digital que armazena e gerencia os documentos arquivísticos digitais, transferidos ou recolhidos de sistemas informatizados de gestão, os quais devem cumprir os requisitos definidos no modelo de referência OAIS [3]. A Resolução nº 43 do Conarq prevê como primeiro padrão de preservação de documentos digitais a ser seguido o modelo OAIS [22].

O repositório digital confiável deve ser capaz de manter autênticos os materiais digitais, de preservando-os e provendo acesso a eles pelo tempo necessário. Para cumprir essa missão, segundo o relatório *Trusted Digital Repositories: attributes and responsibilities* [31], os repositórios digitais confiáveis devem:

- Aceitar, em nome de seus depositantes, a responsabilidade pela manutenção dos materiais digitais.
- Dispor de uma estrutura organizacional que apoie, não somente a viabilidade de longo prazo dos próprios repositórios, mas também dos materiais digitais sob sua responsabilidade.
- Demonstrar sustentabilidade econômica e transparência administrativa.
- Projetar seus sistemas de acordo com convenções e padrões comumente aceitos, no sentido de assegurar, de forma contínua, a gestão, o acesso e a segurança dos materiais depositados.
- Estabelecer metodologias para avaliação dos sistemas que considerem as expectativas de confiabilidade esperadas pela comunidade.
- Considerar, para desempenhar suas responsabilidades de longo prazo, os depositários e os usuários de forma aberta e explícita.
- Dispor de políticas, práticas e desempenho que possam ser auditáveis e mensuráveis.
- Observar os seguintes fatores relativos às responsabilidades organizacionais e de curadoria dos repositórios: escopo dos materiais depositados, gerenciamento do ciclo de vida e preservação, atuação junto a uma ampla gama de parceiros, questões legais relacionadas a propriedade dos materiais armazenados e implicações financeiras.

Uma forma de atestar a confiabilidade de um repositório digital junto à comunidade-alvo dá-se por meio da sua certificação por terceiros. Para esse fim, o [31], em parceria com o NASA, publicou, em 2002, o documento *Trustworthy Repository Audit and Certification: Criteria and Checklist (TRAC)*, que serviu de base para a elaboração da norma ISO 16363 [32].

## **2.3 SEGURANÇA COMPUTACIONAL**

Um importante tema a ser tratado, por ser transversal à todo o escopo deste trabalho, é a segurança computacional. Esse termo é definido como toda proteção inserida em um sistema de informação (software, hardware e afins) para preservar a confidencialidade, integridade e disponibilidade do sistema [33]. Essas três características citadas anteriormente são conhecidos como a tríade CIA (*Confidentiality, Integrity and Availability*) da segurança da informação, também sendo chamados de princípios ou pilares:

- Disponibilidade: garantia de que a informação estará disponível para o seu consumidor legítimo sempre que necessário.
- Confidencialidade: garantia de que a informação só estará disponível para consumidores autorizados.
- Integridade: garantia de que a informação não sofreu alterações indevidas durante o seu ciclo de vida.

Além desses princípios básicos da tríade CIA, ainda existem outros princípios que são de extrema importância:

- Autenticidade: garantia da identidade de quem está enviando a informação.
- Não repúdio: garantia de que o autor de uma informação não poderá negar a sua autoria.

### **2.3.1 Ataques**

Ataques buscam explorar vulnerabilidades dos sistemas computacionais, podendo ser intencionais ou não, mas, em ambos os casos, podem comprometer a confidencialidade, integridade e disponibilidade aos usuários. Segundo Stallings [33], as ameaças e ataques têm como objetivo interferir no fluxo informacional entre a fonte e o destino, podendo ter seus tipos classificados em:

- Interrupção: de alguma forma, algo impede que a mensagem chegue no destinatário;
- Modificação: o agente malicioso intercepta a mensagem, altera-a e, então, reenvia-a ao destinatário, passando-se pelo nó remetente.
- Interceptação: o agente malicioso simplesmente intercepta a mensagem enviada pelo remetente.
- Fabricação: o agente malicioso cria uma nova mensagem e envia ao nó de destinatário, passando-se pelo nó remetente.

Esses ataques remetem a algumas importantes propriedades da segurança da informação, sendo elas:

- Disponibilidade: garantia de que a informação estará disponível para o seu consumidor legítimo sempre que necessário.
- Confidencialidade: garantia de que a informação só estará disponível para consumidores autorizados.

- **Integridade:** garantia de que a informação não sofreu alterações indevidas durante o seu ciclo de vida.
- **Autenticidade:** garantia da identidade de quem está enviando a informação.

O *The Open Web Application Security Project (OWASP)* é uma organização sem fins lucrativos a nível mundial que foca na melhoria da segurança de softwares. Eles organizam e divulgam uma lista com os dez modelos de ataques mais aplicados na atualidade, sendo eles [34]:

- *Injection:* avalia falhas de injeção, como SQL e NoSQL, que ocorrem quando dados não confiáveis são enviados por terceiros como parte de um comando ou consulta. Podem ocasionar a execução de comandos não intencionais ou acessar dados sem autorização.
- *Broken Authentication:* avalia se as funções de autenticação e sessão foram implementadas incorretamente, permitindo que invasores consigam manipular senhas, chaves ou *tokens*.
- *Sensitive Data Exposure:* as aplicações não protegem adequadamente dados confidenciais, como informações de cartão, identidade de usuários, entre outros. Essas falhas permitem que invasores cheguem até essas informações, possibilitando-os fraudar cartões, roubar identidade e outros crimes virtuais.
- *XML External Entities (XXE):* alguns processadores antigos de XML avaliam referências de entidades externas em documentos XML. Os invasores podem utilizar dessa falha para divulgar arquivos internos usando o manipulador de *Uniform Resource Identifier (URI)* de arquivo, varredura de porta interna, execução remota de código, entre outros.
- *Broken Access Control:* este ataque acontece quando as aplicações não configuram a autenticação dos usuários de forma correta. Essas falhas permitem que invasores acessem funcionalidades e/ou dados não autorizados, como acessar contas de outros usuários, visualizar arquivos confidenciais, modificar dados de outros usuários, entre outras informações.
- *Security Misconfiguration:* estes ataques acontecem com muita frequência, onde, ao ocorrer falhas nos sistemas, são exibidas mensagens contendo informações confidenciais.
- *Cross-Site Scripting (XSS):* estes ataques ocorrem sempre que um aplicativo permite a inclusão de dados não confiáveis pelos usuários através de APIs do navegador, onde podem ser gerados scripts JavaScript e criação de HTML como formulários.
- *Insecure Deserialization:* a desserialização insegura de dados podem ser utilizadas para executar ataques de repetição, ataques de injeção e ataques de escalonamento de privilégios.
- *Using Components with Known Vulnerabilities:* este ataque é realizado quando componentes, como bibliotecas e módulos, são executados com os mesmos privilégios do aplicativo, podendo facilitar a perda de dados ou até mesmo a tomada do servidor.

- *Insufficient Logging & Monitoring*: a falta de monitoramento de forma eficiente permite que invasores continuem atacando sistemas, alcançando outros sistemas, adulterando, extraindo ou destruindo dados.

Dessa forma, pode-se dizer que ações que buscam garantir a disponibilidade diminuem o risco de interrupção de um serviço ou uma comunicação. Da mesma forma, medidas de confidencialidade mitigam os impactos de ataques de interceptação e modificação, garantindo autenticidade capaz de proteger contra ataques de modificação e fabricação.

### 2.3.2 Criptografia Simétrica e Assimétrica

Para falar de criptografia, antes de tudo, devem ser definidos alguns termos para posterior explicação de diversos conceitos relacionados à segurança computacional [33]:

- Texto claro: mensagem original.
- Texto cifrado: mensagem codificada.
- Cifragem ou Encriptação: aplicação de um algoritmo de encriptação, realizando substituições ou transformações para a conversão de texto claro em texto cifrado.
- Decifragem ou Decriptação: aplicação de um algoritmo de decriptação, realizando o processo inverso da encriptação para a conversão do texto cifrado em texto claro.
- Chave secreta: um valor independente do texto claro e do texto cifrado, servindo como insumo do algoritmo de encriptação ou decriptação para seus processos de substituição ou transformação. Caso seja modificada a chave secreta, o valor do texto cifrado resultante de um mesmo texto claro será diferente.

A criptografia de chave simétrica foi o primeiro tipo de criptografia utilizado. Ele faz uso de uma mesma chave secreta tanto para a encriptação quanto a decriptação. A Figura 2.3 apresenta o processo de criptografia simétrica ilustrado por Stallings [33].

Os algoritmos simétricos mais famosos existentes são o *Data Encryption Standard* (DES) e o *Advanced Encryption Standard* (AES), no qual o último está em grande parte substituindo o primeiro, que era mais antigo e menos confiável. Além dos dois, também existe o algoritmo RC4, um algoritmo de cifragem de fluxo, enquanto os dois anteriores são algoritmos de cifragem de bloco.

Um dos grandes problemas enfrentados na criptografia tradicional de chave simétrica é a questão da distribuição e gerência das chaves secretas, pois a mesma chave deve estar em posse tanto do remetente quanto do destinatário da mensagem, o que não é uma tarefa trivial e tão segura de se realizar. Por conta desse fator, foi desenvolvida a criptografia assimétrica, onde existem duas chaves para cifragem e decifragem, uma que continua sendo secreta (*i.e.*, privada) e outra que

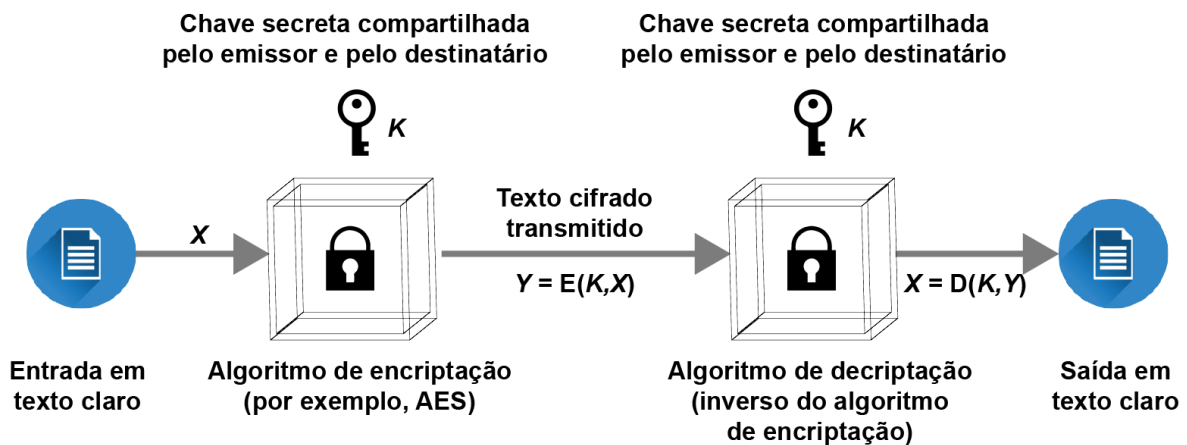


Figura 2.3: Processo de criptografia simétrica. Adaptado de: [33].

será pública. Um texto claro que for cifrado com uma dessas chaves só poderá ser decifrado com o uso da sua outra chave (nunca a mesma), independente se for a pública ou a privada.

A seguir, alguns importantes termos no contexto da criptografia assimétrica:

- Autoridade certificadora: entidade considerada confiável pelas partes envolvidas em uma comunicação segura [35].
- Certificado digital: um registro eletrônico que identifica uma identidade com um conjunto de dados e associa a ela uma chave pública. É emitido por uma autoridade certificadora [35].
- Infraestrutura de chaves públicas: um conjunto de políticas, processos e plataforma tecnológica para administração dos pares de chaves assimétricas, mantendo todo o ciclo de vida das chaves (emissão, manutenção e revogação) [33].

Na Figura 2.4, é ilustrado o processo simplificado da criptografia assimétrica com a realização da cifragem pela chave pública.

### 2.3.3 Assinatura Digital

A assinatura digital é um método criptográfico que permite comprovar a autenticidade e a integridade da informação. Um remetente de uma mensagem consegue obter a *hash* da mensagem e, então, criptografá-la com uma chave privada, que é a assinatura digital da mensagem. Após isso, a mensagem e a sua assinatura são enviadas ao destinatário e, assim, a partir da chave pública do emissor, é possível realizar a decifragem de função resumo (*i.e.*, função *hash*) e validar a integridade da mensagem. Além disso, como o processo de decifragem do dado só pode ser realizado com a chave pública associada a chave privada do remetente, ainda se consegue comprovar a autenticidade da mensagem.



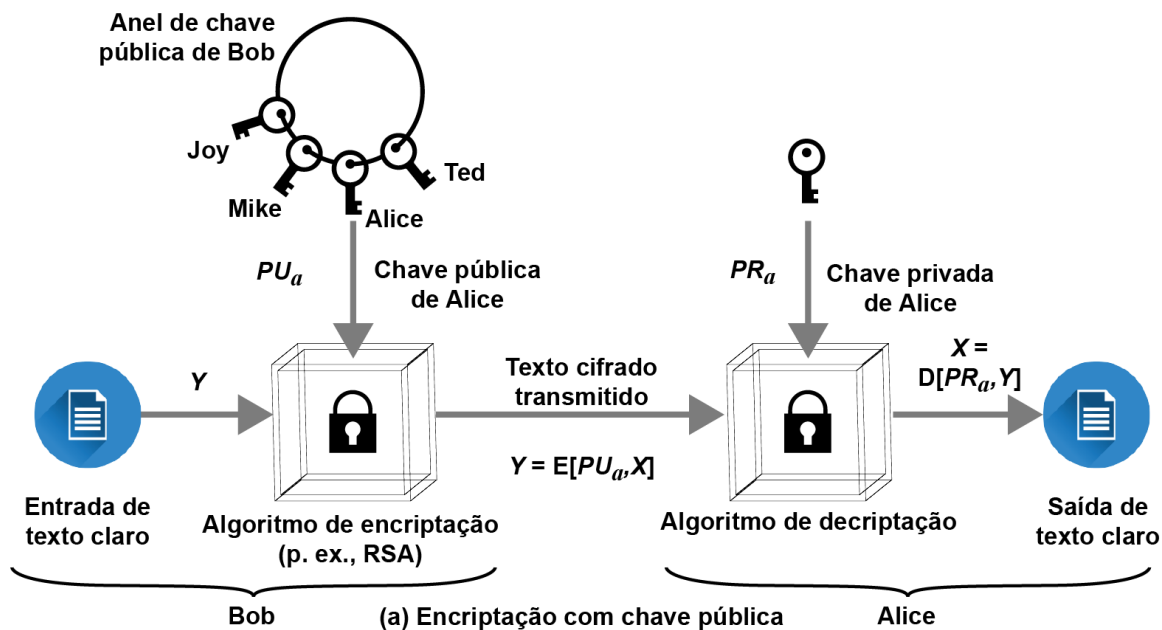


Figura 2.4: Processo de criptografia assimétrica. Adaptado de: [33].

Ao ser aplicada uma função *hash* em qualquer tipo de informação, independente do tamanho que tenha, ela irá gerar um resultado de tamanho fixo [35]. Ademais, uma função *hash* considerada “boa” irá gerar resultados aparentemente aleatórios e únicos a partir de uma simples mudança de 1 *bit* na informação original.

Por conta dessas características, o resultado da função *hash* pode ser utilizado para averiguar a integridade das informações. A partir da informação original e de um *hash* previamente calculado, no momento que se julgar necessário, poderá ser calculado um novo *hash* para comparação e verificação da integridade da informação.

### 2.3.4 Gerenciamento de Eventos e Informações de Segurança

O gerenciamento de eventos que ocorrem no contexto de uma solução, ou até mesmo de uma organização como um todo, é de suma importância quando se trata de segurança da informação. Um efetivo gerenciamento dos registros de log das aplicações pode proporcionar diversos benefícios, entre eles: o armazenamento a longo prazo das informações para futuras auditorias, o cruzamento/correlação de informações entre diversos componentes de uma solução, o monitoramento em tempo real de eventos de segurança, a agregação de informação para análise, a geração de alertas em casos de suspeita de um incidente de segurança, etc. Uma solução composta por softwares livres que realiza a tarefa de gerenciador de logs e eventos é a pilha *Elasticsearch*, *Logstash* e *Kibana* (ELK). O *Elasticsearch* é um banco de dados NoSQL especializado em consultas por texto livre. O *Logstash* é responsável por realizar a coleta, indexação, filtragem e agrupamento dos registros de log. O *Kibana* atua como um *dashboard* de controle para visualização dos

dados armazenados e criação de painéis analíticos sobre os logs gerados.

## 2.4 COMPUTAÇÃO DISTRIBUÍDA

Com o passar dos anos e o avanço nas pesquisas computacionais, a quantidade de dados e a complexidade dos cálculos são cada vez maiores. Dessa maneira, Sloan [36] propõe três abordagens para melhorar a performance: utilizar um algoritmo melhor e mais eficiente; utilizar um computador mais rápido (mais memória e/ou processamento), ou dividir o processamento em múltiplos computadores. Neste contexto, surgem os sistemas distribuídos, que, segundo Coulouris *et al.* [37], os sistemas distribuídos podem ser caracterizados como um conjunto de computadores autônomos interligados em rede com software capaz de produzir um mecanismo integrado de computação.

### 2.4.1 Computação em Nuvem

A computação em nuvem é um modelo de computação distribuída que, segundo Buyya *et al.* [38], pode ser definido como:

Texto um sistema de computação paralela e distribuída que consiste em um conjunto de computadores interligados e virtualizados que são dinamicamente providos e apresentados como um ou mais recursos de computação unificada baseada em contratos de níveis de serviço estabelecidos através de negociação entre o prestador de serviço e os consumidores. [15].

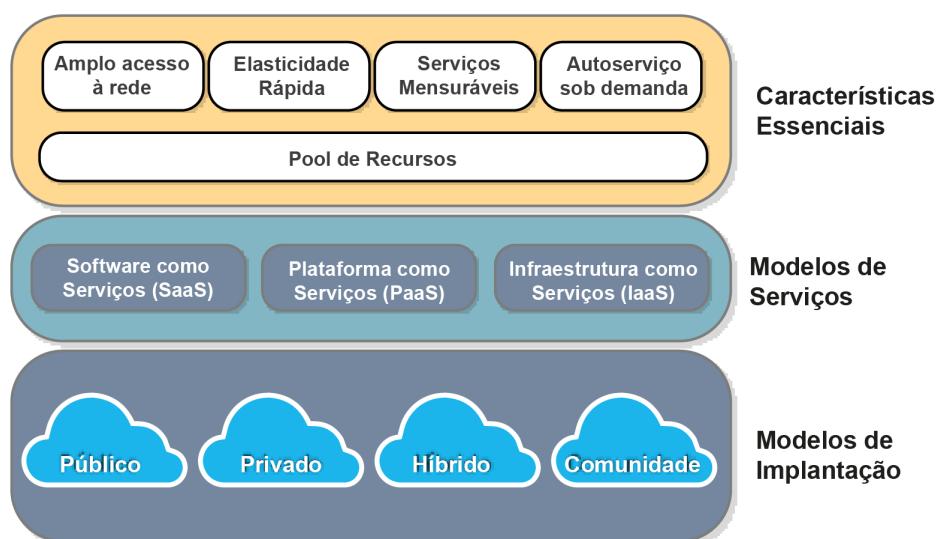


Figura 2.5: Modelo NIST de definição de Computação em Nuvem. Fonte: [39].

Segundo o modelo do National Institute of Standards and Technology (NIST) [39], a computação em nuvem possui algumas características essenciais:

- Auto-atendimento sob demanda (*On-Demand Self-Service*): o consumidor pode consumir os serviços da nuvem, tais como armazenamento e processamento, de forma direta e transparente, sem precisar contatar nenhum administrador da nuvem.
- Amplo acesso à rede (*Ubiquitous Network Access*): todos os serviços da nuvem podem ser acessados por qualquer dispositivo que tenha acesso a internet, facilitando o acesso do consumidor.
- *Pool de Recursos (Resource Pooling)*: os recursos físicos da nuvem são virtualizados para que possam servir a múltiplos usuários utilizando o modelo (*multi-tenancy*).
- Elasticidade Rápida (*Rapid elasticity*): a elasticidade é a capacidade da nuvem se expandir, tanto verticalmente (aumentando as CPU's ou memória) quanto horizontalmente, criando novos servidores, visando a um maior aproveitamento das capacidades disponíveis.
- Serviços Mensuráveis (*Measured Service*): os serviços da nuvem devem ser controlados pelo provedor para que possa haver cobrança pelo uso, permitindo que o cliente seja cobrado exatamente pelo que consumiu.

A computação em nuvem possui três modelos de serviço essenciais:

- Infraestrutura como Serviço (*Infrastructure-as-a-Service - IaaS*): neste modelo, o provedor de serviço oferece a infraestrutura computacional como um serviço, tal como servidores virtuais ou armazenamento. O consumidor tem total controle da nuvem.
- Plataforma como Serviço (*Platform-as-a-Service - PaaS*): este modelo permite ao consumidor da nuvem utilizar serviços e aplicações da nuvem, como ambientes de desenvolvimento ou implantação, diretamente na nuvem. O consumidor pode controlar alguma configuração de aplicação, porém não controla a nuvem.
- Software como Serviço (*Software-as-a-Service - SaaS*): este modelo oferece aplicações totalmente online em forma de serviço, o software funciona totalmente na *Web*. O consumidor não controla nenhuma configuração da nuvem.

## 2.4.2 Edge Computing

A principal ideia da computação de borda (*edge computing*) é trazer o processamento dos dados para geograficamente mais perto da fonte dos dados, pré-processar e depois enviar para a nuvem, permitindo uma resposta mais rápida.

Segundo Shi *et al.* [40], “a computação de borda refere-se a tecnologias de ativação que permitem com que a computação seja executada na borda da rede, em dados *downstream* em nome

de serviços em nuvem e *upstream* dados em nome dos serviços IoT”. A Figura 2.6 ilustra esse processamento em duas vias na computação de borda.

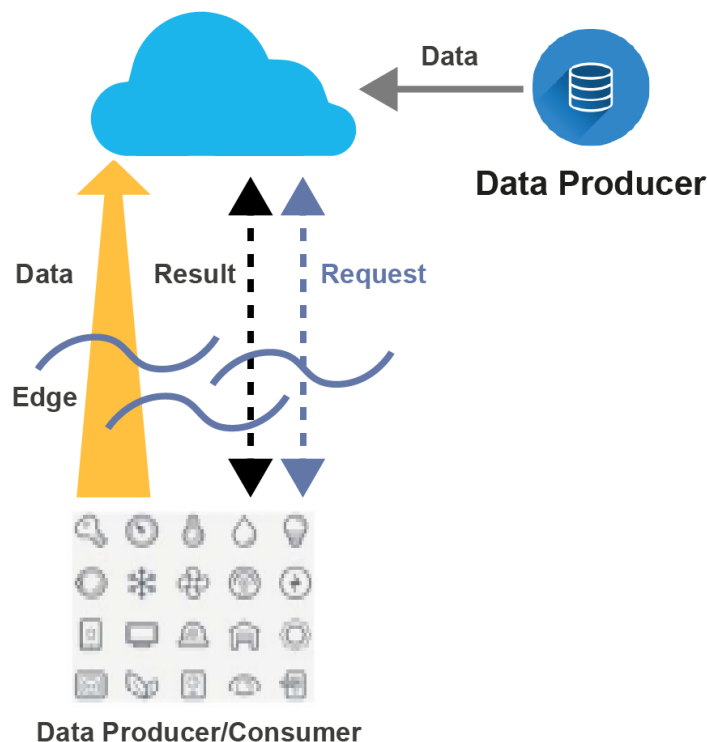


Figura 2.6: O paradigma de computação de borda. Fonte: [40].

Nesse paradigma, “as coisas” não são necessariamente os consumidores. Funcionam também como produtores de dados, não só requisitando algum serviço, mas também executando tarefas da nuvem, como armazenamento de dados, *caching* ou distribuindo as requisições dos usuários.

### 2.4.3 Microsserviços

Microsserviços são uma abordagem arquitetural de desenvolvimento de aplicação pequenas e autônomas, fáceis de manter e que possuem uma única responsabilidade. Segundo Nerman [41], “microsserviços são pequenos serviços autônomos que trabalham em conjunto, e toda comunicação entre eles deve ser feita via chamadas de rede para forçar a separação entre os serviços e evitar os perigos de forte acoplamento”.

A implementação de microsserviços tem grandes vantagens, como a possibilidade da escrita dos microsserviços em diferentes linguagens de programação, possibilitando aplicar o que de melhor cada linguagem oferece quanto aos seus objetivos, além da simplificação de *deploy* dos microsserviços e fácil escalabilidade do sistema, pois, caso necessário, é possível escrever novos microsserviços que adicionam funcionalidades novas ou intermediárias, sendo mapeadas quanto à infraestrutura.

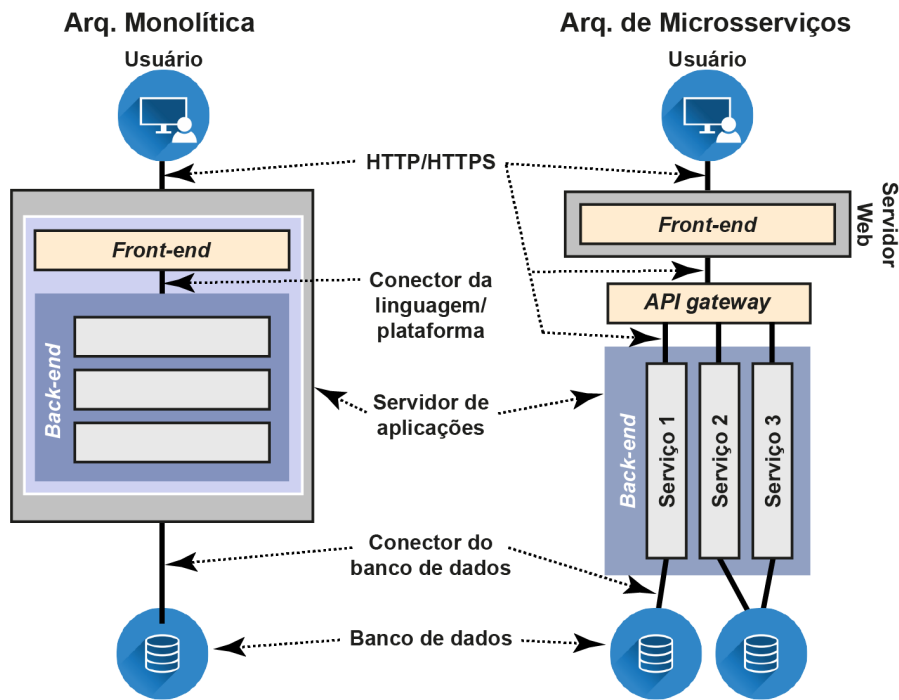


Figura 2.7: Comparação entre as arquiteturas Monolítica e Microserviços. Fonte: [42].

A Figura 2.7 destaca a diferença entre a abordagem monolítica e a de microserviços. Enquanto a estrutura monolítica utiliza a mesma infraestrutura para todas as funcionalidades, a estrutura de microserviços utiliza estruturas separadas, minimizando, assim, o ponto único de falha, já que cada microserviço utiliza sua própria infraestrutura. Um exemplo é uma aplicação que tenha a funcionalidade de cadastro de usuários e de login caso o serviço de cadastro de novos usuários deixe de funcionar, o serviço de login continua funcionando normalmente, diferentemente da estrutura monolítica, em que todos os serviços deixariam de funcionar.

#### 2.4.4 Internet das Coisas e Industria 4.0

Segundo a União Internacional de Telecomunicações (UIT): “IoT é uma infraestrutura global para a sociedade da informação, permitindo serviços avançados através da interconexão de coisas (físicas ou virtuais) com base nas tecnologias de informação e comunicação interoperáveis existentes e/ou em desenvolvimento” [43].

A IoT proporciona uma nova possibilidade de uso de várias tecnologias já existentes, principalmente as relacionadas a sensores e atuadores, sejam eles físicos e virtuais, processados ou não, mas capazes de capturar, processar, armazenar, transmitir e fornecer informações de forma que o sistema resultante seja infinitamente maior que sua soma, sendo possível através da inteligência extraída dos dados.

De acordo com Ferreira *et al.* [44], as interações entre objetos permitirão que prestem serviços complexos para as pessoas sem necessariamente requerer intervenção humana. Dessa forma, os

objetos inteligentes sairão do foco do cotidiano do homem, ficando em execução de segundo plano, e, por isso, tornar-se-ão pervasivos.

Com isso, pode-se concluir que, para que um sistema seja considerado IoT, é necessário algum nível de inteligência, seja de um aprendizado de máquina ou outro método que busque responder questões ou orientar ações através do processo de transformar dados brutos em informações que, por sua vez, podem se transformar em aprendizado, gerando algum tipo de conhecimento.

Segundo o entendimento de Silva *et al.* [45], a Internet das Coisas se dá através de entidades que atuam como fornecedores e/ou consumidores de dados relacionados com o mundo físico, dando ao IoT um foco maior em dados e informação. Apesar da faceta em tecnologias de comunicação e seus protocolos, não há perda de características de integração e de funcionalidades e recursos fornecidos pelos objetos inteligentes.

Essa nova rede de alta escalabilidade de bilhões de objetos inteligentes interligados que se multiplicam organicamente pode proporcionar a pessoas, empresas e governos o aumento da eficiência na administração de suas áreas de atuação, tais como saúde, educação, segurança, indústria e transporte, de forma sustentável, ou seja, preservando o meio ambiente através do uso racional e inteligente dos recursos naturais.

Com o advento das mesmas tecnologias disruptivas que produziram a IoT [46], o chamado paradigma Indústria 4.0 se tornou uma realidade e tem auxiliado a produção com sistemas integrados e inteligentes. Dispositivos industriais inteligentes e conectados auxiliam na produtividade, nos custos, no controle e monitoramento de processos e ativos operacionais.

Conforme citado por Khan *et al.* [47], a IoT Industrial (IIoT) é um subconjunto da IoT que exige maior segurança, proteção e comunicação confiável sem interromper as operações industriais em tempo real, dada sua aplicação em ambientes industriais de missão crítica. De acordo com Esposito *et al.* [48], as soluções IoT existentes precisam de adaptações para focar nas peculiaridades da notificação de incidentes e nas características dos dispositivos com recursos limitados que são frequentemente usados na indústria.

A conversão de documentos arquivísticos do suporte analógico para o digital requer uma linha de produção de processamento de documentos (limpeza, preparação, digitalização e indexação). Em seguida, esses documentos, bem como outros documentos nato digitais, estão sujeitos a transferências subsequentes e processamento como aqueles ilustrados na Figura 2.2. A adição de sensores e controles, tanto na forma de dispositivos quanto de módulos de software, para gerenciar essa linha de produção constitui em uma encarnação da IIoT dedicada à cadeia de produção documental que é considerada neste trabalho.

#### **2.4.5 Middleware IoT**

Segundo Cruz [49], a arquitetura elementar de um sistema IoT é composta por três camadas: camadas de aplicação, Middleware IoT e a camada física. A interoperação entre essas camadas e

seus componentes ocorre por intermédio do Middleware IoT, que recebe as mensagens da camada física, executa o processamento e encaminha para a aplicação, sendo que o mesmo fluxo pode ocorrer de forma invertida. A Figura 2.8 ilustra esse formato simplificado da arquitetura IoT.

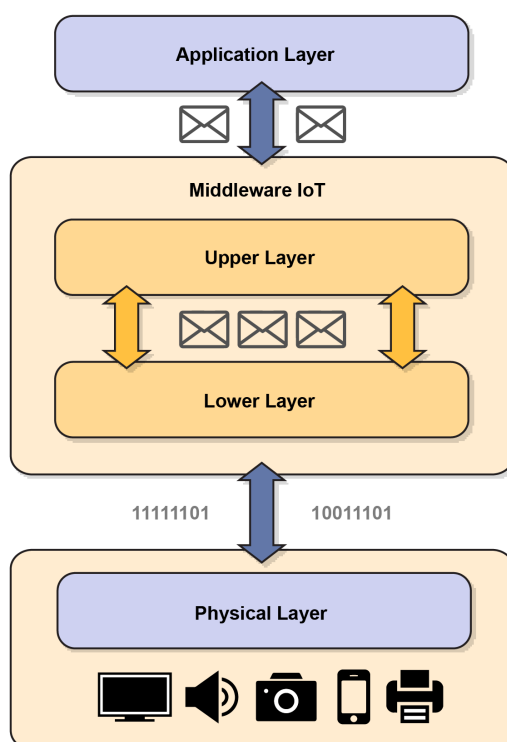


Figura 2.8: Estrutura do Middleware IoT. Fonte: Adaptado de: [50].

Segundo Ferreira *et al.* [51], as Arquiteturas de Middleware para IoT são estruturas físicas e lógicas complexas que devem ser modularizadas e dotadas de todas as tecnologias relevantes disponíveis. A função principal do Middleware, segundo o autor, é integrar dispositivos e aplicações, através de rotinas de controle e de monitoramento de estado de dispositivos.

O Middleware também pode ser utilizado com múltiplas instâncias em arranjo hierárquico. Em seu trabalho, Menezes *et al.* [52] utiliza uma arquitetura de modelo hierárquico de processamento, escalonando instâncias do Middleware em diversos níveis para fins de distribuição do processamento (com fundamentos de *edge computing*) e, conseqüentemente, melhor utilização de recursos computacionais, entre outras funcionalidades adicionais que são executadas pelo novo módulo denominado Engine.

O modelo hierarquizado com uso de múltiplas instâncias do Middleware serve para alcançar o melhor uso de recursos computacionais da solução. De acordo com Tong *et al.* [53], a distribuição do processamento nas bordas (i.e. *edge computing*) da rede evita a sobrecarga e latência da solução, ou seja, ao invés do envio de todos os dados para a instância raiz do Middleware na nuvem principal, os dados são previamente processados em instâncias intermediárias que estão geograficamente mais próximos dos dispositivos inteligentes.

Uma representação simplificada da arquitetura hierárquica com duas instâncias do Middleware pode ser visualizada na Figura 2.9.

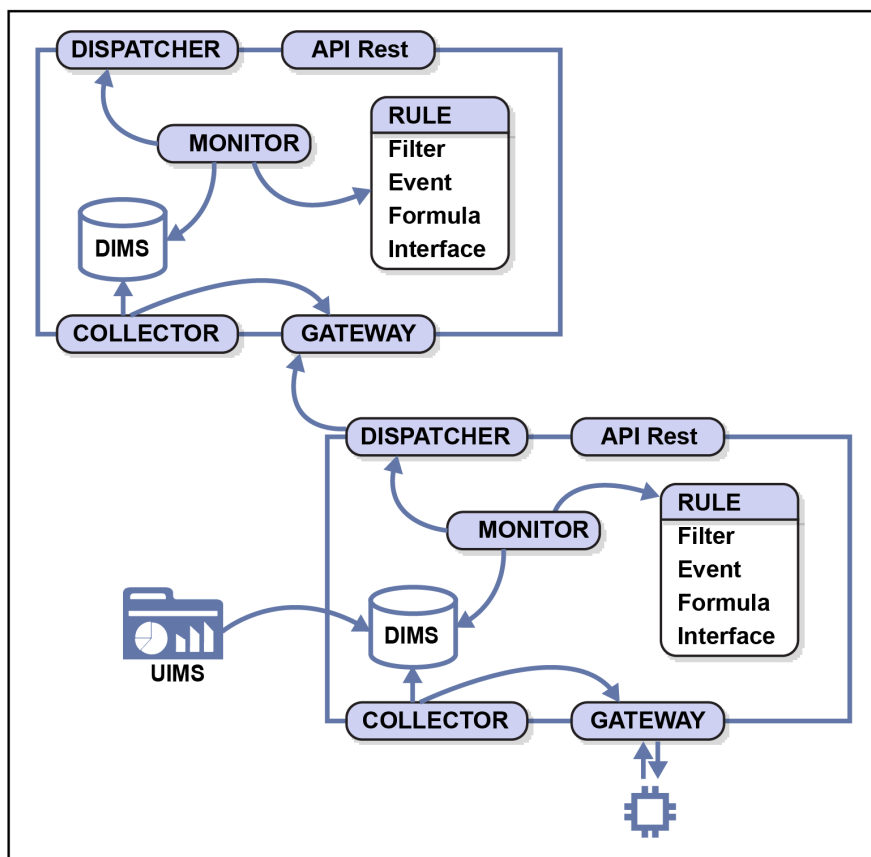


Figura 2.9: Arquitetura do Middleware hierárquico do UIoT. Fonte: Adaptado de: [52].

Esse Middleware possui como principais entidades de sua lógica os seguintes elementos:

- Client: representação do dispositivo que está inserido na rede IoT, sendo identificado pelo seu MAC Address e Chipset.
- Service: representação do serviço que um determinado dispositivo tem da capacidade de enviar dados para o Middleware, e.g. “Temperatura”. Cada client pode possuir de zero a vários serviços associados.
- Data: representação do dado que um determinado dispositivo envia para o Middleware, ficando associado à um determinado serviço, e.g. “10 °C”.
- Rules: representação das regras de processamento que podem ser executadas pelo componente Engine, e.g. “A partir de todos os dados de temperatura maiores que 40 °C, calcule a média dessas temperaturas e a envie para o Middleware UIoT na hierarquia superior”.

Para a persistência de dados do UIoT, é utilizando um modulo chamado de DIMS, que é uma aplicação web Python desenvolvida com o Framework Flask. Ela fornece uma API REST com



operações para criação e consulta de *client's*, *service's*, *data's* e *rule's*. O banco de dados utilizado para armazenamento das entidades é o banco NoSQL MongoDB. Originalmente, o UIoT não oferecia suporte para manipular dados textuais, trabalhando basicamente com dados numéricos, mas, no contexto do trabalho, ele foi evoluído para trabalhar com esse tipo de dado. Essa alteração será melhor descrita na seção 5.3.

#### 2.4.5.1 Gateway do Middleware UIoT

Segundo Martins *et al.* [54], o gateway do Middleware UIoT é o principal módulo dentro da arquitetura do UIoT. Trata-se de uma aplicação Python responsável por fornecer os protocolos de comunicação para integração com os dispositivos inteligentes, possuindo uma gama variada de integrações inteligíveis pelo módulo. A integração pode ser realizada por meio de: *sockets* TCP ou UDP; protocolo de mensageiria MQTT; protocolo de redes sem fio ZigBee; ou por meio de dispositivos com maior poder de processamento até mesmo por protocolo HTTP via sua API REST. Todas essas formas de integração são chamadas de *listeners* dentro do módulo.

Por conta do seu modelo hierarquizado, um gateway também poderá receber dados provenientes de uma outra instância inferior do UIoT (via chamadas REST), não ficando restrito somente a receber dados de dispositivos inteligentes, mas também de outras instâncias do Middleware da solução. Após o recebimento dos dados, o Gateway UIoT, em sua implementação original [54], envia-os para o *Data Interface Management System* (DIMS), via sua API REST.

Para a visualização dos dados, é utilizado um módulo responsável chamado de *User Interface Management System* (UIMS). Esse módulo consiste em um *dashboard* básico de visualização das entidades cadastradas no UIoT, sendo um módulo opcional na arquitetura do Middleware. Essa aplicação *Single Page Applications* (SPA) acessa o DIMS diretamente para possibilitar a visualização dos dados.

#### 2.4.5.2 Engine

O Middleware Hierárquico UIoT de Menezes *et al.* [52] incluiu em seu funcionamento um novo módulo, uma aplicação Python que age como um motor de execução de regras de processamento chamado Engine.

Nessa arquitetura, o gateway, ao receber dados dos dispositivos inteligentes, armazena-os em um *buffer*. O componente interno do Engine, denominado Collector, periodicamente coleta os novos dados no *buffer* do gateway e os insere no DIMS, para que os demais componentes possam consultar os dados.

Outro componente interno do Engine, denominado Monitor, executa periodicamente as regras cadastradas no DIMS, sendo que, se um evento definido na regra ocorrer, a regra é executada e, então, os dados que atendem a ela podem ser enviados para um novo local via o componente interno Dispatcher, por exemplo, enviando os dados para o middleware da hierarquia superior.

### 2.4.5.3 Autenticação, Comando e Controle

A autenticação permite a formação da rede IoT e, por consequência, a comunicação M2M (*Machine to Machine*). A partir do ingresso dos dispositivos na rede, é possível o comando e controle desses dispositivos, possibilitando o consumo de dados e serviços.

A autenticação pode ocorrer essencialmente de duas maneiras: centralizada ou descentralizada. Na abordagem centralizada, um conjunto de nós na rede tem a função da identificação dos dispositivos e serviços, bem como o gerenciamento destes. Um desses nós também poderia assumir essa responsabilidade sozinho, sendo necessário, para tal, uma maior escalabilidade de recursos computacionais.

A trabalho, Silva *et al.* [45] apresenta uma proposta de autenticação descentralizada, em que o próprio dispositivo possui controle sobre os serviços que fornece, atualizando a rede sempre que existirem novas informações relativas a seus serviços. Tal característica proporciona o aumento da escalabilidade da rede, assim como melhor controle e integração entre dispositivos heterogêneos. Em sua abordagem, através de comunicação com uso de JSON, os dispositivos devem enviar um documento com as informações necessárias para três tipos primários de interação:

- Documento de identidade

Código 2.1: Padrão do documento JSON para Identidade.

```
id_cps : 4C4C4544004731108047B4C04F4C3232,  
id_prc : 51 06 04 00 FF FB EB BF,  
hd_srl : W761TTGL, driver : ethernet,  
mac    : 74:e6:e2:ce:23:6d, host : rpy-iot
```

- Documento de registro

Código 2.2: Padrão do documento JSON para registro.

```
id_token : 11768768, mac : 74:e6:e2:ce:23:6d,  
services:[{name : get_temp, type : float,  
unit : celsius, desc : temp sensor}]
```

- Documento de atualização

Código 2.3: Padrão do documento JSON para Identidade.

```
id_token : 11768768, mac : 74:e6:e2:ce:23:6d  
services:[{id_serv : 001, value: 24}]
```

## 2.5 INTELIGENCIA ARTIFICIAL E APRENDIZADO DE MÁQUINA

A Inteligência Artificial (IA) é uma técnica que simula os processos de inteligência humana através de máquinas que processam cálculos matemáticos e produzem um resultado similar ou igual a de um ser humano. Algumas das áreas que utilizam essa técnica vem crescendo e se destacando como o reconhecimento de fala, visão computacional e *deep learning* [55]. A Inteligência Artificial baseada em aprendizado de máquina permite que as entidades que possuem grandes quantidades de dados disponíveis extraiam conhecimentos que seriam praticamente impossíveis para qualquer pessoa descobrir. As quatro grandes categorias da IA são: máquinas que pensam como seres humanos, máquinas que se comportam como seres humanos, máquinas que pensam racionalmente e máquinas que se comportam de forma racional [56].

O aprendizado supervisionado ou de máquina é uma técnica de aprendizado para deduzir uma função de dados de treinamento, os quais consistem em pares de objetos de entrada (tipicamente vetores) e saídas desejadas. A saída da função pode ser um valor contínuo (chamado de regressão) ou pode prever um rótulo de classe do objeto de entrada (chamado classificação). A tarefa do aprendizado supervisionado é prever o valor da função para qualquer objeto de entrada válido, depois de ter visto um número de exemplos de treinamento. Para conseguir isso, o aprendizado supervisionado tem que replicar, a partir dos dados apresentados, para situações não vistas de uma forma razoável.

### 2.5.1 Redes Neurais

Uma rede neural é uma estrutura que aplica uma série de algoritmos que trabalham para reconhecer padrões em um conjunto de dados através de um processo que imita a maneira de como o cérebro humano trabalha. Uma rede neural é composta por camadas conectadas que tem como base a camada de entrada, onde são passados os dados a serem processados; uma ou mais camadas ocultas, que são responsáveis pelos cálculos matemáticos; e camada de saída, que apresenta o resultado do processamento. A composição de uma rede neural possui, ainda, nós interconectados chamados de neurônios, os quais transmitem a informação entre si, assim como o cérebro humano [57].

#### 2.5.1.1 Multi Layer Perceptron

Uma Multi Layer Perceptron - MLP é um tipo de rede neural que utiliza a técnica de aprendizado supervisionado para realizar o seu treinamento. É também um tipo de rede neural *feed-forward*, ou seja, a rede é alimentada em um único sentido, da camada de entrada até a camada de saída. Para ser caracterizada como uma rede MLP, ela precisa ter no mínimo três camadas (camada de entrada, uma camada oculta e uma camada de saída), além disso, cada nó das camadas ocultas e das camadas de saída devem possuir uma função de ativação não linear [58].

### 2.5.1.2 *Deep Learning* e Visão Computacional

*Deep learning* é um subcampo da IA proveniente do aprendizado de máquina relacionado a algoritmos que projetam redes neurais artificiais e aprendem com grandes quantidades de dados, ou seja, é uma rede neural que possui muitas camadas ocultas, dando a ideia de profundidade. No modelo de aprendizagem profunda, o algoritmo aprende a executar tarefas como classificação de imagens, texto ou som, podendo alcançar uma precisão bastante alta, superando o desempenho humano [59].

### 2.5.1.3 Classificação de Imagens

A classificação de imagens é um campo da IA que pertence à visão computacional, sendo esta é capaz de analisar uma imagem e identificar a qual classe ela pertence. Seu aprendizado é supervisionado e, para realizar seu treinamento, o modelo recebe um conjunto de imagens e suas classes de destino, possibilitando ao modelo identificar novas imagens ainda não conhecidas por ele.

### 2.5.1.4 Detecção de Objetos

Uma das formas de valer-se da detecção de objetos é utilizando o algoritmo *You look once* (YOLO) [60], no qual é utilizado conceitos de redes neurais convolucionais profundas para realizar as tarefas de detecção e classificação de objetos.

O conceito de redes neurais é descrito neste trabalho na subseção 2.5.1. A arquitetura YOLO possui 75 camadas de convolução e camadas auxiliares de *upsampling* e *downsampling*.

O conceito de *bounding boxes* é apresentado por Redmon e Farhadi (2017) [61], e pode ser utilizado em vários outros domínios, como veículos autônomos, detecção de pessoas e também em detecção de termos chaves em documentos.

A detecção de objetos é baseada em *bounding boxes*, e cada uma dessas caixas tem cinco elementos conhecidos, sendo eles  $x, y, w$  e  $h$ , sendo estes os atributos de posicionamento dessa caixa de detecção e, por último, temos a variável que nos informa o resultado dessa detecção, um booleano chamado de score de confiança, ou seja, apenas fala se existe um objeto ou não. Utiliza-se das redes neurais convolucionais para realizar a redução espacial da imagem e, após isso, combina-se a regressão linear usando outro tipo de redes neurais para que seja feita a predição. Para esse tipo de arquitetura, a probabilidade acima de 0.5 é considerada como correta para existência de objetos nesse *grid*. A Tabela 2.1 contém as principais equações para essa detecção.

A forma utilizada é da seguinte maneira: a rede neural divide a imagem de entrada em  $Z \times Z$  *grids*. Cada *grid* desse *frame* de imagem prevê apenas um único objeto. A configuração padrão é de  $7 \times 7$  *grids*, duas *bounding boxes* (B) e 20 classes de imagens. Portanto o resultado após o uso é um tensor de predição com o formato de  $(Z, Z, B \times 5 + 20) = (7, 7, 30)$ .

Tabela 2.1: Equações de predição do YOLO [62].

Descrição	Equações
Escore de Confiança da caixa (B)	$P_r(object) \cdot IoU$
Probabilidade Condicional para cada classe por objeto (PC)	$P_r(class_i object)$
Escore de Confiança da Classe (C)	$P_r(class_i) \cdot IoU$
Equação final	$B \cdot C$

Na Tabela 2.1,  $P_r(object)$  é a probabilidade da caixa conter o objeto.  $IoU$  é a interseção sobre a união entre a caixa detectada com o objeto de fato.  $P_r(class_i|object)$  é a probabilidade de um objeto pertencer à determinada classe.  $Class_i$  é obtido através da presença de um objeto  $P_r(class_i)$ .

Na Administração Pública, a grande maioria dos documentos não são padronizados, o que dificulta o desenvolvimento de um algoritmo que seja genérico. Portanto, é necessário fazer uso de técnicas de inteligência artificial, onde um modelo de aprendizado de máquina é capaz de se adaptar para detecção de estruturas que podem ser previamente definidas. Essa detecção ocorre após a extração das informações de um arquivo em formato PDF e, com o devido tratamento, é possível detectar uma tabela e extrair essa informação. Segundo os autores de *A yolo-based table detection method* [63], é necessária a utilização de pré e pós processamento desse texto para que haja uma uniformidade no momento da detecção.

## 2.6 TRABALHOS RELACIONADOS

Durante a fase de pesquisa bibliográfica, utilizando repositórios internacionais de pesquisa acadêmica, até onde sabemos, não foi encontrado nenhum trabalho de pesquisa abrangendo, de modo geral, uma proposta similar usando uma IIoT para fornecer serviços de segurança para garantir a cadeia de custódia de uma transferência documental, acesso e gerenciamento de linha de produção. No entanto, em relação aos módulos e componentes da solução proposta, há trabalhos correlatos interessantes, apesar de aplicados a outras situações, como será discutido nas subseções seguintes.

### 2.6.1 Confiança e Cadeia de Custódia

Como confiança computacional, modelos de confiança e gestão de confiança têm sido objetos de diversos esforços de pesquisa [64]. A definição de confiança utilizada neste trabalho tem o objetivo de formalizar e esclarecer aspectos de confiança em protocolos de comunicação. Foi adotado o conceito de confiança de Yahalom *et al.* (1993) [65], o qual afirma que a noção de confiança é fundamental para a compreensão das interações entre entidades como dispositivos de computação, seres humanos, organizações, nações e outros. De acordo com essa referência, o fato de que uma entidade A confia em uma entidade B em algum aspecto, informalmente, significa

que A acredita que B se comportará de uma certa maneira e executará alguma ação sob certas circunstâncias específicas. De acordo com Adnane *et al.* [64], essa definição implica que uma relação de confiança é estabelecida a partir da possibilidade de conduzir uma ação protocolar, a qual é avaliada pela entidade A com base no que ela conhece sobre a entidade B e as circunstâncias da operação. Considerando a ação em questão e as circunstâncias de execução, é necessário distinguir diferentes classes de confiança [65]. No que diz respeito a uma CoC de documentos, conforme tratado no presente trabalho, as classes adequadas para as ações realizadas por nossa proposta incluem a confiança em outra entidade para autenticar a origem dos documentos e atestar a integridade desses documentos.

Para uma lista detalhada de outras classes de confiança úteis na CoC, a referência [66] apresenta o uso da CoC em perícia, argumentando que um dos problemas em todas as investigações forenses é fornecer informações claras sobre como as evidências foram coletadas, preservadas, analisadas, apresentadas, além de comprovar que o material não foi alterado durante as etapas de investigação. Este trabalho discute brevemente que vários estudiosos defendem a ideia de uma cadeia de custódia (CoC) e sugerem a utilização de *blockchain* para garantir a segurança dos dados trafegados pela cadeia. Embora cite o *blockchain*, o artigo não comenta como integrá-lo a dispositivos IoT para atuar em uma cadeia de produção documental, faltando também uma análise das questões relacionadas à segurança cibernética, dois temas que são abordados neste trabalho.

## 2.6.2 Industry 4.0

Geralmente, um processo industrial automatizado precisa minimizar atrasos, principalmente aqueles relacionados às rotinas de segurança. Em Aazam *et al.* [67], os desafios relacionados à IoT Industrial são discutidos, e a computação em névoa é usada localmente com suporte de middleware entre o ambiente industrial e os serviços web de controle.

Os desafios da Indústria 4.0 e da IIoT também são abordados por Matthyssens (2019) [68], apresentando uma breve pesquisa de tópicos discutidos nos últimos anos e destacando ligações entre tecnologias emergentes no domínio. Neste trabalho também são abordados os recursos para inovação de valor na Indústria 4.0

A proposta de Boyes *et al.* [69] consiste de uma nova estrutura que associa conceitos de IIoT e sistemas ciberfísicos de acordo com o paradigma da Indústria 4.0. O artigo citado propõe uma taxonomia da IoT e enumera algumas características de dispositivos IIoT. Acerca da segurança, o artigo analisou algumas vulnerabilidades e identificou brechas na literatura sobre o assunto.

No artigo de Garrido-Hidalgo *et al.* [70], um estudo de caso é feito utilizando dispositivos Iot e sensores com foco em comunicação *wireless* a fim de avaliar seu desempenho. O artigo citado discute tendências gerais a respeito de automação e troca de dados em tecnologias de fabricação. Alguns conceitos de Indústria 4.0 e IIoT são usados no trabalho, e os resultados mostram entraves que precisam ser discutidos para melhorar a IIoT. O artigo apresenta um panorama dos protocolos de comunicação para Indústria 4.0 e a noção da fábrica inteligente de estrutura modular, onde

sistemas ciberfísicos monitoram os procedimentos físicos, criam uma cópia virtual do mundo físico e tomam decisões descentralizadas. Usando a IoT, os sistemas ciberfísicos comunicam e cooperam até com humanos em tempo real, e a noção de Internet dos Serviços (IoS) possibilita que serviços internos e interorganizacionais sejam oferecidos e utilizados pelos participantes da cadeia de valor.

O presente trabalho também emprega computação em névoa e de borda, além de microsserviços de software para integrar uma cadeia de produção documental modular, considerando a preservação da autenticidade e integridade dos documentos trafegados. Diferentemente dos artigos citados, este trabalho introduz um conjunto de serviços de segurança operando especificamente para monitorar a CoC da linha de produção documental, considerando um abrangente conjunto de vulnerabilidades que foram previstas em um modelo adversarial. Nossa validação experimental utiliza um protótipo completo desenvolvido, o qual permite a efetiva verificação dos serviços de segurança baseados em IoT propostos.

### **2.6.3 Segurança em IoT**

Os requisitos de segurança em IIoT são apresentados no trabalho de Hansch *et al.* (2019) [71], definindo os requisitos técnicos, funcionais e processuais para a segurança nesse domínio. Os requisitos técnicos se referem às propriedades técnicas dos dispositivos e ao sistema IIoT em geral, o que implica a utilização de criptografia simétrica e assimétrica, assinatura digital, bem como proteções relacionadas à negação de serviços (DoS). Os requisitos funcionais estão relacionados ao acesso externo ao sistema, exigindo gestão de usuário com credenciais, autorização e autenticação. Os requisitos de processo incluem o processo de desenvolvimento e o ciclo de vida do produto, estando relacionados às medidas de segurança, como testes de penetração, segmentação de rede e restrições de acesso.

Similarmente, o presente trabalho introduz um conjunto de serviços de segurança operando especificamente para monitorar a DCoC da linha de produção documental, considerando um abrangente conjunto de vulnerabilidades que foram previstas em um modelo adversarial. Nossa proposta segue o princípio de medidas de segurança para IoT completamente distribuídas [46] e é desenvolvida posteriormente como um protótipo para verificar a possibilidade efetiva de aplicação em um linha de produção IIoT.

### 3 ESTUDO DO PROBLEMA E HIPÓTESES DA TESE

O documento arquivístico possui elementos diplomáticos que devem ser considerados para comprovar sua autenticidade e veracidade. Esses elementos incluem assinaturas humanas, identificação da instituição, marcas d'água, número do protocolo e número do identificador. O aumento da produção de registros documentais arquivísticos, quase exclusivamente em meio digital, impõe desafios a diversas áreas do conhecimento, principalmente no que se refere à gestão, à preservação e ao acesso a tais documentos em diferentes etapas do ciclo de vida. As funcionalidades relacionadas geralmente devem ser fornecidas usando diferentes softwares interoperáveis.

Desde a década de 1990, a comunidade internacional tem desenvolvido iniciativas para orientar a modelagem e implementação de repositórios digitais, apontando os requisitos para garantir confiabilidade a esses repositórios. Implementar um repositório digital confiável é essencial para garantir a preservação, o acesso e a autenticidade de longo prazo dos materiais digitais.

Além dos ambientes de preservação gerenciados (RDC), o modelo OAIS recomenda o uso de outro ambiente associado, destinado exclusivamente à disseminação e ao acesso a documentos digitais e informações relacionadas, preservados como seguros e autênticos. Esse ambiente, ou plataforma de acesso, tem sido cada vez mais discutido e adotado mundialmente para publicação de arquivos digitais na internet.

Espera-se que esse sistema de produção forneça, desde a origem do documento até seus destinos, uma CoC previsível e controlada do produtor ao consumidor. A estruturação de ambientes de gestão, preservação e acesso visa apoiar a criação, manutenção e divulgação de documentos autênticos e completos ao longo do seu ciclo de vida, proporcionando, assim, um acesso fácil e rápido à documentação. Além disso, essas estruturas articuladas contribuem e promovem o acesso à informação de interesse geral (transparência ativa), uma vez que preservam e fornecem informações demandadas por funcionários e gestores de organizações públicas e privadas e pelo cidadão comum.

Portanto, o objeto de pesquisa deste trabalho é garantir a autenticidade e as propriedades de integridade da CoC documental, dadas as vulnerabilidades e os problemas de rastreabilidade existentes desde a origem do documento até seu armazenamento seguro em um repositório RDC, a partir do qual o documento pode ser acessado abertamente. Dessa forma, o presente trabalho propõe o IoTSec2DCoC, uma estrutura para melhoria da CoC documental, utilizando medidas de segurança abrangentes, operando com suporte de IoT e estruturado de acordo com a arquitetura de software de microsserviços.



### 3.1 CADEIA DE CUSTÓDIA DOCUMENTAL CONVENCIONAL

Deve ser dada precisão ao conceito de uma CoC documental convencional, uma vez que constitui a base para comparar a estrutura da nossa proposta e seus resultados como um protótipo de prova de conceito validado.

A CoC documental, que denotamos como convencional e serviu como nossa base de referência, é a do modelo OAIS, ilustrado na Figura 3.1. Embora esse modelo conceitualmente forneça garantias de preservação de documentos, a verificação prática dessas garantias ficará para estudos posteriores.

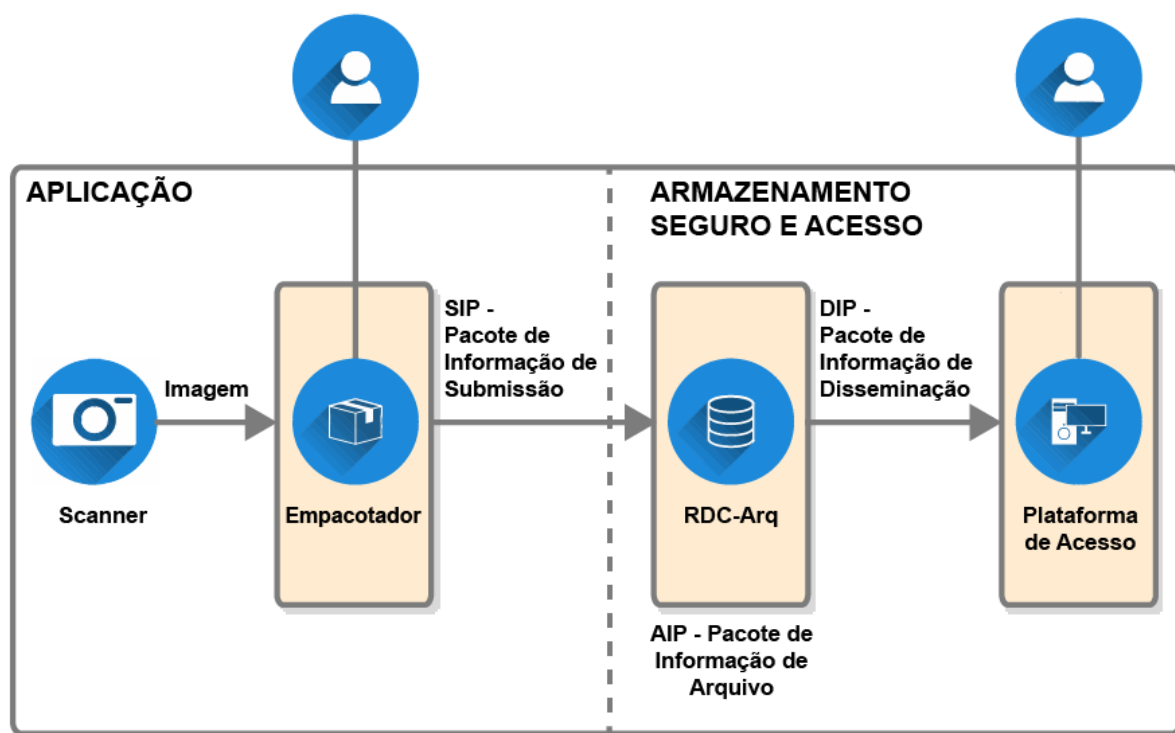


Figura 3.1: Arquitetura de uma solução de cadeia de custódia convencional. Fonte: Adaptado de: [23].

De acordo com o modelo OAIS, a CoC deve ser ininterrupta e preservada durante todo o processo, desde a admissão do documento digital pelo sistema, seja ele uma cópia digital convertida (escaneada) de uma contraparte analógica (documento impresso) ou um documento nato digital. Segundo Force (2002) [19], existem ameaças a essa cadeia para os documentos armazenados e para os documentos que são processados e transferidos através do espaço ou do tempo, seja entre pessoas ou aplicativos. Especificamente, essas ameaças estão relacionadas a vulnerabilidades em várias situações, como as seguintes:

- A infraestrutura de tecnologia da informação que dá suporte ao armazenamento é suscetível a vários eventos contínuos de operação incorreta, inclusive na manutenção preventiva e evolutiva, afetando formatos de dados e informações referenciais de banco de dados.

- Os fluxos na Figura 3.1 apresentam, como primeira fonte, dispositivos que utilizam software proprietário de terceiros, atuando durante o processo de captura de imagens, nos quais não se deve confiar cegamente, uma vez que existem riscos de esses dispositivos e softwares serem adulterados para gerar documentos comprometidos antes que entrem no fluxo para o repositório seguro.
- Durante a transferência de pacotes, existem riscos consideráveis de integridade, pois esses pacotes não têm nenhuma proteção de criptografia e podem ser facilmente capturados, modificados ou substituídos.
- O modelo carece de rastreabilidade efetiva, pois não há monitoramento nem logs fora do repositório seguro. Portanto, em caso de incidente de segurança, é difícil identificar o dano e repará-lo, dificultando também a análise forense.
- No que se refere à disponibilidade, como há forte acoplamento de serviços nas principais aplicações, isso pode ser um problema caso alguma dessas aplicações fique indisponível, prejudicando o funcionamento regular. Não está claro se é possível balancear o modelo para proteger e otimizar a execução, e nenhuma contingência é observada em caso de latência de comunicação ou perda de informações no repositório.

Vale lembrar que essas observações consideram que os componentes da Figura 3.1 seguem os padrões do modelo OAIS. Como essa conformidade não ocorre necessariamente na infraestrutura de tecnologia da informação real, há um problema crítico em confiar cegamente nessas estruturas, ou seja, sem um sistema independente de monitoramento e detecção de mau comportamento.

Portanto, a CoC documental convencional deve ser complementada com um conjunto abrangente de medidas de segurança para proteger e validar todos os seus fluxos de informação, garantindo, assim, a integridade e os metadados do documento digital. Essa necessidade é o que motiva nossa proposta IoTSec2DCoC apresentada a seguir.

## 3.2 HIPÓTESES

Em função das vulnerabilidades e deficiências do modelo OAIS e a necessidade de evolução para modelos industriais inteligentes, sustentáveis e seguros, como foi visto anteriormente, organizamos o trabalho para testar uma série de hipóteses:

- Hipótese 1 - as assinaturas digitais e a criptografia simétrica permitirão criar meios de autenticação e integridade aplicáveis em todos os componentes da cadeia de custódia e verificáveis no fluxo dos documentos de um componente para o outro.
- Hipótese 2 - quanto à cadeia de custódia digital de documentos, a IIoT vai permitir uma solução em tempo real de monitoração a eventos de segurança.

- Hipótese 3 - a IIoT tem a plasticidade e a flexibilidade para diversas configurações adequadas a cada ambiente específico de produção documental.
- Hipótese 4 - a IIoT, os protocolos e medidas de segurança podem ser meios com pequena pegada ambiental, ou seja, aplicações de técnicas sustentáveis que contribuam para os ambientes urbanos, apresentando, ainda, o desempenho industrial adequado.
- Hipótese 5 - a estrutura de IIoT permite incluir medidas de segurança totalmente distribuídas para monitorar e reagir a eventos de segurança e ataques.
- Hipótese 6 - um dispositivo de IoT pode ser usado especificamente para a interface com a linha de digitalização de documentos analógicos.

## 4 PROPOSTA DA ARQUITETURA IOTSEC2DCOC PARA GARANTIR A CADEIA DE CUSTÓDIA DOCUMENTAL

Este capítulo apresenta a estrutura geral e o funcionamento da proposta IoTSec2DCoC, enquanto os detalhes descritivos sobre os componentes e as operações são tratados nas respectivas subseções.

### 4.1 VISÃO GERAL DA ARQUITETURA PROPOSTA

Como o IoTSec2DCoC é projetado como um sistema ciberfísico para atuar como controlador da linha de produção documental, é composto por dispositivos e componentes de software que operam juntos sob a coordenação de um Middleware IoT distribuído, isto é, um conjunto de software distribuído em uma instância de rede IoT. Nessa estrutura, os componentes mostrados na Figura 3.1 da CoC documental convencional, ou seja, o dispositivo de scanner e os componentes de software do Empacotador, o RDC-Arq e a Plataforma de Acesso, ficam rodeados pelos elementos de controle que darão suporte à CoC com um conjunto totalmente distribuído de medidas de segurança [46], conforme o esquema mostrado na Figura 4.1.

Nesse esquema proposto, um Dispositivo Inteligente de ingestão de documentos é introduzido para cuidar dos documentos vindos do scanner que são verificados para fins de segurança, permitindo a coleta de metadados e a execução de operações de pré-processamento, como reconhecimento óptico de caracteres, classificação de documentos baseada em aprendizado de máquina e possível futuro processamento avançado de linguagem natural, tratamento de imagem e reconhecimento de padrões. Após a ingestão, o documento é transferido pelos componentes IoT (dispositivos virtuais, Gateway de API, Middleware Edge, Middleware Cloud) para o sistema Empacotador e, posteriormente, para o RDC-Arq e a Plataforma de Acesso. Cada transferência é supervisionada por medidas de segurança implantadas em pares de pontos de saída e entrada, conectando cada componente ao próximo, na forma de dispositivos IoT virtuais verificadores, assim chamados por se tratarem de módulos de software embutidos que se comportam como dispositivos IoT de hardware-software integrados. Tais dispositivos verificadores agem como sensores capazes de detectar autenticação de documentos e falhas de integridade, de forma análoga a outros sensores da indústria.

Esse posicionamento permite que esses agentes de segurança verifiquem a autenticidade de origem e a integridade do documento em toda a linha de produção. Como os atores de segurança dialogam com controladores gerais baseados em nuvem, para fins de configuração de serviços de segurança (registro de serviço, assinatura digital, certificados), a estrutura também permite informações de segurança e monitoramento de eventos, incluindo detecção e mitigação de ataques,

também mostrado na Figura 4.1.

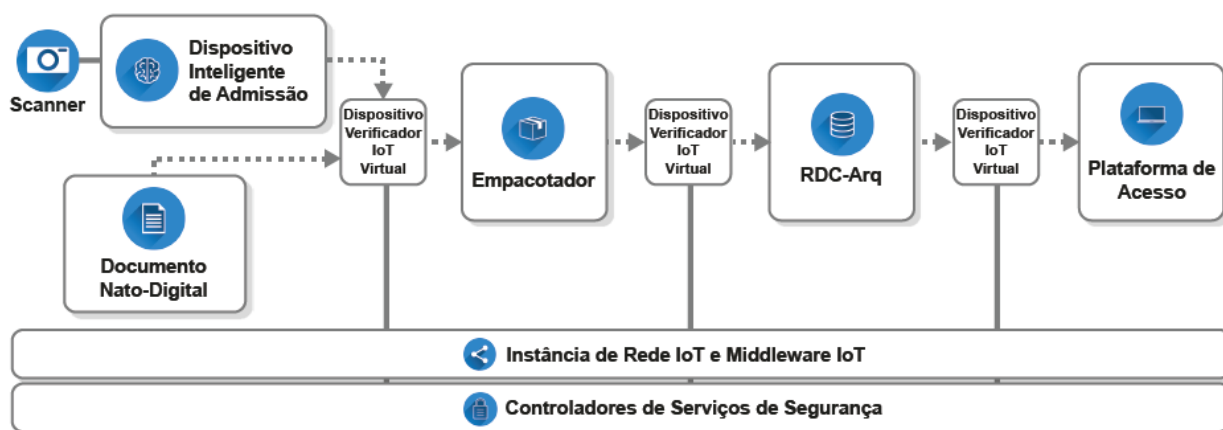


Figura 4.1: Visão geral da estrutura IoTSec2DCoC.

A arquitetura proposta pretende, além disso, apresentar a plasticidade e flexibilidade para se adaptar às variações nas linhas de produção de documentos, semelhante às plantas industriais. Dadas as diferentes configurações possíveis de Middleware IoT, o suporte da instância de rede IoT permite acomodar um número variável de scanners que depende do número de documentos a serem processados e do desempenho desejado. A mesma plasticidade se aplica ao caso de fontes de documentos nato digitais que apresentam tipos heterogêneos e cujo número é variável.

A instanciação de Middlewares Edge traz a flexibilidade de adaptar dinamicamente à planta de produção de documentos e, como essa arquitetura baseada em IoT também pode ser instanciada em outros locais, permite a implementação de várias plantas de produção de documentos.

No que diz respeito à integração e supervisão do Empacotador, do RDC-Arq e da plataforma de acesso, do ponto de vista da estrutura da IoT, esses sistemas são considerados como aplicações IoT cujo comportamento, em relação à autenticidade de origem e à integridade dos documentos, deve ser supervisionado por meio da intermediação de módulos de software embarcados que constituem dispositivos virtuais de IoT verificadores.

Essa estrutura moldável é, então, concebida para proporcionar medidas totalmente distribuídas que estejam em conformidade com as práticas internacionais e as normas legais que apoiam a segurança da informação e a CoC documental.

Outra característica essencial da IoTSec2DCoC é que os componentes de software de interação são especificados de acordo com a arquitetura de microsserviços [72], permitindo o desenvolvimento de software seguro e operações responsáveis pelo mapeamento dos módulos codificados para a infraestrutura de implantação dos vários microsserviços que compõem cada aplicativo.

## 4.2 ESTRUTURA DA ARQUITETURA DE SOFTWARE IOTSEC2DCOC

De acordo com a arquitetura de software mostrada na Figura 4.2, os módulos IoTSec2DCoC, estruturados como conjuntos de microsserviços, são definidos em termos das funções que devem desempenhar, seguindo o princípio de modularização, o qual recomenda agrupar pequenos módulos e componentes para resolver um problema específico. Uma consideração importante é que a codificação de microsserviços em diferentes linguagens de programação permite aplicar o melhor oferecido por cada um em relação aos objetivos de cada microsserviço.

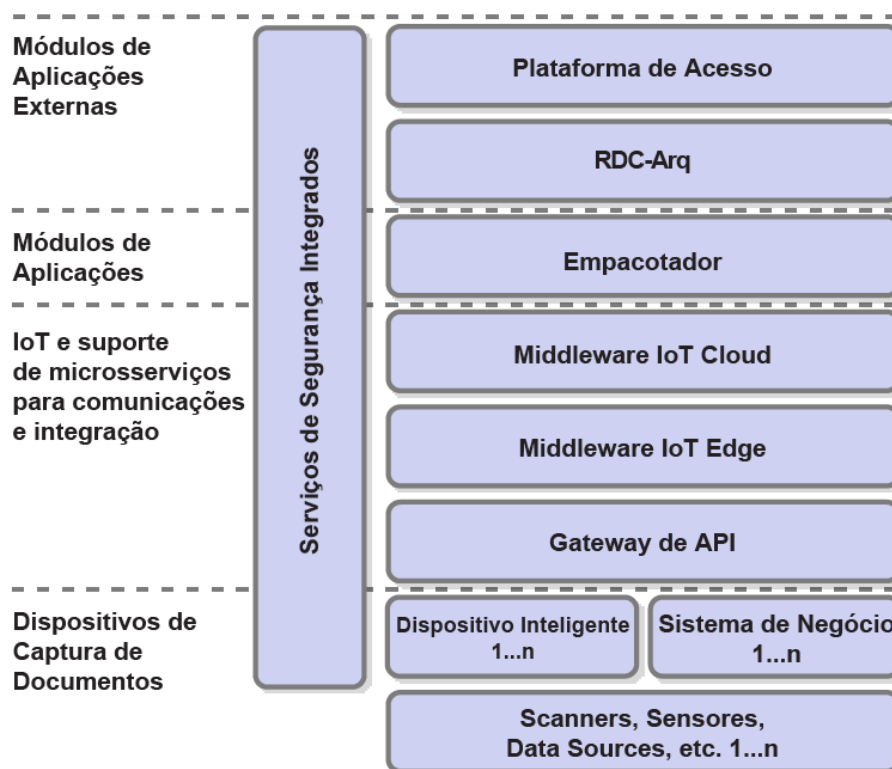


Figura 4.2: Arquitetura de Software IoTSec2DCoC.

Como todas as comunicações entre os vários microsserviços são feitas por meio da API REST [73], isso também contribui para a plasticidade e a flexibilidade, pois o REST é um protocolo comprovado que simplifica a implantação dos microsserviços. Além disso traz escalabilidade ao *framework*, bem como possibilidades de evolução, já que, se necessário, é possível escrever novos microsserviços que agregam funcionalidades novas ou intermediárias.

Por uma questão de validação da proposta, a escolha de uma arquitetura de microsserviço contribuiu para o desenvolvimento do protótipo de prova de conceito (PoC) IoTSec2DCoC, cuja arquitetura de microsserviço é desenvolvida e implementada dentro do *framework* Spring Boot [74], com um conjunto de ferramentas do Spring Cloud junto a uma série de componentes programados com a linguagem de programação Java. Outros componentes do protótipo PoC são codificados na linguagem Python usando o Framework Flask [75]. O funcionamento integrado desses módulos

é descrito na subseção 5.3.

Considerando esses princípios de microsserviços e a escolha da IoT como meio ciberfísico de controle da CoC documental, a arquitetura de software da IoTSec2DCoC compreende os componentes especificados na Tabela 4.1.

Tabela 4.1: Componentes da IoTSec2DCoC.

<b>Componentes</b>	<b>Descrição</b>
Dispositivo Inteligente	Conjunto de aplicativos implantados no hardware de um dispositivo, como o Raspberry Pi, responsável por interagir com o mundo físico e transmitir os documentos digitalizados para as camadas superiores.
Dispositivo Virtual	Componentes de software embutidos que se comportam como dispositivos IoT de hardware-software integrados, como o dispositivo verificador que supervisiona as transferências de documentos entre outros módulos na CoC documental.
Gateway de API	Componente que é um ponto único de comunicação para serviços em nuvem, balanceamento de carga (Round-robin) e tolerância a falhas.
Serviços de Segurança Integrados	O conjunto abrangente de serviços responsável por apoiar medidas de segurança e integração para toda a solução, envolvendo autenticação, assinatura digital, gestão de certificados digitais (autoridade certificadora), centralização de logs (serviço de monitoramento) e registo de serviços.
Middleware Edge	Instância hierárquica de limite inferior do Middleware IoT, facilitando as comunicações entre as outras camadas e fornecendo armazenamento de dados para processamento futuro. As instâncias de Middleware Edge são aquelas que se integram diretamente com dispositivos inteligentes.
Middleware Cloud	Instância hierárquica de limite superior do Middleware IoT, com a mesma funcionalidade da instância de borda (Middleware Edge), embora seja responsável por integrar e agregar dados das instâncias inferiores.
Sistema Empacotador	Componente para consultar os novos documentos e seus metadados do Middleware IoT da nuvem e preparar um pacote para enviar o SIP ao RDC-Arq.
RDC-Arq	Componente para ingestão dos pacotes SIP, gerando os pacotes AIP e DIP, sendo o primeiro para preservação em repositório seguro, e o último para disponibilidade de documentos na Plataforma de Acesso.
Plataforma de Acesso	Componente para consumir os pacotes DIP e para publicar os metadados e documentos digitais contidos nos pacotes para fornecer acesso aos usuários.

### 4.3 DETALHAMENTO DOS PROCESSOS E FLUXOS DO MODELO

Na codificação conforme a Figura 4.1, observa-se que os microsserviços são consistentes quanto às funções as quais eles devem executar, ou seja, módulos e componentes pequenos que resolvem um problema específico. A codificação dos microsserviços em diferentes linguagens de programação possibilita aplicar o que de melhor cada linguagem oferece quanto aos seus objetivos. Toda a comunicação entre os diversos microsserviços é feita através de API REST.

Como vantagens, têm-se a simplificação de *deploy* dos microsserviços e a fácil escalabilidade do sistema. Então, caso necessário, é possível escrever novos microsserviços que adicionam funcionalidades novas ou intermediárias.

A Figura 4.3 apresenta a solução IoTSec2DCoC, detalhando seus fluxos de processos - cada passo está referenciado no diagrama com uma numeração dentro de círculos pontilhados em vermelho -, os quais realizam as seguintes tarefas:

- **Passo 1:** o usuário humano se autentica com seu *token* criptográfico no Dispositivo Inteligente, contendo seu certificado digital. Se o usuário for validado, o scanner é acionado pelo Dispositivo Inteligente para realizar a digitalização dos documentos via uma API de integração com o scanner.
- **Passo 2:** o Dispositivo Inteligente se autentica no Serviço de Autenticação com as suas credenciais existentes no Raspberry Pi, adquirindo um *token* JWT [76], o qual, posteriormente, será utilizado na comunicação com o Middleware Edge. O componente de OCR, dentro do Dispositivo Inteligente, realiza os seguintes processos:
  - Recupera imagens digitalizadas.
  - Realiza o processamento de OCR.
  - Gera o TIFF multi-páginas [77] do documento e um PDF com OCR embutido.
  - Envia esses arquivos juntamente com os metadados técnicos relativos ao documento.

Ainda dentro do Dispositivo Inteligente, o componente interno de classificação e extração de metadados executa os seguintes processos:

- Realiza a classificação do tipo de documento a partir do texto extraído do OCR utilizando técnicas de *deep learning*.
  - Executa extração de metadados essenciais sobre o documento via técnicas de *deep learning* com visão computacional.
  - O componente cria um arquivo compactado (formato ZIP) com as imagens em formato TIFF e o documento em PDF juntamente com todos os metadados.
  - O Dispositivo Inteligente acessa o serviço de check-out. Esse componente recebe o arquivo ZIP, assina-o digitalmente e o repassa para o Middleware Edge.
- **Passo 3:** o item 4 trata dos documentos arquivísticos analógicos que foram digitalizados, mas a solução também prevê a possibilidade futura de suporte à cadeia de custódia de documentos nato digitais. Nesses casos, os sistemas de negócio poderiam consumir um Serviço de Admissão de Arquivos Nato Digitais, informando o documento e seus metadados em um arquivo compactado, assinado digitalmente pelo componente de checkout, o qual também deverá ser utilizado pelos sistemas de negócio. Esse serviço está previsto no modelo, mas ainda não foi implementado.



- **Passo 4:** O Middleware Edge recebe os documentos dos dispositivos inteligentes e do Serviço de Admissão de Arquivos Nato Digitais, caso esse exista, e executa os seguintes processos:
  - A partir do componente de Check-in, é validada a assinatura digital do arquivo compactado recebido no Serviço de Validação de Assinaturas Digitais.
  - Em caso de sucesso na validação da assinatura digital, os dados da assinatura são enviados para o Middleware Cloud.
- **Passo 5:** O Middleware Cloud recebe o arquivo, realiza a validação via seu componente de Check-in, o processa e o armazena. Em uma certa periodicidade, o middleware envia um conjunto de arquivos armazenados para o Empacotador (assinados novamente). Então, o Empacotador realiza os seguintes processos:
  - Realiza a solicitação de validação da assinatura digital para cada um.
  - Extrai os dados do arquivo comprimido e os processa de acordo com suas regras, criando um novo registro de documento para serem realizados os passos de revisão humana da digitalização, a depender da criticidade e norma vigente para a amostra.
- **Passo 6:** o Empacotador, então, realiza os seguintes passos para envio do documento à próxima camada da solução:
  - Gera o SIP com todos os dados do documento, incluindo os que podem ter sido alterados com a revisão humana, e realiza uma nova assinatura digital a partir do componente de Check-out.
  - Envia o SIP assinado para o RDC-Arq.
  - O componente de Check-in no RDC-Arq recebe o arquivo de assinatura e solicita a sua validação.
  - Se a assinatura digital for validada, o SIP é extraído do arquivo de assinatura digital e enviado ao RDC-Arq (Repositório Arquivístico Digital Confiável), denominado Archivematica, o qual será detalhado na subseção 5.5.2.
- **Passo 7:** o RDC-Arq realiza todos os procedimentos necessários de um Repositório Arquivístico Digital Confiável, guarda o AIP criptografado utilizando uma chave simétrica GPG e, então, envia um pacote DIP para a Plataforma de Acesso da solução (Atom), novamente realizando a assinatura digital e sua validação.
- **Etapa 8:** o usuário final pode acessar os documentos e metadados na interface da plataforma de acesso. Se o usuário deseja fazer o download de todos os dados do documento, o componente Check-out pode gerar uma nova assinatura a ser baixada com o conteúdo. A posteriori, o usuário pode validar a assinatura digital anexada ao conteúdo baixado através da interface web do serviço de validação de assinatura digital para verificar a integridade e autenticidade dos dados baixados.

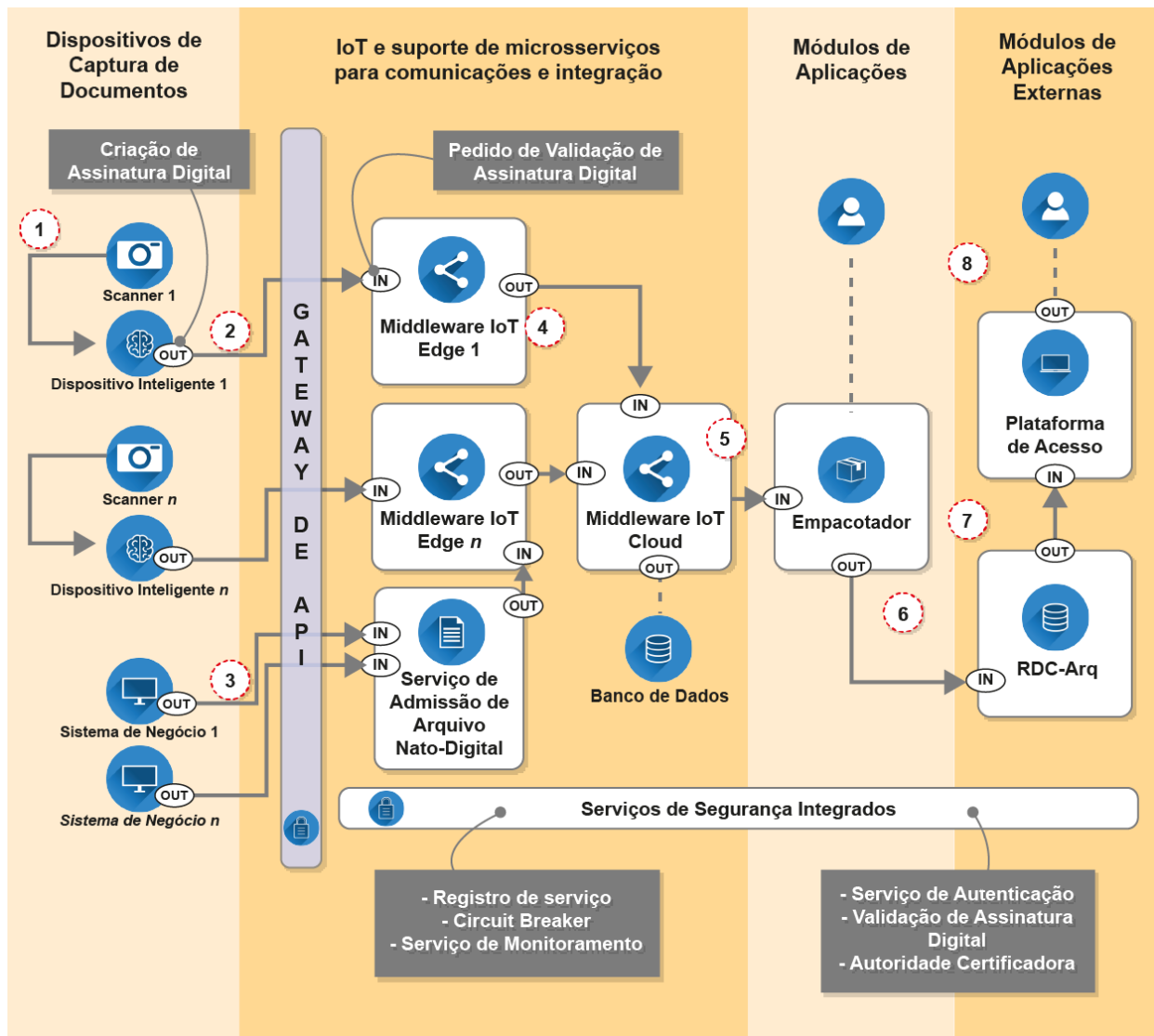


Figura 4.3: Arquitetura da solução IoTSec2DCoC e seus fluxos.

Em paralelo, são executadas outras tarefas constantes durante todo o fluxo de dados da solução. Por exemplo, o monitoramento do status dos serviços pelo Circuit Breaker [72], além do envio de todos os registros de log de todos os módulos do sistema para Serviço de Monitoramento, para futura rastreabilidade e auditoria. Como esse serviço foi implementado na PoC, será descrito na subseção 5.4.3.

Também em um momento inicial, todos os componentes da solução que fornecem serviços HTTP e/ou realizam assinatura digital solicitam um certificado digital válido ao Serviço de Autoridade Certificadora, caso ele esteja expirado ou não possuam um. A descrição da implementação na PoC pode ser vista na subseção 5.4.1.

Na Figura 4.4, são destacados todos os componentes que fazem parte da rede IoT. Da esquerda para a direita, temos os seguintes componentes: o Dispositivo Inteligente na camada de dispositivos; a instância do Middleware Edge e sua instância superior o Middleware Cloud com seu banco de dados e todos os serviços de segurança integrados, todos na camada de IoT e suporte

de microserviços para comunicações e integração; e, por fim, o Empacotador.

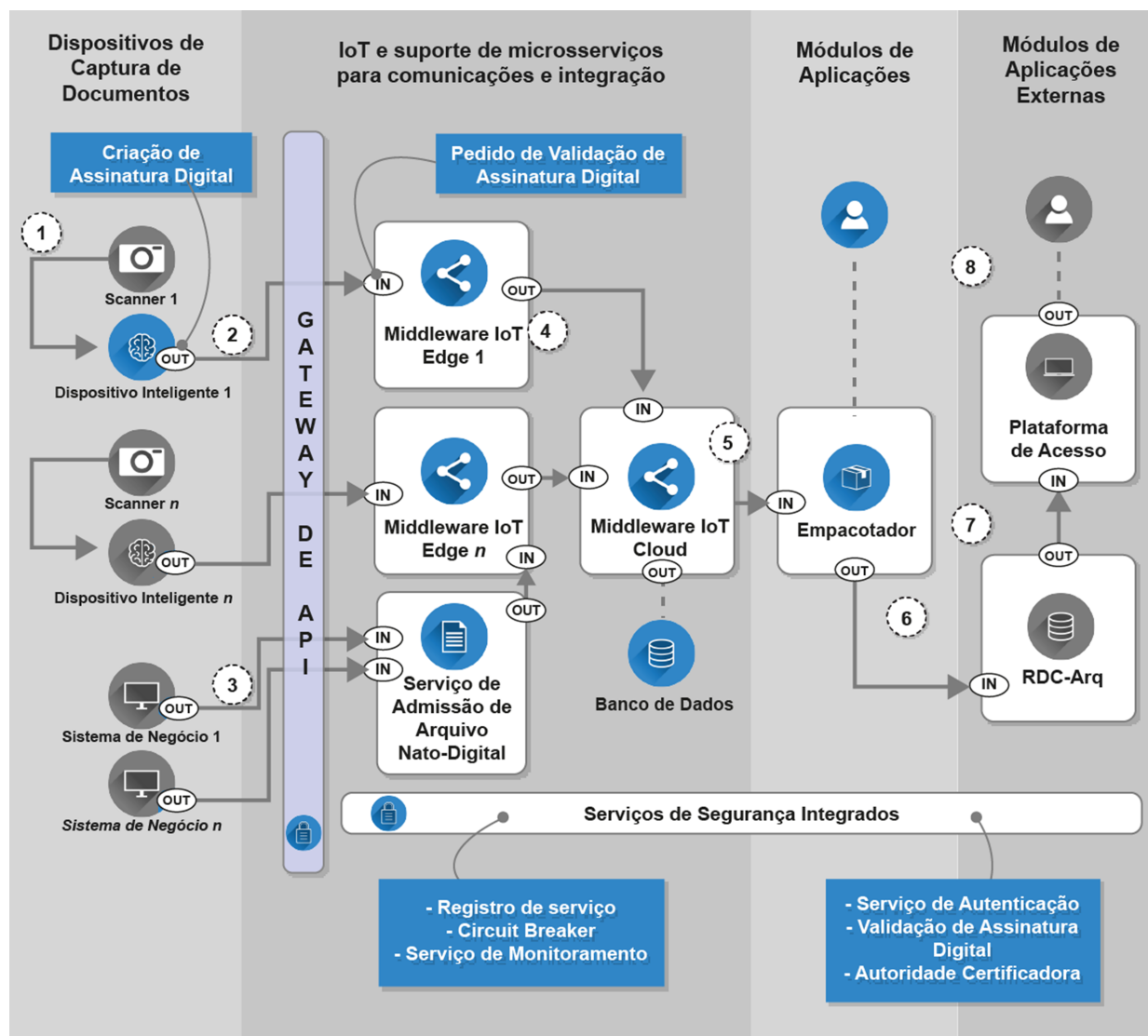


Figura 4.4: Identificação de Dispositivo Inteligente e outros componentes IoT.

Todos esses componentes em destaque na Figura 4.4 estão interligados através da IoT, pois foram devidamente autenticados e tiveram seus serviços registrados.

Não podemos deixar de observar o diagrama da arquitetura na Figura 4.5, no qual os componentes da cadeia de custódia convencional estão presentes no novo modelo proposto, porém, estão securitizados por diversas medidas e contramedidas de segurança. Da esquerda para a direita, temos: o scanner, agora gerenciado pelo Dispositivo Inteligente; o Empacotador, agora serviço da IoT; o RDC-Arq e a Plataforma de Acesso. Mesmo sendo aplicações externas, utilizam os dispositivos IoT virtuais Check-in e Check-out que estão conectados à IoT e garantem que somente sistemas de negócio, componentes e outros, previamente cadastrados e validados se conectem e trafeguem dados pela IoT.

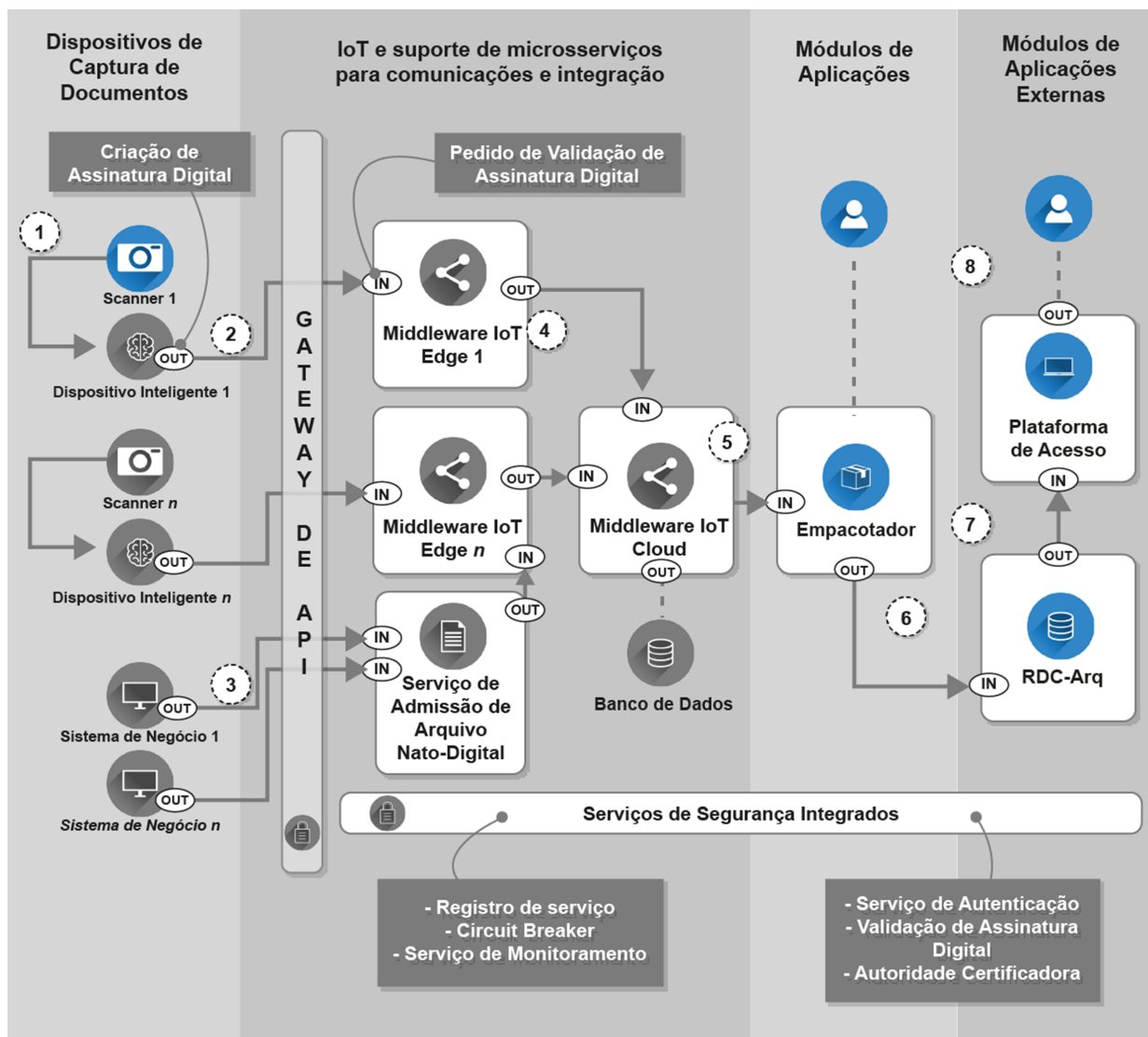


Figura 4.5: Identificação dos componentes da cadeia de custódia convencional.

Na Figura 4.6, ficam destacados os componentes de Check-out e Check-in da solução. Eles estão presentes em todo o fluxo principal da cadeia de custódia, pois fornecem os mecanismos de garantia de autenticidade e integridade dos documentos trafegados. Na saída de todo componente da solução responsável pelo envio de documentos arquivísticos, deve existir o componente de Check-out. Este, a partir do seu certificado digital e chave privada configurados, realiza a assinatura digital do documento (TIFF, PDF e metadados) e, então, envia para o próximo componente do fluxo o documento, a assinatura digital e o certificado digital do assinante. Quem recebe os dados enviados pelo Check-out no próximo componente do fluxo é sempre o seu Check-in, o qual solicita a validação da assinatura digital do documento enviado e também a validação da confiança do certificado digital, garantindo a integridade e autenticidade documental durante toda a cadeia de custódia.

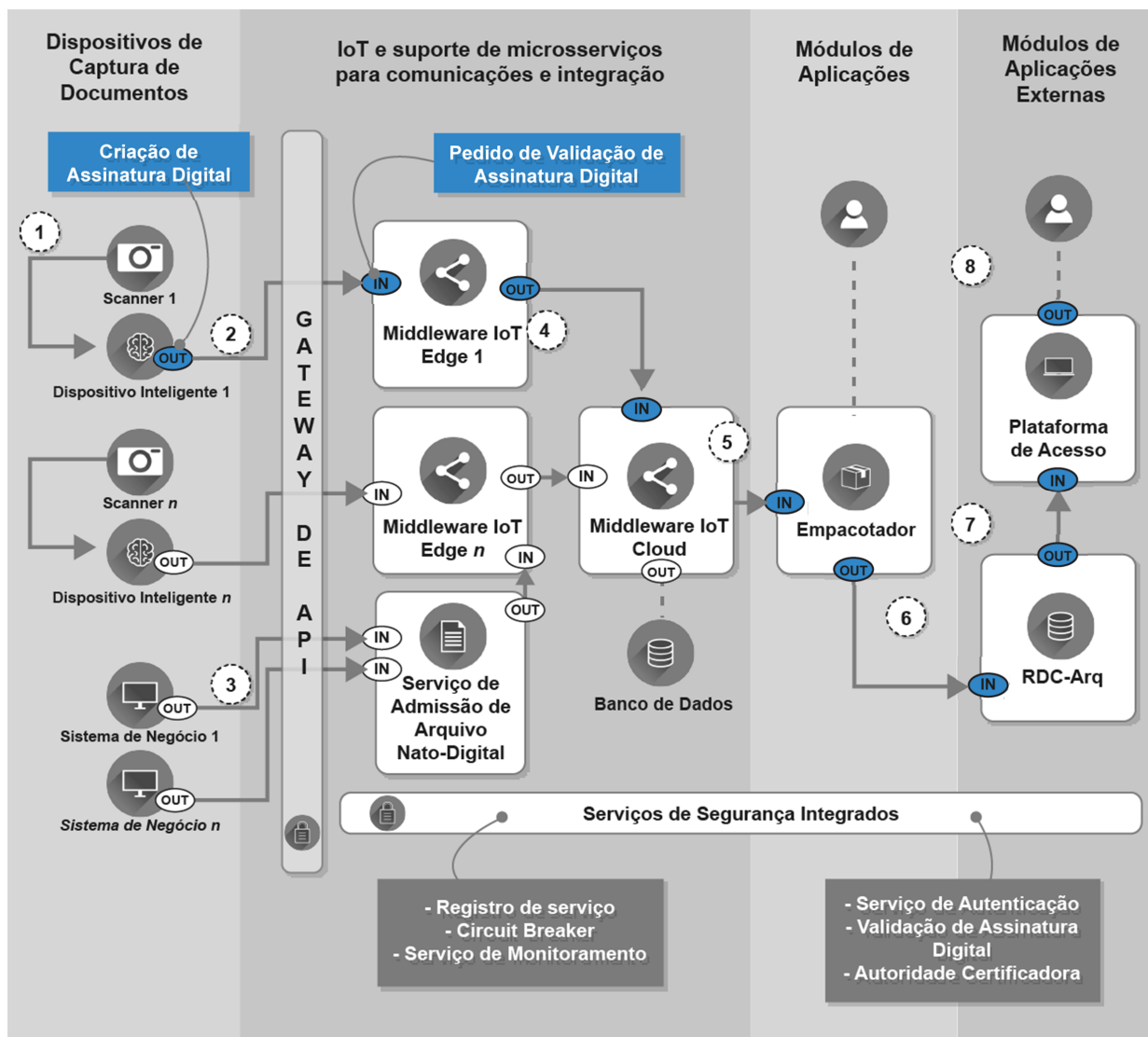


Figura 4.6: Identificação dos componentes Check-in e Check-out.

Na Figura 4.3, estão destacados todos os serviços de segurança integrados, os quais atuam paralelamente em todo o fluxo, fornecendo microsserviços que atuam no facilitamento de integração e na securitização da solução. Segue a descrição resumida de cada um (estes serão melhor descritos nas seções subsequentes):

- **Registro de Serviço:** é o responsável por centralizar o registro de todos os componentes da solução, fornecendo nomes DNS que facilitam a comunicação entre os componentes.
- **Circuit Breaker:** auxilia no controle da perturbação do funcionamento normal da aplicação, fazendo com que os serviços sejam tolerantes a falhas, diminuindo o impacto de indisponibilidades de um serviço nos outros paralelos.
- **Serviço de Monitoramento:** oferece mecanismos de rastreabilidade e auditoria para toda a solução.

- **Serviço de Autenticação:** realiza a autenticação de elementos externos à rede IoT, os Dispositivos Inteligentes e os Sistemas de Negócio.
- **Serviço de Validação de Assinaturas Digitais:** é o serviço consumido pelos componentes de Check-in para validar assinaturas e a confiança de certificados digitais.
- **Serviço de Autoridade Certificadora:** oferece mecanismos de gestão de todo o ciclo de certificados digitais (emissão, manutenção e revogação) para todos os componentes da solução que necessitem.

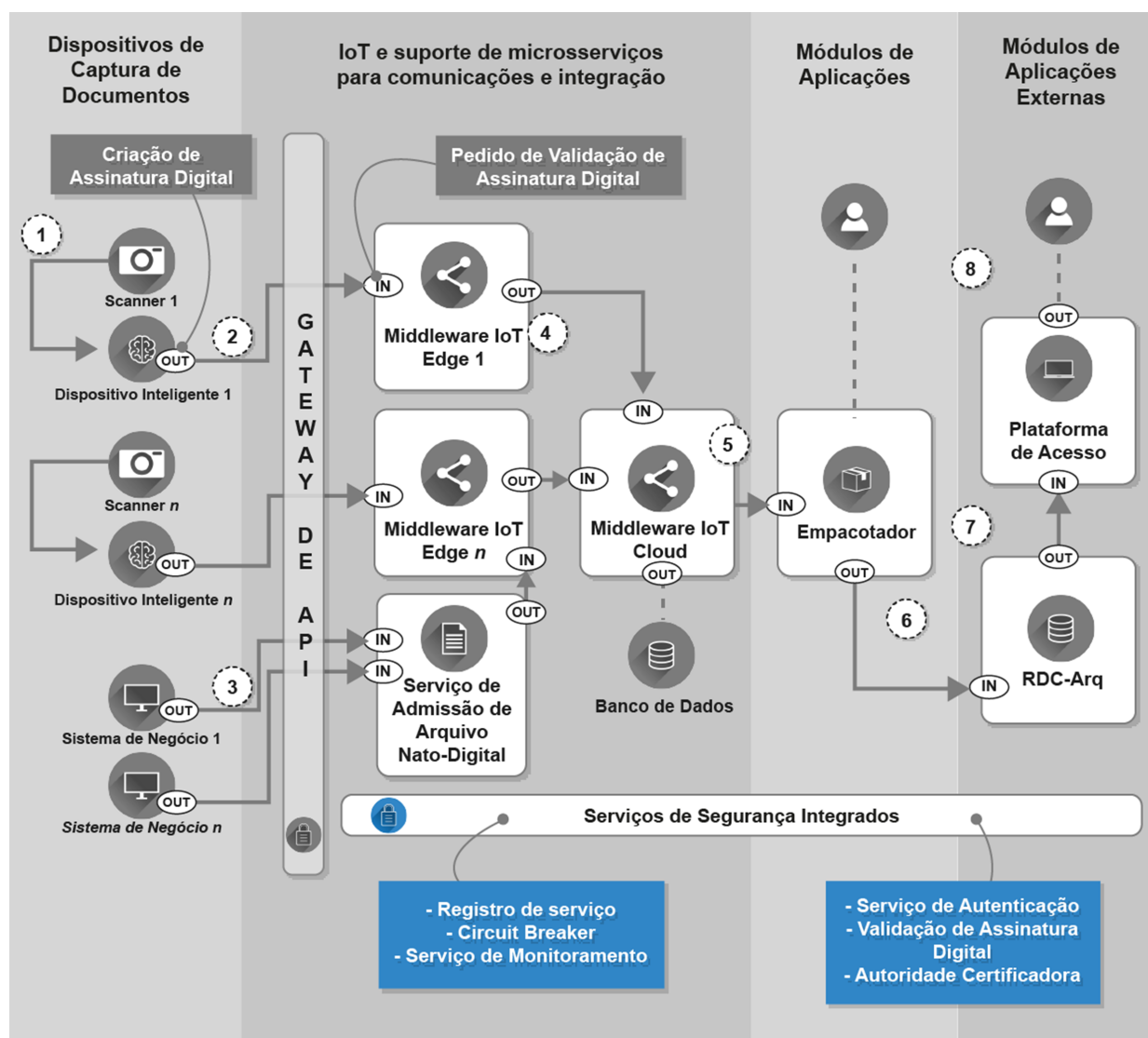


Figura 4.7: Identificação dos serviços de segurança.

## 5 PROVA DE CONCEITO

Para a validação da proposta, um protótipo PoC funcional foi desenvolvido e testado em ambiente de laboratório controlado. A execução dos testes foi organizada para permitir a observação sistemática da estrutura e funcionalidades propostas do IoTSec2DCoC, bem como para obter registros do funcionamento e da interoperação das instâncias da arquitetura.

Os requisitos do protótipo desenvolvido incluíram:

- A necessidade de verificar a autenticação e autorização dos dispositivos para coletar documentos, incluindo o scanner, o Dispositivo Inteligente e artefatos de Middleware IoT.
- A operação correta do scanner e do Dispositivo Inteligente em relação à qualidade do documento digitalizado, integridade de metadados e precisão na classificação de documentos por um módulo de aprendizado de máquina.
- A verificação da plasticidade e flexibilidade de implantação do IoTSec2DCoC por meio dos níveis do Middleware IoT distribuído (dispositivo, Edge e Cloud).
- O desempenho em relação a protocolos, latência de comunicação, processamento de equilíbrio e resiliência diante de falhas operacionais.

Portanto, o desenvolvimento da PoC envolveu a configuração de um scanner comercial e o desenvolvimento de um dispositivo inteligente desde o início, incluindo o hardware e o artefato de aprendizado de máquina. Para a rede IoT e middleware, foram desenvolvidos adaptadores para o middleware existente do UnB Internet das Coisas – UIoT [46, 50, 51]. Além disso, todos os recursos de segurança relacionados a autenticação, criptografia, validação, monitoramento e orquestração foram implementados por meio do conjunto de ferramentas do Spring Cloud, integrando a arquitetura de microsserviço com a rede IoT e fornecendo autenticação e interoperação entre dispositivos e serviços, conforme projetado para o fluxo seguro e monitorado de informações na IoTSec2DCoC.

Esse protótipo PoC foi utilizado para testar os serviços e as funcionalidades do IoTSec2DCoC com o objetivo de validar a proposta de trabalho em alguns cenários diferentes, apresentados e com seus resultados discutidos nas subseções seguintes.

### 5.1 DESCRIÇÃO DO AMBIENTE DE TESTES

Com o objetivo de realizar uma simulação dessa arquitetura, foi criado um ambiente de testes na Cloud Amazon EC2 - AWS. Foram criadas 10 máquinas virtuais conforme a Tabela 5.1.

Tabela 5.1: Máquinas Virtuais na Amazon EC2.

Name	Instance Type	Core	Memory	SO
ELK	t2.medium	02	04 Gb	Ubuntu Server v18.04
Serviços de Segurança	t2.medium	02	04 Gb	Ubuntu Server v18.04
Empacotador	t2.small	01	02 Gb	Ubuntu Server v18.04
Gateway de API	t2.small	01	02 Gb	Ubuntu Server v18.04
Middleware Edge 1	t2.small	01	02 Gb	Ubuntu Server v18.04
Middleware Edge 2	t2.small	01	02 Gb	Ubuntu Server v18.04
Middleware Cloud	t2.small	01	02 Gb	Ubuntu Server v18.04
Registro de Serviços	t2.small	01	02 Gb	Ubuntu Server v18.04
Banco de Dados	t2.small	01	02 Gb	Ubuntu Server v18.04
Archivematica /Atom	t2.xlarge	04	16 Gb	Ubuntu Server v18.04

Além disso, para o propósito de automação da implantação, o que viabilizou o desenvolvimento dos componentes da PoC usando um modelo de integração contínua, foi utilizada a tecnologia Docker para containers e Docker Compose [78] para orquestração destes, possibilitando o desenvolvimento paralelo e testes de tantos componentes que a solução possui.

### 5.1.1 Descrição do Hardware Utilizado

A configuração de hardware para a validação incluiu os seguintes dispositivos usados no protótipo da PoC:

- Scanner Fujitsu Fi – 7280.
- Dispositivo Inteligente.
- Notebook com processador Pentium I7 2.8/800/1Mb, com 32GB de memória RAM, 1 Tb SSD e placa de vídeo do GTX 1050TI 4GB.

### 5.1.2 Descrição do Ambiente de Software

A seguir, na Tabela 5.2, estão sumarizados todas as tecnologias que foram aplicadas na PoC para cada componente da arquitetura descritos na seção 4. As tecnologias serão melhor explicadas no decorrer desta seção.



Tabela 5.2: Sumarização de tecnologias aplicadas na PoC da IoTSec2DCoC.

Componentes	Tecnologias aplicadas na PoC
Scanner	Qualquer scanner comercial suportado pelo SANE. Na PoC foi utilizado o Fujitsu FI-7280.
Dispositivo Inteligente	Aplicação web e linha de comando desenvolvida usando Python e o Framework Flask, utilizando a API SANE para integração com scanners, a biblioteca OpenCV para processamento de imagem, a ferramenta Tesseract para OCR, e as bibliotecas Tensorflow, Keras e o Framework YOLO para algoritmos de <i>deep learning</i> . A aplicação é implantada em um Raspberry Pi.
Gateway de API	Ferramenta Zuul do Spring Cloud.
Middleware Edge	Middleware UIoT configurado em um nível hierárquico inferior, com a implementação adaptada ao escopo do projeto, mantido pelo laboratório LATITUDE/UnB. Foi desenvolvido usando Python e o Framework Flask, fazendo uso de uma base de dados não relacional MongoDB.
Middleware Cloud	Middleware UIoT configurado em um nível hierárquico superior, com a implementação adaptada ao escopo do projeto, mantido pelo laboratório LATITUDE/UnB. Foi desenvolvido usando Python e o Framework Flask, fazendo uso de uma base de dados não relacional MongoDB.
Empacotador	Aplicação web desenvolvida no escopo do projeto usando Java e o Framework Spring Boot.
RDC-Arq	Aplicação Archivematica, mantida pela Artefactual.
Plataforma de Acesso	Aplicação Atom, mantida pela Artefactual.
Serviço de Admissão de Arquivo Nato-Digital	Aplicação web desenvolvida no escopo do projeto usando Python e o Framework Flask.
Check-out, Check-in e Serviço de Validação de Assinaturas Digitais	Aplicações web desenvolvidas usando Python e o Framework Flask, utilizando a biblioteca cryptography para funções de assinatura digital e pyOpenSSL para validação de cadeia de confiança de certificados digitais.
Serviço de Autenticação	Aplicação web desenvolvida no escopo do projeto usando Java e o Framework Spring Boot, utilizando a biblioteca JWT para geração de <i>tokens</i> JWT e banco de dados MySQL para persistência de credenciais.
Serviço de Autoridade Certificadora	Ferramenta CFSSL da CloudFare.
Registro de Serviços	Ferramenta Eureka do Spring Cloud.
Circuit Breaker	Ferramenta Hystrix do Spring Cloud.
Serviço de Monitoramento	Conjunto de ferramentas da pilha ELK (ElasticSearch, Logstash e Kibana), as quais são um banco de dados de pesquisa textual, uma ferramenta de processamento de logs/eventos e um painel de visualização de dados, respectivamente.

### 5.1.3 Protótipo do Dispositivo Inteligente com Computador de Placa Única

A Figura 5.1 exhibe o protótipo de hardware do Dispositivo Inteligente, que é parte integrante deste trabalho, e foi criado, a partir da impressão de um *case* personalizado em impressora 3D. Sua montagem conta com um computador de placa única Raspberry PI 3 B+ com 1Gb de memória

RAM, 64GB de disco sólido e um display LCD TFT Touch 3.5.



Figura 5.1: Foto do protótipo de dispositivo montado.

### 5.1.4 Amostra de Dados

A caracterização dos dados se dá a partir da extração da informação de três tipos de documentos: Boletim de Pessoal, Despacho e “Ofício. As informações da classe Boletim de Pessoal vieram de documentos digitalizados do período de origem, do extinto MP, entre 1942 a 2008, no formato TIFF através de um scanner, totalizando 27.6 gigabytes. O texto dessas imagens é extraído a partir de um serviço de OCR, que tem como característica a extração do texto contido em uma imagem. As demais classes utilizadas para compor o *dataset*, Despacho e Ofício, foram extraídas do Sistema Eletrônico de Informações (SEI), do Ministério da Economia, antigo Ministério do Planejamento, entre o período de 2014 e 2017. Por fim, foi criado um arquivo CSV que contém as três classes com a quantidade de 1000 registros cada.

Tabela 5.3: Descrição da amostra de dados.

Documento	Classe	Base de dados
Tamanho	14.7 gigabytes	BPS
	30.0 gigabytes	SEI
Período de coleta	1942 a 2008	BPS
	2014 a 2017	SEI
Documentos	1000	BPS
	2000	SEI
Categorias	01	BPS
	02	SEI

A Tabela 5.3 contém os registros dos documentos circulados no sistema, conforme descrito, sendo que a base de dados resultante possui 44.7 gigabytes e um total de 3000 documentos digitais em formato PDF/A-1B, contendo 3 classes para classificação: Boletim de Pessoal, Despacho e Ofício. Dentro do sistema, o documento é apresentado em forma de página web e pode ser exportado como PDF.

Um subconjunto dessa amostra contendo 100 documentos é selecionado para os testes de validação em relação às medidas de segurança da proposta, conforme discutido a seguir.

## 5.2 DISPOSITIVOS INTELIGENTES DE ADMISSÃO

O Dispositivo Inteligente para admissão de documentos é projetado como um módulo de hardware e software integrado, devendo interagir com o mundo físico para capturar, processar e transmitir os documentos digitalizados para os próximos componentes da linha de produção documental, o que significa também enviar dados para as camadas superiores da rede IoT que dá suporte à CoC.

Assim, esses dispositivos de captura de documentos representam a primeira etapa no suporte à CoC de documentos digitais, autenticando a origem do documento físico digitalizado e executando a tarefa de OCR após o pré-processamento da imagem capturada para melhorar sua precisão.

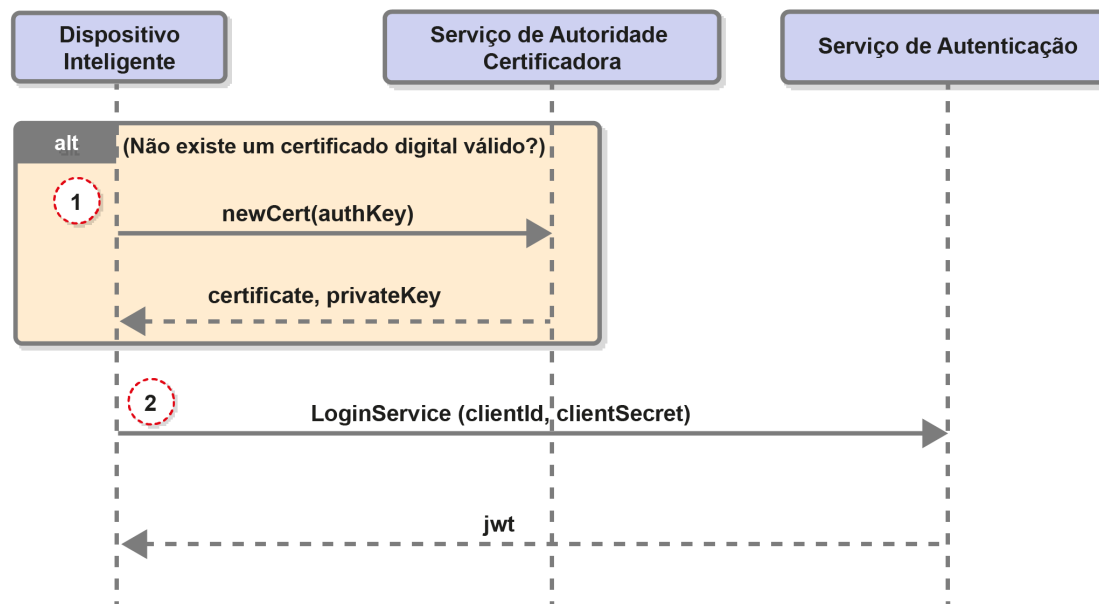


Figura 5.2: Fluxo de login do Dispositivo Inteligente na rede IOT.

Conforme mostrado na Figura 5.2, na operação desse Dispositivo Inteligente, a etapa 1 fornece a configuração dos parâmetros de autenticação e autorização para esse dispositivo operar com os

outros componentes IoT do IoTSec2DCoC. Primeiro, o componente de verificação do Dispositivo Inteligente obtém um certificado digital válido para ser usado no dispositivo executando assinaturas digitais. O dispositivo solicita ao serviço gerenciador de certificado digital um novo certificado digital e uma chave privada. Essa solicitação é permitida apenas se um administrador fornecer uma chave de autenticação previamente configurada nos Dispositivos Inteligentes.

Na etapa 2 da Figura 5.2, o Dispositivo Inteligente deve se autenticar com controladores IoT-Sec2DCoC informando suas credenciais para o Serviço de Autenticação, o qual responde fornecendo um *token* JWT.

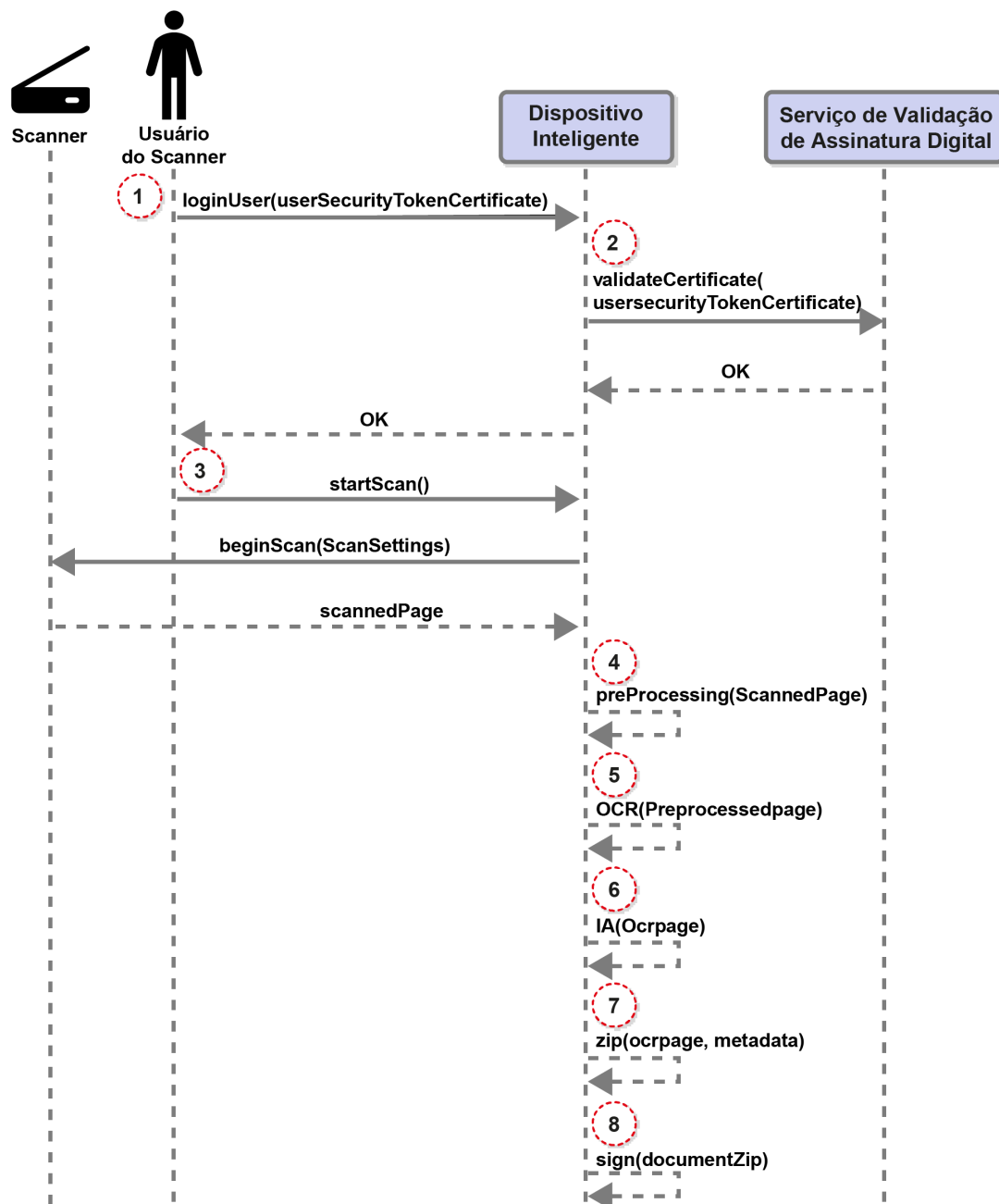


Figura 5.3: Fluxo de processamento do componente OCR.

Depois que o Dispositivo Inteligente recebe seu certificado digital e *token* JWT, o processo de digitalização é iniciado e seguido por outras etapas de processamento, conforme identificado pelos círculos vermelhos pontilhados na Figura 5.3, incluindo:

- **Etapa 1 - Login de usuário:** depois que o Dispositivo Inteligente recebe seu certificado digital e *token* JWT, o usuário pode fazer login nos controladores IoTSec2DCoC fornecendo seu *token* de certificado digital.
- **Etapa 2 - Validação do certificado digital do usuário:** durante a tentativa de login, o certificado digital do Dispositivo Inteligente é validado pelo Serviço de Validação de Certificado.
- **Etapa 3 - Digitalização:** se o certificado for aceito, o processo de digitalização começa, e o Dispositivo Inteligente comanda a digitalização das páginas. Essa integração com o scanner é realizado com a API SANE [79], o qual fornece comandos para configuração e digitalização.
- **Etapa 4 - Pré-processamento:** para o pré-processamento de imagens, para efeito de validação da proposta, utiliza-se a biblioteca de visão computacional OpenCV [80] com basicamente três procedimentos de processamento de imagens. Primeiro, há a conversão da imagem digitalizada em tons de cinza. Em seguida, o aplicativo Median Blur opera para suavizar a imagem. Por fim, existe a binarização da imagem. Essas operações visam melhorar a precisão do processo de OCR, porque esse pré-processamento dá maior ênfase ao texto do documento, removendo ruídos indesejados da imagem do documento de entrada que podem prejudicar o reconhecimento de caracteres.
- **Etapa 5 - OCR:** para fins de validação, a ferramenta Tesseract [7] é usada com módulos Python, os quais executam funcionalidades por meio do *shell* do sistema operacional. O mecanismo operacional usa *deep learning*, na forma de uma rede neural recorrente de memória de curto e longo prazo (LSTM), gerando um arquivo PDF combinado com a imagem original e um texto pesquisável na imagem combinado com um bloco de texto simples codificado em UTF-8 para inclusão no arquivo de metadados. Nesta etapa, o Tesseract abstrai a utilização do modelo LSTM, não sendo necessário criar um modelo de treinamento, pois ele utiliza modelos pré-treinados, dando suporte aos diferentes idiomas, incluindo o Português.
- **Etapa 6 - Classificação:** um algoritmo simples de classificação de documentos é proposto neste trabalho usando *deep learning* a partir do texto do documento, utilizando a ferramenta Keras [81]. Além disso, também é utilizado um outro algoritmo de extração de metadados, usando *deep learning* e técnicas de visão computacional com o Framework YOLO [82] (realiza a detecção de objetos em tempo real), onde todos os metadados extraídos são registrados em um arquivo JSON. Para fins de validação, esses algoritmos foram treinados com um conjunto de dados de documentos brasileiros oficiais pertencentes a três classes, as quais serão mais detalhadas na seção 6. Para outros conjuntos de dados documentais, esse

módulo deve ser retreinado ou substituído por um classificador equivalente. O classificador é detalhado na subseção 5.5.4.

- **Etapa 7 - Compactação:** finalmente, os dados são compactados em um único arquivo, permitindo que sua assinatura digital cubra as páginas com OCR, e o arquivo JSON com os dados.
- **Etapa 8 - Assinatura Digital:** o componente de Check-out executa a assinatura digital do arquivo compactado. Para efeito de validação da proposta, a assinatura é processada usando a biblioteca de criptografia Python [83], criando um *hash* SHA512 do arquivo compactado, criptografando-o com a chave privada do Dispositivo Inteligente.

```
(scanner) latitude@latitude-jetson:~/Documents/dev/iotsec2coc/scanner$ python -m iotbot2coc_scanner
28/06/2021 21:24:16 - scanner - INFO: [SCANNER] [CHECKPOINT] Starting processing Package_ID=0d40e3d4-06fb-4a6c-ae68-84f293baf61a Module=smart-device-1
28/06/2021 21:24:16 - scanner - INFO: SANE version = 1.0.27
28/06/2021 21:24:16 - scanner - INFO: iotbot2coc-scanner version = 0.0.1
28/06/2021 21:24:16 - scanner - INFO: Tesseract version = 4.0.0-beta.1
28/06/2021 21:24:16 - scanner - INFO: OpenCV version = 4.4.0
28/06/2021 21:24:20 - scanner - INFO: Device selected = (('device_name', 'fujitsu:fi-7280:1201817'), ('device_vendor', 'FUJITSU'), ('device_mode', 'fi-7280'), ('device_type', 'scanner'))
28/06/2021 21:24:20 - scanner - INFO: Scanner configuration = (('scanner_mode', 'Gray'), ('scanner_source', 'ADF Duplex'), ('scanner_resolution', 150), ('scanner_depth', 8))
28/06/2021 21:24:20 - scanner - INFO: Scanning started
28/06/2021 21:24:22 - scanner - INFO: Scanning finished, total pages = 2
28/06/2021 21:24:22 - scanner - INFO: [SCANNER] [CHECKPOINT] Scanning finished Package ID=0d40e3d4-06fb-4a6c-ae68-84f293baf61a Module=smart-device-1
28/06/2021 21:24:22 - scanner - INFO: (Page 1) Starting page processing
28/06/2021 21:24:22 - scanner - INFO: (Page 2) Starting page processing
28/06/2021 21:24:22 - scanner - INFO: (Page 2) Image saved
28/06/2021 21:24:22 - scanner - INFO: (Page 2) Starting image preprocessing for OCR
28/06/2021 21:24:22 - scanner - INFO: (Page 1) Image saved
28/06/2021 21:24:22 - scanner - INFO: (Page 1) Starting image preprocessing for OCR
28/06/2021 21:24:23 - scanner - INFO: (Page 1) Image preprocessing for OCR finished
28/06/2021 21:24:23 - scanner - INFO: (Page 1) Making OCR
28/06/2021 21:24:23 - scanner - INFO: (Page 2) Image preprocessing for OCR finished
28/06/2021 21:24:23 - scanner - INFO: (Page 2) Making OCR
28/06/2021 21:24:25 - scanner - INFO: (Page 1) Tesseract output = Tesseract Open Source OCR Engine v4.0.0-beta.1 with Leptonica
28/06/2021 21:24:25 - scanner - INFO: (Page 1) OCR finished (.txt and .pdf)
28/06/2021 21:24:25 - scanner - INFO: (Page 1) Page processing finished
```

Figura 5.4: Log de processamento do Dispositivo Inteligente implantado na PoC.

Na Figura 5.4, pode ser observado uma parte do log de processamento do Dispositivo Inteligente da PoC. Ao final da execução de todo o processamento, são gerados os três arquivos a seguir:

- Arquivo TIFF multi-páginas com todas as imagens digitalizadas.
- Arquivo PDF com imagens digitalizadas e texto pesquisável do OCR.
- Arquivo JSON com metadados técnicos (tanto do scanner quanto das imagens) e o texto do OCR (e.g. Código 5.1).

Código 5.1: Amostra do arquivo JSON com metadados técnicos.

```

1  "id": 1cf83a3e-aabf-4a43-8911-f7f273319770,
2  "created_at": 2019-05-28 12:03:36.687671,
3  "source": scanner,
4  "scanner_api_software_name": SANE,
5  "scanner_api_software_version": 1.0.25,
6  "capture_software_name": iotbot2coc-scanner,
7  "capture_software_version": 0.0.1,
8  "ocr_software_name": Tesseract,
9  "ocr_software_version": tesseract_version,
10 "image_processing_software_name": OpenCV,
11 "image_processing_software_version": 3.4.4,
12 "os_release": Raspbian GNU/Linux 9.8 (stretch),
13 "os_system_info": Linux raspberrypi 4.14.98-v7+ #1200 SMP Tue Feb 12 20:27
    :48 GMT 2019 armv7l GNU/Linux,
14 "device_name": fujitsu:fi-7280:1201852,
15 "device_type": scanner,
16 "scanner_mode": Gray,
17 "scanner_source": ADF Duplex,
18 "scanner_resolution": 300,
19 "scanner_depth": 8,
20 "pdf_hashsum": {SHA256}
    eb5edac81ee26fa5842b93b17f2df6b75101e27fc79f68a61e44d7b0d71bf953,
21 "pdf_file_size": 10061369,
22 "tiff_hashsum": {SHA256}
    fde3072c1814d34f88303de1f04b17e5f4f821a9c0852b85c8aa10ead0b89c3c,
23 "tiff_file_size": 16830368,
24 "tiff_dimension_width": 2550,
25 "tiff_dimension_height": 3300,
26 "text_pages": [\n\nMP\n nSUM\u00c1RIO\n\n13.11.2007 11.16\n\nBOLETIM DE
    PESSOAL E SERVI\u00c7O\nEDI\u00c7\u00c3O ESPECIAL\n\nORIGEM\nSECRETARIA
    DE OR\u00c7AMENTO FEDERAL (SOF)\n\nPALAVRA CHAVE\nPORTARIA SOF N\u00ba
    053 DE 08 DE NOVEMBRO DE 2007\n\nAtos Administrativos\n\nf, SoF\n
    nPORTARIA N\u00ba 55 DE 08 DE NOVEMBRO DE 2007.\n\nDisp\u00f5e sobre a
    participa\u00e7\u00e3o de servidores da\nSecretaria de Or\u00e7amento
    Federal em programas de\ncapacita\u00e7\u00e3o de curta, m\u00e9dia e
    longa dura\u00e7\u00e3o,\n\nO SECRET\u00c1RIO DE OR\u00c7AMENTO FEDERAL
    , SUBSTITUTO, no uso das atribui\u00e7\u00f5es\n\u00e7\u00e3o pelo
    inciso V do art. 13 do Cap\u00edtulo IV do Anexo IX da Portaria MP n\u00ba
    232, de 3 de agosto de\n2005 observado o disposto no Decreto n\u00ba
    5.707, de 23 de fevereiro de 2006, resolve:\n\nCAP\u00cdTULO I -
    DAS DISPOSI\u00c7\u00d5ES PRELIMINARES\n\nArt. 1\u00ba A participa\u00e7\u00e3o
    de servidores da Secretaria de Or\u00e7amento Federal \u2014
    u2014 SOF em eventos de\ncapacita\u00e7\u00e3o de curta, m\u00e9dia e
    longa dura\u00e7\u00e3o, passa a ser regulamentada por esta Portaria
    com os seguintes\nobjetivos:\n\n1 - promover a melhoria da efici\u00eancia,
    efic\u00e1cia e qualidade dos servi\u00e7os prestados pela
    SOF;\n\n1 - valorizar os servidores, por meio de sua capacita\u00e7\u00e3o
    permanente e adequa\u00e7\u00e3o aos novos\nperfis
    profissionais requeridos no setor p\u00fablico; e\n\n1 - racionalizar
    os investimentos com capacita\u00e7\u00e3o.\n\nArt. 2\u00ba A capacita\u00e7\u00e3o
    visa proporcionar ao servidor:\n\n1 -
    oportunidades de aquisi\u00e7\u00e3o de conhecimentos, habilidades e
    atitudes necess\u00e1rias ao seu\ndesempenho profissional, dentro de
    sua \u00e1rea de atua\u00e7\u00e3o;\n\n1 - acesso a novas tecnologias,
    m\u00e9todos e procedimentos, estimulando a pesquisa e atualiza\u00e7\u00e3o
    \nprofissional;\n\n1 - incentivo ao seu autodesenvolvimento e ao
    desenvolvimento profissional cont\u00ednuo; e\n\n1 - aumento de sua
    efici\u00eancia e desempenho.\n\nArt. 3\u00ba Para fins desta Portaria,
    consideram-se a\u00e7\u00f5es de capacita\u00e7\u00e3o:\n\n1 - cursos
    presenciais e \u00e0 dist\u00e2ncia;\n\n1 - aprendizagem em servi\u00e7o
    ;\n\n1 - semin\u00e1rios;\n\n1 - congressos;]
```

Como o Dispositivo Inteligente irá interagir com o Middleware Edge, esse middleware precisa de detalhes sobre o dispositivo e os serviços que ele fornece. Portanto, o dispositivo deve registrar seus serviços com o middleware conforme mostrado na Figura 5.5. Considerando essa figura, se o Dispositivo Inteligente ainda não estiver registrado no Middleware Edge, ele inicia seu registro fornecendo os dados do cliente junto com o JWT, previamente adquirido (etapa 1). Em seguida, o Middleware Edge valida o JWT (garantindo sua autenticidade) e permite que o dispositivo inteligente forneça seus dados de descrição de serviço, ou seja, o serviço de envio de documentos digitalizados e metadados relacionados (etapa 2).

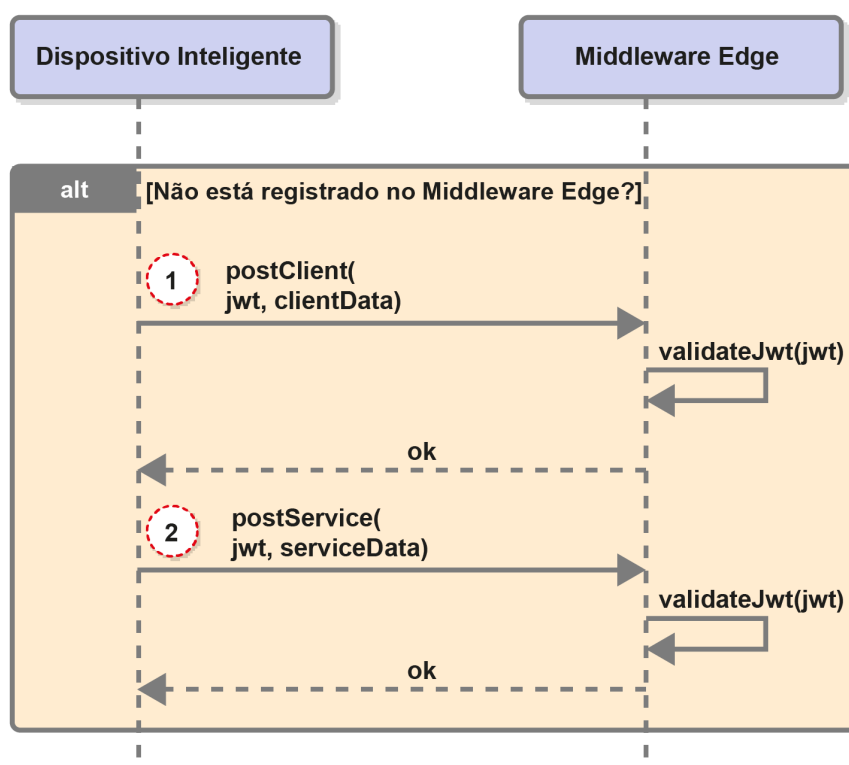


Figura 5.5: Registro do Dispositivo Inteligente na rede IOT.

Após seu registro no Middleware Edge, o Dispositivo Inteligente é capaz de enviar documentos digitalizados, pré-processados e classificados para a rede IoT, de acordo com o processo mostrado na Figura 5.6. O Dispositivo Inteligente envia os documentos para o Middleware Edge, o qual encaminha os documentos recebidos para o Middleware Cloud. Na etapa 1, o Dispositivo Inteligente envia o pacote de assinatura, composto pelo documento digitalizado, seus metadados e a assinatura digital, junto com o JWT, o qual é validado pelo Middleware Cloud. Em seguida, na etapa 2, o Middleware Edge valida a assinatura digital do Dispositivo Inteligente, acessando o Serviço de Validação de Assinatura Digital e, na etapa 3, envia os dados para o Middleware Cloud, que verifica novamente a assinatura digital do documento, verificando, assim, a autenticidade da origem e a integridade do documento. Essa verificação emparelhada da assinatura do documento é repetida em cada transferência subsequente do documento, de modo que as propriedades de autenticação e integridade sejam mantidas em toda a CoC documental.



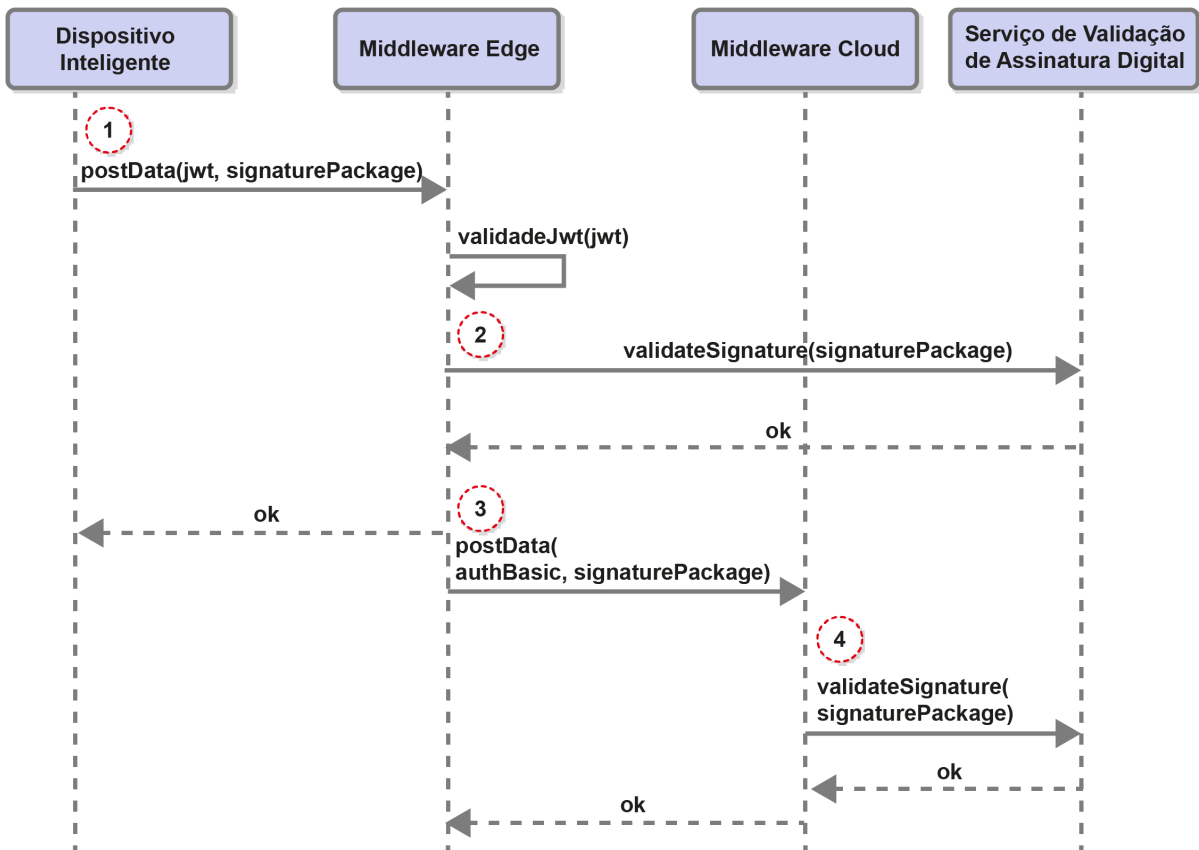


Figura 5.6: Transferência de pacote de documentos do Dispositivo Inteligente para o Middleware IoT Edge.

### 5.3 IOT E SUPORTE DE MICROSERVIÇOS PARA COMUNICAÇÕES E INTEGRAÇÃO

A IoTSec2DCoC associa uma instância de rede IoT com microsserviços integrados para fornecer autenticação de dispositivo, registro de serviços e roteamento inteligente para pacotes de documentos, os quais sistematicamente verificados quanto à autenticidade e integridade durante as operações de transferência.

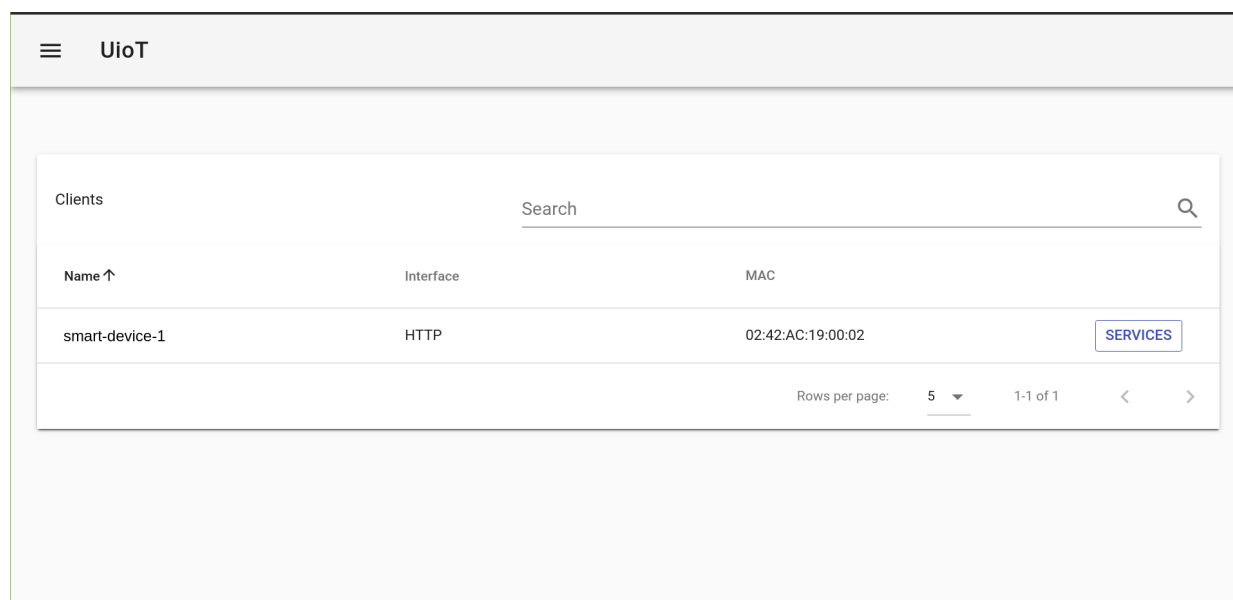
Depois de coletados pelo Middleware Cloud, os documentos ficam disponíveis para as fases subsequentes da linha de produção documental. O Middleware periodicamente envia os documentos recém-processados para o Empacotador, permitindo a revisão humana da digitalização. Os documentos validados são encaminhados para o RDC-Arq e, posteriormente, para a Plataforma de Acesso para fins de difusão. A verificação de autenticidade e integridade é realizada nas transferências do Empacotador para o RDC-Arq e para a Plataforma de Acesso.

Na arquitetura proposta, duas partes de Middleware IoT, ou seja, respectivamente o Middleware Edge e o Middleware Cloud, operam em um arranjo hierárquico: o primeiro módulo atuando na borda, próximo aos dispositivos para os quais fornece serviços de comunicação e in-

tegração; enquanto o segundo é responsável pelo gerenciamento de toda a rede IoT, fornecendo a interface vinculada aos aplicativos IoT.

Para o protótipo da PoC, usado para fins de validação neste trabalho, o Middleware Hierárquico denominado UIoT [52] é usado, pois esse middleware compreende um conjunto escalável de componentes, como um gateway IoT projetado para a interação com coisas, serviços de dados distribuídos com técnicas de otimização de armazenamento e comunicação, e um módulo de painel tanto para visualização de dados como para interface administrativa de usuário.

Originalmente, o Middleware UIoT possuía somente suporte para trabalhar com dados numéricos, sendo estes enviados pelos seus clientes (e.g. Dados de temperatura, pressão, distância, etc.), o que é o mais usual para dispositivos IoT. Uma das funções do Middleware Hierárquico é justamente realizar algum tipo de agregação dos dados de um nível inferior para um nível superior, realizando operações de soma, média e afins. No caso do IoTSec2DCoC, os dados transmitidos pelos dispositivos inteligentes são textuais, o que originou a necessidade de alteração do UIoT para suportar tais dados, principalmente no componente Engine (exemplificado na Figura 2.9), sendo criada uma entidade de "Rule" específica para dados textuais, fora outras alterações na manipulação dos dados. Além disso, teve que ser desenvolvido uma nova interface para o componente de "Dispatcher" para envio dos dados do Middleware Cloud ao Empacotador, semelhantemente como ocorre na integração da instância do Middleware Edge com a instância *cloud*. Outra alteração também ocorrida para suportar dados textuais foi no componente DIMS. Como este usa o banco de dados MongoDB em sua utilização tradicional, suporta documentos de até 16MB. Além disso, deve ser utilizada uma outra técnica de manipulação de dados chamada GridFS. Na Figura 5.7, pode ser visualizada a tela do *dashboard* do Middleware UIoT Edge utilizado na PoC.



The screenshot shows a web interface for 'UioT'. At the top left, there is a hamburger menu icon and the text 'UioT'. Below this is a search bar labeled 'Clients' with a search icon on the right. The main content is a table with the following data:

Name ↑	Interface	MAC	
smart-device-1	HTTP	02:42:AC:19:00:02	<a href="#">SERVICES</a>

At the bottom right of the table, there is a pagination control showing 'Rows per page: 5' and '1-1 of 1' with navigation arrows.

Figura 5.7: Dashboard para visualização de dados do UIoT implantado na PoC.

Em relação aos microsserviços, esses artefatos de software são integrados à instância da rede

IoT para fornecer um conjunto de serviços adequado às medidas de segurança que dão suporte à CoC documental. Esses serviços incluem autenticação de dispositivo, registro de recursos de dispositivo, registro de módulos de Middleware IoT e controle de roteamento em toda a hierarquia desses módulos. O funcionamento integrado desses serviços é discutido com mais detalhes na subseção 5.4.

Os microsserviços são fornecidos por meio do Spring Cloud, a qual opera em um ambiente de nuvem distribuída e possui regras de implantação de configuração simples. A Figura 5.8 mostra a arquitetura Spring Cloud com seu Gateway de API, usada na IoTSec2DCoC. É utilizado um único ponto de comunicação entre diferentes terminais de API e vários clientes que precisam acessá-los, formando uma coleção de serviços fracamente acoplados.

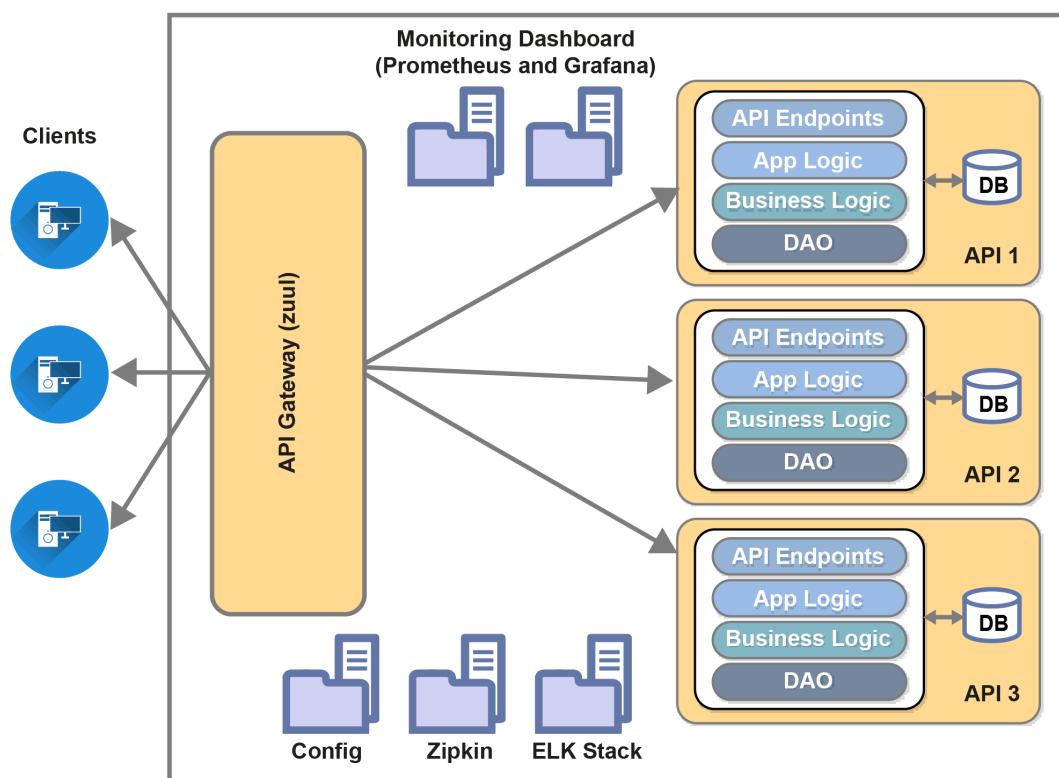


Figura 5.8: Arquitetura com um servidor Edge. Adaptado de: [84].

Essa forma de implantação de microsserviços é útil para a IoTSec2DCoC lidar com medidas de segurança e monitoramento de rede IoT, pois, por meio de um único ponto de entrada, os serviços de software necessários estão disponíveis, ao invés de entregar e monitorar cada serviço em um canal separado. Outra característica é a tolerância a falhas e balanceamento de carga. A partir da utilização de algoritmos como o Round-robin, o Gateway de API (Zuul) define qual serviço deve receber a solicitação e clona novas instâncias para qualquer serviço que esteja recebendo requisições excessivas.

Essas tarefas são apoiadas por um processo de controle de descoberta de serviço, o qual regis-

tra todos os serviços que fazem parte da rede Spring Cloud. Devido ao escalonamento automático de serviços e mudanças constantes de IP, todo serviço, assim que entra na rede Spring Cloud, precisa se registrar com descoberta de serviço para que o Gateway de API possa encaminhar a rota para o serviço correto. O desempenho desses serviços é monitorado por um módulo (Zipkin) que rastreia todas as rotas direcionadas pelo Gateway de API, apresentando, assim, um log detalhado de todas as chamadas, se tiveram sucesso ou não, bem como o tempo gasto com cada solicitação. Essa ferramenta compõe o que foi chamado no modelo de Serviço de Monitoramento

Essa configuração de microsserviços em uma instância de rede IoT orquestra os serviços de segurança integrados IoTSec2DCoC que reforçam a CoC documental, conforme descrito a seguir.

## **5.4 SERVIÇOS DE SEGURANÇA INTEGRADOS**

Esta subseção detalha os componentes e processos que visam aumentar a segurança da arquitetura proposta IoTSec2DCoC, considerando os problemas relatados na seção 3.1.

### **5.4.1 Serviço de Autoridade Certificadora**

As comunicações entre todos os módulos e componentes devem ser configuradas para utilizar um canal seguro do protocolo TLS [85], operando em todos os serviços disponíveis. Esse protocolo requer certificados digitais X.509 para a criptografia de chave pública, a qual permite a autenticação do servidor e que é usada na fase inicial da conexão TLS para realizar a troca de uma chave de criptografia simétrica gerada aleatoriamente, sendo esta, por sua vez, usada posteriormente para criptografar todos os dados durante as comunicações. Essa é uma medida comum de segurança padrão para evitar que um agente malicioso interprete as mensagens trocadas pelos componentes da estrutura IoTSec2DCoC.

A confidencialidade das comunicações em TLS é suportada por uma infraestrutura de chave pública (Public Key Infrastructure - PKI) para fornecer confiança em relação aos certificados digitais usados. Esses certificados também servem no IoTSec2DCoC para realizar e validar as assinaturas digitais dos documentos processados, garantindo sua autenticidade de origem e integridade, serviço este descrito na seção 5.4.2.

A PKI também é essencial para gerenciar o ciclo de vida dos certificados digitais, incluindo estados e transições, desde sua emissão até sua revogação, além de procedimentos de recuperação, validação e renovação. O gerenciamento das chaves associadas aos certificados digitais é uma tarefa demorada e arriscada para ser realizada manualmente pelos administradores do sistema. Assim, idealmente, essa tarefa deve ser totalmente automatizada em um ambiente dinâmico como a IoT, usando uma Autoridade Certificadora (AC), a qual já está integrada no modelo proposto.

Em relação a essa proposta de protótipo para PoC, a ferramenta escolhida para a PKI é o CFSSL [86], um módulo usado internamente pela Cloudflare para suas atividades de AC. Essa

ferramenta fornece um utilitário de linha de comando para lidar com operações de certificados digitais X.509, além de uma API HTTP que permite realizar as mesmas operações remotamente, incluindo: criação de chaves privadas, solicitações de assinatura de certificado digital (CSR), certificados digitais, cadeias de certificados digitais de confiança e assinatura de certificados digitais.

Implantar a AC como entidade local na proposta da IoTSec2DCoC é interessante, pois permite que a CoC documental opere de forma autônoma, ou seja, independente de uma PKI externa, como ICP-Brasil [87] ou WebTrust [88], uma vez que essa AC fornece à IoTSec2DCoC um certificado de autoridade autoassinado e uma API HTTP disponível para os outros componentes. Assim, esses componentes podem solicitar seus certificados e consultar uma lista de certificados revogados que está disponível para a validação das propriedades de autenticidade e integridade dos documentos processados. Portanto, dentro da IoTSec2DCoC, os certificados digitais são emitidos para a AC raiz, AC intermediária e certificados folha para dispositivos inteligentes, Middleware IoT e qualquer outro aplicativo autorizado.

Dentro dessa proposta de protótipo PoC, a AC CFSSL usa arquivos de configuração JSON para definir os atributos de cada certificado gerado. Por exemplo, o arquivo usado para a CA raiz especifica que o algoritmo de autenticação é Algoritmo de Assinatura Digital de Curva Elíptica (ECDSA) com uma chave de 384 bits [89]. Aos serviços que irão receber esse certificado e fornecem algum serviço HTTP é configurada a conexão TLS 1.3 [85], o algoritmo ECDHE para troca de cifras, AES para encriptação de mensagens e SHA como função de *hash*.

A identificação do certificado digital seguirá o padrão CN = <Nome>, OU = LATITUDE, O = UnB, L = Brasilia, ST = DF, C = BR, tendo validade de dez anos.

Código 5.2: Arquivo de configuração JSON para a AC Raiz do CFSSL.

```
{ "CN": "AC Raiz IoTSec2DCoC",
  "key": {
    "algo": "ecdsa",
    "size": 384 },
  "names": [{
    "C": "BR",
    "ST": "DF",
    "L": "Brasilia",
    "O": "UnB",
    "OU": "LATITUDE"
  }],
  "ca": {
    "expiry": "87600h"}}
```

Quando uma solicitação de certificado vem de um componente, são criados dois arquivos de configuração JSON. Um para o servidor, definindo a validade e o uso do certificado digital, sendo os possíveis usos: assinatura digital, codificação de chaves, autenticação de cliente e servidor. O outro arquivo JSON é enviado ao cliente, o qual precisará aceitar o certificado digital emitido. Segue, no Código 5.3, esses exemplos de configurações.

Código 5.3: Arquivo de Configuração JSON para a Solicitação de Certificado Digital “Folha”.

```
# Server JSON
{"signing": {
  "default": {
    "usages": [
      "digital signature",
      "key encipherment",
      "server auth",
      "client auth"
    ],
    "expiry": "8760h"}
}
}
# Client JSON
{"CN": "smart-device1",
  "key": {
    "algo": "ecdsa",
    "size": 384
  },
  "names": [
    {
      "C": "BR",
      "ST": "DF",
      "L": "Brasilia",
      "O": "UnB",
      "OU": "LATITUDE" }]
}
```

#### 5.4.2 Serviços de Assinatura Digital e Validação de Documentos (Checkout e Check-in)

Na estrutura IoTSec2DCoC, os componentes da solução geram ou transformam documentos, operando no conteúdo PDF, TIFF e metadados relacionados, devendo assinar digitalmente o arquivo compactado resultante contendo os documentos. Por outro lado, os mesmos componentes devem validar as assinaturas digitais com os arquivos compactados, oriundos de componentes anteriores no fluxo

Para os nós da IoTSec2DCoC que precisam realizar a geração da assinatura digital e a verificação da assinatura digital correspondente, dois módulos genéricos são projetados, respectivamente denominados componentes de Check-out e de Check-in. Ambos foram desenvolvidos para o protótipo da PoC como aplicativos web Python com o Framework Flask. O artefato de Check-out, operando no lado da saída de um nó, gera um *hash* SHA de 512 bits do arquivo compactado que contém o documento digitalizado e, em seguida, assina digitalmente o documento, criptografando esse *hash* com uma chave privada, usando a biblioteca de criptografia Python. Por outro lado, o componente de Check-in, operando no lado da entrada de um nó, verifica a assinatura digital adicionada ao arquivo compactado solicitando um serviço da web chamado Serviço de Validação de Assinatura Digital, que pode realizar a validação usando a chave pública de Check-out anterior. O Serviço de Validação de Assinatura Digital também valida o certificado digital do signatário, verificando a cadeia de confiança das ACs que emitiram esse certificado. Conceitualmente, os pares que compreendem um componente Check-out e um componente Check-in constituem o denominado Dispositivo IoT Virtual Verificador, mostrado na Figura 4.1, o qual impõe criptografia e

controles de *endpoint* em todos os dispositivos e artefatos de software que transferem documentos na IoTSec2DCoC.

### 5.4.3 Serviço de Monitoramento

Para o gerenciamento de eventos e informações de segurança, foi utilizada uma solução composta por softwares livres que realiza a tarefa de gerenciador de logs e eventos: a pilha ELK [90], acrônimo para Elasticsearch, Logstash e Kibana. O Elasticsearch é um banco de dados NoSQL especializado em consultas por texto livre. O Logstash é responsável por realizar a coleta, indexação, filtragem e agrupamento dos registros de log, o Kibana atua como um *dashboard* de controle para visualização dos dados armazenados e criação de painéis analíticos sobre os logs gerados. Além desses três componentes ainda existe um quarto chamado Filebeat, que atua como um cliente de envio de logs, implantado em cada módulo da solução que terá os registros coletados, sendo responsável por encaminhá-los para o serviço Logstash. Dentro da solução do IoT-Sec2DCoC, existe uma instalação da pilha ELK para centralização dos logs de todos os serviços da solução, inclusive de todos os Dispositivos Inteligentes, justamente para atingir os benefícios anteriormente citados.

### 5.4.4 Serviço de Autenticação (JWT)

O Serviço de Autenticação da IoTSec2DCoC autentica usuários, como os operadores do Dispositivo Inteligente e do Empacotador. Os usuários devem apresentar um identificador individual, como CPF ou CNPJ, e senha para acesso a um certificado digital armazenado em banco de dados específico e, em seguida, um *token* criptográfico no formato JWT. Uma autenticação análoga é fornecida para Dispositivos Inteligentes os quais permitem que esses dispositivos operem na rede IoT, uma vez que a chave simétrica usada para assinar o JWT é entregue ao Middleware IoT, permitindo, assim, que o Middleware Edge autentique Dispositivos Inteligentes na rede. No exemplo do formato do JWT no Código 5.4, é possível identificar que cada uma das três seções do JWT são codificadas em Base64 e são unidas pelo caractere “.”, sendo esse o JWT codificado que é enviado no cabeçalho das requisições HTTP. O exemplo no Código 5.5 é a versão codificada referente ao Código 5.4.

Código 5.4: Estrutura do JWT de um Dispositivo Inteligente.

```
# HEADER: ALGORITHM & TOKEN TYPE
{"alg": "HS512",
 "typ": "JWT"}
# PAYLOAD: DATA
{"sub": "smart-device1",
 "iat": 1516239022}
# VERIFY SIGNATURE
HMACSHA512(
  base64UrlEncode(header) + "." +
  base64UrlEncode(payload), <PRIVATE_KEY_HERE>)
```

Código 5.5: JWT codificado em Base64 de um Dispositivo Inteligente.

```
eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9  
.eyJzdWIiOiJzbWVydC1kZXZpY2UxIiwiaWF0IjoxNTE2MjM5MjYyOTQ  
.4XeLftDFeJdfjUJy84bkm3GbTDfNZsn7SznwZmSQvc8
```

## 5.5 MÓDULOS DE APLICAÇÃO

Os componentes da CoC convencional mostrados na Figura 3.1 são integrados à IoTSec2DCoC como aplicativos que consomem recursos de IoT e microsserviços de segurança. Esses aplicativos incluem o Empacotador, o RDC-Arq e a Plataforma de Acesso. O Empacotador é considerado um aplicativo interno na proposta, pois participa diretamente do processo de produção documental, principalmente nas tarefas de controle de qualidade; enquanto o RDC-Arq e a Plataforma de Acesso são considerados aplicativos externos plugáveis que operam para fins de preservação e divulgação sem interferir na produção de documentos. Esse arranjo é uma variação possível do modelo OAIS, uma vez que tal modelo oferece apenas uma referência para suas entidades funcionais, permitindo, então, serem combinadas ou divididas e distribuídas em diferentes aplicações [91]. As subseções a seguir apresentam os detalhes dessa distribuição entre os aplicativos propostos para a IoTSec2DCoC.

### 5.5.1 Empacotador

Este componente de empacotamento permite a revisão de documentos digitalizados por arquivistas humanos, ou seja, após esses documentos serem digitalizados, são tratados pelo Dispositivo Inteligente. A tarefa de revisão pode recusar arquivos digitalizados de baixa qualidade e pode determinar que o documento seja digitalizado novamente, contudo, para os documentos aceitos na revisão, mais metadados são incluídos. Consequentemente, a principal funcionalidade do empacotador é permitir ao arquivista adicionar os metadados do documento aos arquivos e inserir todas as partes do documento em um pacote a ser enviado ao repositório confiável. Na Figura 5.9, pode ser visualizada a tela de revisão de pacotes do Empacotador utilizado na PoC. A construção desse pacote segue um método de arquivamento denominado “Normalização Manual” [92], o qual está comumente disponível em RDC-Arq’s, como Archivematica, RODA [93] e LOCKSS [94], mas que foi integrado à IoTSec2DCoC, pois, para proteger a tarefa conforme exigido, é necessário ler arquivos criptografados, uma função que normalmente não é encontrada, nem facilmente implementável em repositórios confiáveis.



Módulo de Captura de Documentos Arquivísticos - MCD-Arq  
 MINISTÉRIO DA ECONOMIA - ME

Administrador (administrador) 

Início / Revisar

Digitalizador   Validador   Cadastrador   **Revisor**   Conferente

Revisões pendentes

Nome do pacote	Identificador	Localização funcional	Motivo	Arquivo(s)	Ações
IoTSec2CoC  Boletim de serviço de 07.01.2007 (digitalizado em 30-05-2021 às 18:56:59)	c65d2800-	ME - SEDE		<a href="#"> IoTSec2CoC </a>	REVISAR
	9378-4fc9-			<a href="#">Boletim de serviço de 07.01.2007</a>	RECUSAR
	b594- 698bc52d6342			<a href="#">(digitalizado em 30-05-2021 às 18:56:59).pdf</a>	ACEITAR

Figura 5.9: Tela de revisão de pacotes do Empacotador implantado na PoC.

O Empacotador produz o SIP contendo documentos e informações de descrição relacionadas, os quais são enviados para a entidade de admissão do RDC-Arq (ingestão). Como pode ser observado na Figura 5.10, a estrutura do pacote SIP tem um diretório denominado “*reservation*”, o qual contém os arquivos que serão utilizados para criar o pacote AIP (que são preservados no repositório), e um diretório denominado “*access*”, o qual contém os arquivos que serão utilizados para criar o pacote DIP (que serão disseminados para o Atom). O AIP irá conter o arquivo TIFF, o qual é destinado à preservação, e o DIP conterà o arquivo PDF destinado ao acesso.

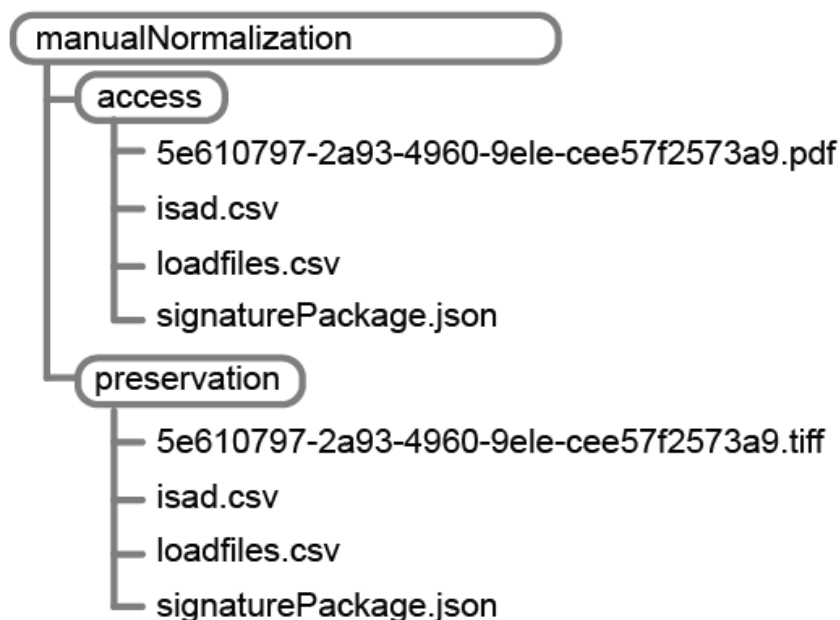


Figura 5.10: Estrutura do SIP transferida do Empacotador para o RDC-Arq.

Dentro do SIP, o arquivo “isad.csv” contém metadados ISAD(G) [28] organizados nos se-

guintes campos: o código de referência é obtido a partir do Dispositivo Inteligente e representa os itens de classificação do documento, incluindo título, nível de descrição, dimensão, suporte de material, nome do produtor, tipo de documento, entre outros. O arquivo “loadfiles.csv” serve unicamente para indicar qual será a localização do documento na árvore de categorias de documentos do Atom, já o “signaturePackage.json” contém os dados de assinatura digital para validação pelo componente de Check-in.

### 5.5.2 Repositório Arquivístico Digital Confiável - RDC-Arq

A função central deste aplicativo é fornecer armazenamento seguro para os documentos produzidos. Na proposta IoTSec2DCoC, essa função é modelada para um componente plugável, uma vez que existem vários sistemas que fornecem as funcionalidades necessárias. Os requisitos para esses sistemas serem conectados à IoTSec2DCoC são três: devem fornecer um módulo de ingestão para receber pacotes SIP do empacotador; transformar cada SIP em um AIP para armazenamento no RDC-Arq; e, finalmente, o aplicativo deve ser capaz de recuperar os pacotes AIP armazenados e transformá-los em pacotes DIP a serem fornecidos para a Plataforma de Acesso.

Para o protótipo IoTSec2DCoC PoC, o Archivematica [95] foi escolhido por ele atender às exigências mencionadas, sendo capaz de armazenar o AIP em um diretório local ou remoto, conforme definido pelo administrador. Na Figura 5.11, pode ser visualizada a tela de transferência de pacotes do Archivematica utilizado na PoC, já na Figura 5.12, pode ser visualizada a tela de administração de pacotes, onde são listados o AIP e o DIP do exemplo de uso.

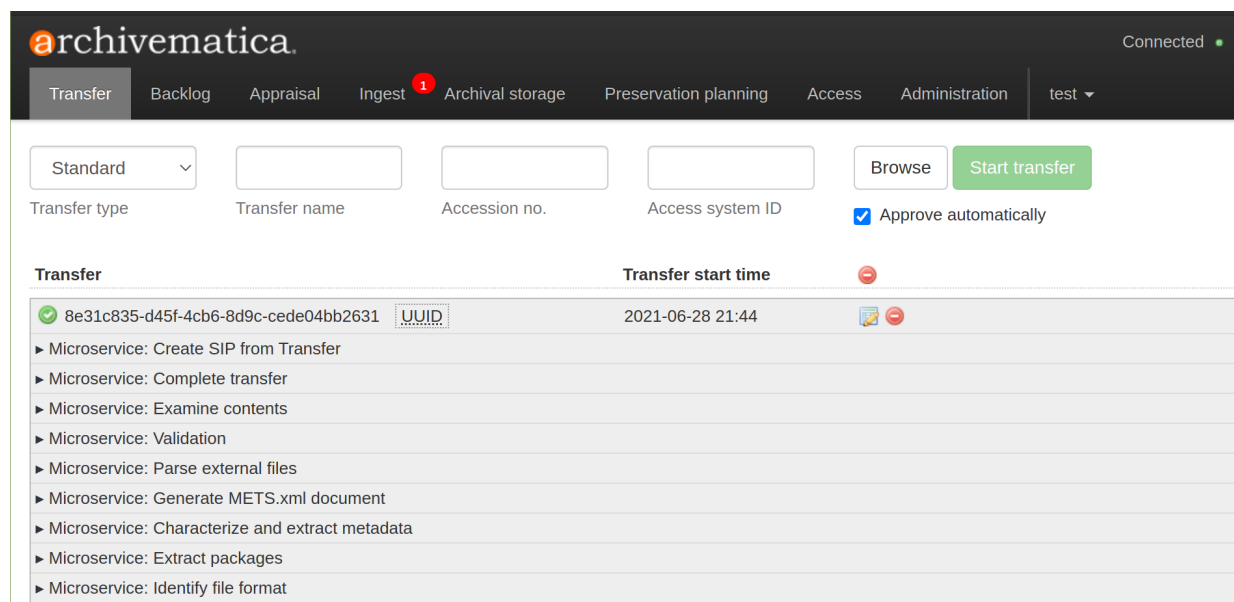


Figura 5.11: Tela de transferência de pacotes do Archivematica implantado na PoC.

Archivematica Storage Service Home Pipelines Spaces Locations Packages Administration Log out

## All Packages

Packages are Transfers, SIPs, DIPs and AIPs uploaded to a Location managed by the storage service.

[View recovery requests](#) | [View delete requests](#)

Show  entries Search:

UUID	Originating Pipeline	Current Location	Size	Type	Replica Of
f2504772-afe2-4554-95c1-4e7ada42b70d	Archivematica on b47f5f2c7f17 (9b081965-a3d3-4cf9-a93c-fc4121a17b24)	/var/archivematica/sharedDirectory/www/DIPsStore/f250/4772/afe2/4554/95c1/4e7a/da42/b70d/8e31c835-d45f-4cb6-8d9c-cede04bb2631-d1483351-e44c-484f-8afb-75f5b75b62f1	200.7 KB	DIP	
d1483351-e44c-484f-8afb-75f5b75b62f1	Archivematica on b47f5f2c7f17 (9b081965-a3d3-4cf9-a93c-fc4121a17b24)	/var/archivematica/sharedDirectory/www/AIPsStore/d148/3351/e44c/484f/8afb/75f5/b75b/62f1/8e31c835-d45f-4cb6-8d9c-cede04bb2631-d1483351-e44c-484f-8afb-75f5b75b62f1.7z	367.3 KB	AIP	

Showing 1 to 2 of 2 entries (filtered from 38 total entries) ◀ Previous Next ▶

Figura 5.12: Tela de administração de pacotes (AIP e DIP) do Archivematica implantado na PoC.

O Archivematica é um Repositório Arquivístico Digital Confiável, de acordo com os requisitos do Conselho Nacional de Arquivos, pois trata-se de um sistema de preservação digital projetado para manter os dados baseados em padrões de preservação digital e acesso em longo prazo para conjuntos de objetos digitais.

Conforme especificações, o Archivematica utiliza os seguintes padrões para identificação e codificação dos objetos digitais:

- *Metadata Encoding and Transmission Standard* (METS) [96]: padrão de transmissão e codificação de metadados, aplicado à codificação de metadados por meio de um esquema XML padronizado.
- *Preservation Metadata* (PREMIS) [97]: esquema XML desenvolvido para metadados de preservação baseado em entidades e unidades semânticas que abrigam informações sobre um objeto digital necessário para apoiar e registrar ações de preservação digital.
- *Dublin Core* (DCMI) [98]: esquema de metadados que visa descrever objetos digitais, tais como vídeos, sons, imagens, textos e sites na web.
- *BagIt* [99]: conjunto de convenções de sistema de arquivos hierárquicas projetadas para suportar armazenamento baseado em disco e transferência de rede de conteúdo digital arbitrário, utilizado pelo Archivematica para gerar os Pacotes de Informações de Arquivamento (AIPs), sendo confiáveis, autênticos, seguros e independentes do sistema para armazenamento em seu repositório preferido.

O repositório utiliza uma abordagem de microsserviços em arquitetura Pyton. Essa arquitetura é orquestrada de forma integrada a ferramentas livres e de código aberto, permitindo aos usuários

processar objetos digitais para a preservação digital através do processamento dos pacotes de submissão, preservação e disseminação do modelo OAIS [3] e as demais normas de preservação digital e melhores práticas.

Outra característica da solução é o Registro da Política de Formatos (FPR), onde suas políticas de formato padrão são registradas. A FPR é flexível para a identificação de formato, extração de pacote, transcrição e normalização. O FPR é integrado com o PRONOM, um projeto do Arquivo Nacional Britânico que possui um repositório de registro técnico de formato de arquivos digitais com acesso via *webservice*, consumindo mais de 1000 formatos de arquivos registrados, possibilitando, ainda, atualizar as ferramentas, regras e comandos do FPR local, conforme a Figura 5.13.

	Description	Group	Actions
Rules			
Commands	3D Studio (x-fmt/19)	Image (Raster)	View   Edit
Format policy registry	3D Studio Shapes (x-fmt/102)	Model	View   Edit
Tools	3DM (x-fmt/432, x-fmt/433, x-fmt/434, ...)	Model	View   Edit
Characterization	3DS Max (fmt/978)	Model	View   Edit
Rules	3GPP Audio/Video File (fmt/357)	Video	View   Edit
Commands	3MF 3D Manufacturing Format (fmt/829)	Text (Markup)	View   Edit
Event Detail	4X Movie File (fmt/1150)	Video	View   Edit
Rules	7 Zip (fmt/484)	Package	View   Edit
Commands	7-bit ANSI Text (x-fmt/21)	Text (Plain)	View   Edit
Extraction	7-bit ASCII Text (x-fmt/22)	Text (Plain)	View   Edit
Rules	8-bit ANSI Text (x-fmt/282)	Text (Plain)	View   Edit
Commands	8-bit ASCII Text (x-fmt/283)	Text (Plain)	View   Edit
Normalization	AAE Sidecar Format (fmt/980)	Text (Markup)	View   Edit
Rules	AbiWord Document (fmt/890)	Word Processing	View   Edit
Commands	AbiWord Document Template (fmt/891)	Word Processing	View   Edit
Transcription	ACBM Graphics (x-fmt/301)	Image (Raster)	View   Edit
Rules	AccessData Custom Content Image (archivematica-fmt/2, archivematica-fmt/3, fmt/842, ...)	Disk Image	View   Edit
Commands	Acrobat Catalog Cat File (fmt/452)	Portable Document Format	View   Edit
Validation	Acrobat Language definition file (x-fmt/427)	Text (Source Code)	View   Edit
Rules	Acrobat PDF 1.0 - Portable Document Format (fmt/14)	Portable Document Format	View   Edit
Commands	Acrobat PDF 1.1 - Portable Document Format (fmt/15)	Portable Document Format	View   Edit
Verification	Acrobat PDF 1.2 - Portable Document Format (fmt/16)	Portable Document Format	View   Edit
Rules			

Figura 5.13: Registro da Política de Formatos.

A arquitetura da solução utiliza o conceito de empacotamento, ou seja, todo o fluxo de processamento tem como objetivo a geração do Pacote de Informações de Arquivo (AIP), o qual garante o armazenamento de longo prazo.

Nesse sistema, o armazenamento de dados é protegido por meio de software OpenPGP, o qual combina chave pública forte e criptografia simétrica para fornecer serviços de segurança para comunicações eletrônicas e armazenamento de dados [100], sendo apenas adequado para garantir a CoC documental. O Archivematica usa criptografia de chave pública e gera chaves

RSA públicas e privadas, sendo a chave pública usada para criptografar os arquivos armazenados que também são compactados no formato 7-Zip.

A transferência do repositório de confiança para a plataforma de acesso também é supervisionada pelo dispositivo virtual verificador, permitindo uma visão uniforme de toda a CoC documental.

### 5.5.3 Plataforma de Acesso

A função central desta aplicação, mediante solicitação, é entregar e divulgar informações sobre os documentos produzidos e preservados. Na proposta IoTSec2DCCoC, essa função também é modelada para um componente plugável, uma vez que existem vários sistemas que fornecem as funcionalidades necessárias. Portanto, o sistema deve recuperar os pacotes DIP do RDC-Arq e verificar se há alguma discrepância entre o *hash* dos dados armazenados para fins de entrega e a referência de valor correspondente nos pacotes AIP.

Para a validação da proposta, o protótipo PoC é configurado com o sistema Atom [101] como plataforma de acesso, pois tal sistema é capaz de compartilhar com o Archivematica um diretório local ou remoto, permitindo assim a transferência automática de pacotes DIP, os quais são importados usando um *script* PHP Symfony. Conforme necessário, outro *script* PHP Symfony é executado para verificar se existe alguma discrepância entre o *hash* dos dados armazenados no Atom e a referência nos pacotes AIP Archivematica, garantindo, assim, a integridade das informações divulgadas e dos documentos acessados da web.

A Figura 5.14 demonstra a tela de visualização de documentos do Atom utilizado na PoC.

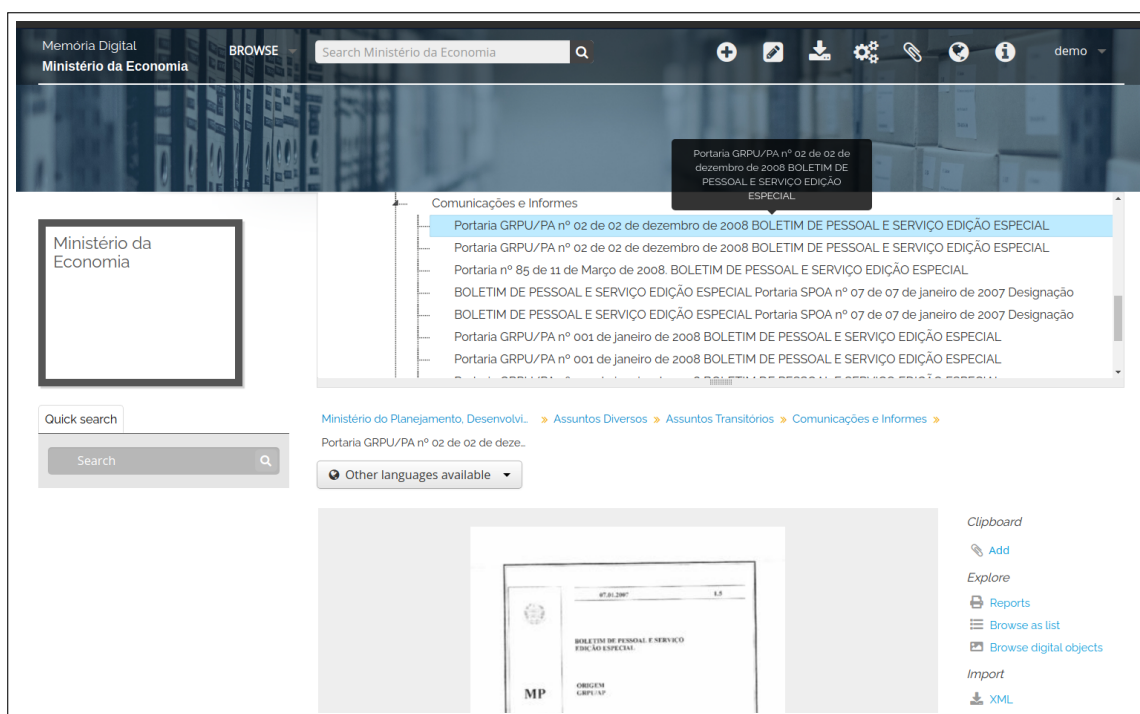


Figura 5.14: Tela de visualização de documentos do Atom implantado na PoC.

### 5.5.4 Componente de Classificação e Extração de Metadados

É o componente que representa a etapa de aplicação de técnicas de inteligência artificial para classificação de documentos digitalizados e extração de metadados sobre o documento. Esse componente atua basicamente como um serviço web dentro do Dispositivo Inteligente, o qual possui uma API REST a ser consumida pelo componente de captura.

Foi desenvolvida uma rede neural para classificar os documentos digitalizados, criada a partir do *framework* chamado Keras, o qual tem como *background* o Framework TensorFlow. Ele foi elaborado para facilitar o desenvolvimento de técnicas que utilizam inteligência artificial.

A Figura 5.15 representa a rede neural desenvolvida para treinar o modelo de classificação dos documentos digitalizados.

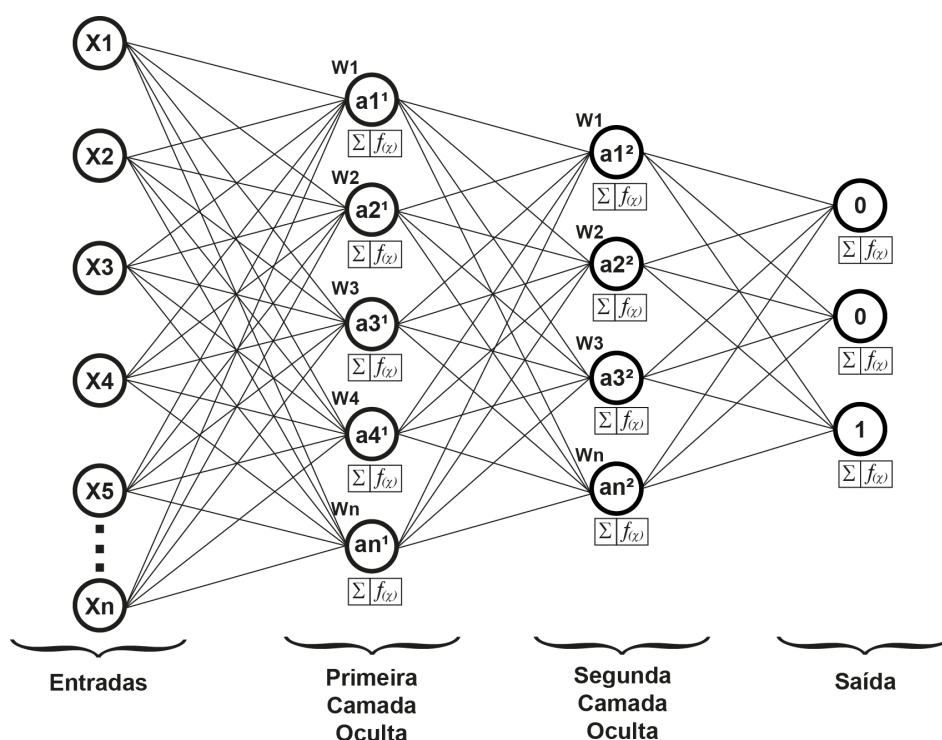


Figura 5.15: Arquitetura da rede neural com 2 camadas.

A rede neural possui um *input* de 40.000 mil palavras que foram inseridas em um vocabulário, o qual representa o cenário em que os documentos fazem parte. Também conta com duas camadas ocultas, sendo que a primeira possui 32 neurônios e a segunda, 16 neurônios. Por fim, a rede neural conta com um *output* de 3 neurônios, que correspondem a quantidade de classes as quais os documentos podem pertencer, representado pela Equação (5.1):

$$sum = \sum_{i=1}^n xi * wi \tag{5.1}$$

A Equação (5.1) é representada pela soma de  $i$  a  $n$ , no qual  $i$  é o primeiro item de entrada e  $n$  o último. A multiplicação entre a entrada e o peso, para cada elemento da rede, é representada pelo elemento  $x_i * w_i$ . Caso a rede neural fosse representada por 3 entradas, a soma da função de ativação seria a seguinte:  $x_1w_1 + x_2w_2 + x_3w_3$ . De acordo com a Figura 4.1, após a soma, a função step é usada para ativar ou desativar o próximo neurônio em uma forma binária, na qual o resultado será sempre 0 ou  $x$ , representado pela Equação (5.2):

$$f(n) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \quad (5.2)$$

O treinamento e o teste do modelo, construído a partir de um processamento que executou 30 épocas, pode ser observado no gráfico da Figura 5.16. Por volta da quarta época, o modelo já se ajusta e demonstra uma acurácia de 99%. Esse percentual de assertividade é considerado ótimo, porém foi identificado que é preciso melhorar a captura do texto digitalizado por meio do processo de OCR, pois, em documentos antigos ou mais desgastados, o OCR extrai caracteres que juntos não representam uma palavra válida para o vocabulário português.

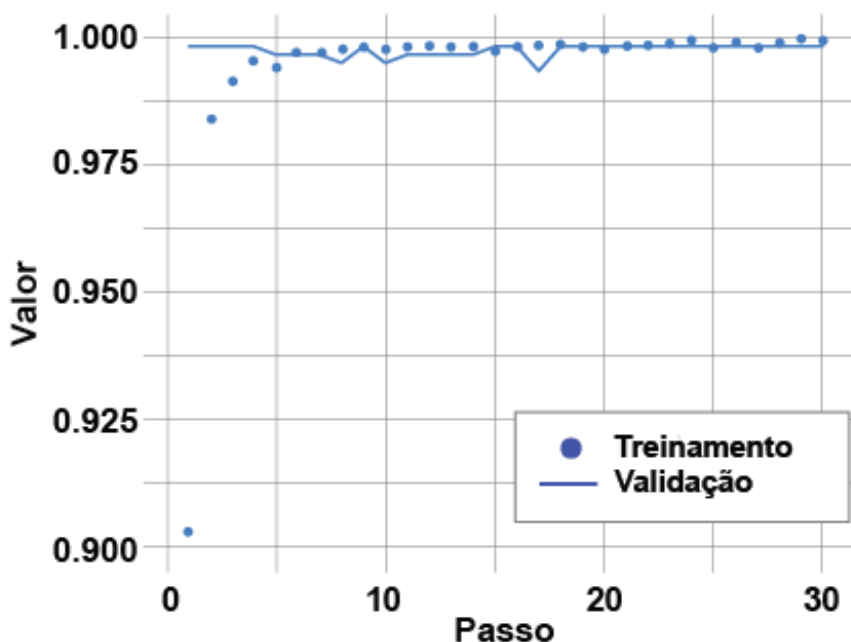


Figura 5.16: Gráfico do resultado de treinamento e validação de acurácia.

Com validação do modelo por meio da matriz de confusão da Figura 4.3, podemos observar que não existe discrepância na classificação dos documentos. Observada a classe “Boletim de Pessoal”, nota-se que dois documentos foram classificados como despacho, sendo que, na verdade, eram do tipo documental Boletim de Pessoal. No entanto, a paridade entre os documentos classificados é bem próxima, validando o modelo como ótimo.

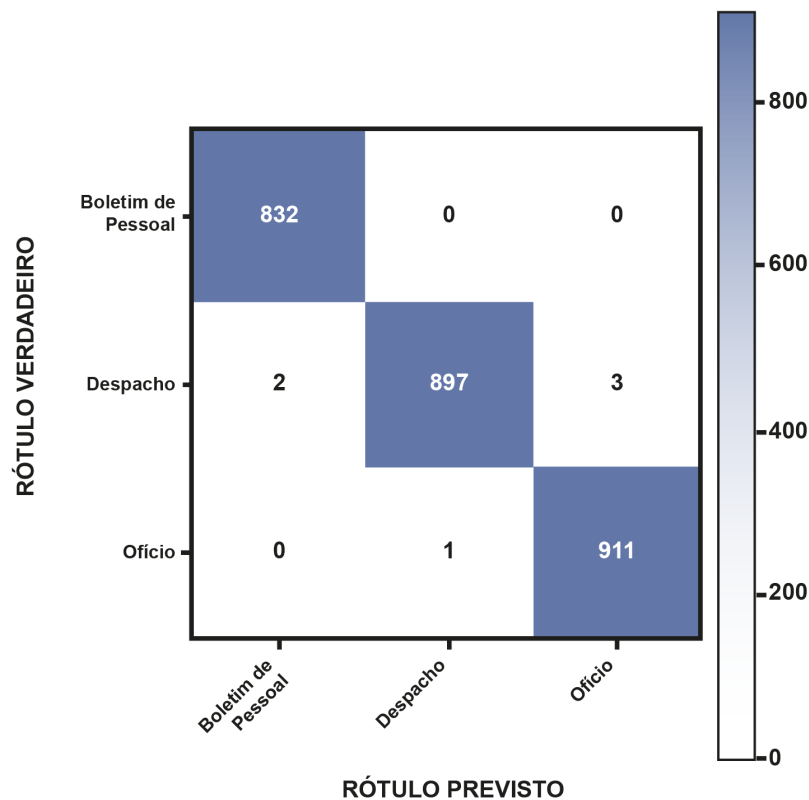


Figura 5.17: Matriz de confusão.

Após a classificação documental utilizando o mesmo arquivo da extensão JSON, são extraídos os seguintes metadados:

- numeroPaginas: quantidade de páginas do documento PDF.
- numeroEdicao: versão do documento no formato x.xx.
- nomeDocumento: nome do documento físico.

Para a extração dos demais metadados do documento classificado, optou-se por utilizar técnicas de visão computacional a partir de imagens para encontrar a região onde está localizada aquela entidade, detectando os quatro pontos cartesianos que representam a localização do termo na imagem.

O modelo utilizado para treinar o detector de objetos é o Yolov5s [82]. Trata-se de um modelo pré-treinado que possui todos os parâmetros pré-definidos e abstraídos pelo próprio *framework*. Além disso, o *framework* carrega um conjunto de pesos pré-definidos que calibram a rede neural convolucional. Seu treinamento é realizado a partir de uma interface chamada Makesense.ai, a qual carrega cada imagem definindo sua *label* e marcando os *bouding boxes* dos objetos a serem detectados. A partir daí, é necessário apenas carregar no modelo Yolov5s [82], os pesos e informar qual o diretório de arquivos de treinamento e os *labels*. Por fim, é realizada a validação do modelo com imagens desconhecidas.



Foram utilizadas 100 imagens para treinamento com os seus respectivos *bouding boxes*, conforme o exibido na Figura 5.18

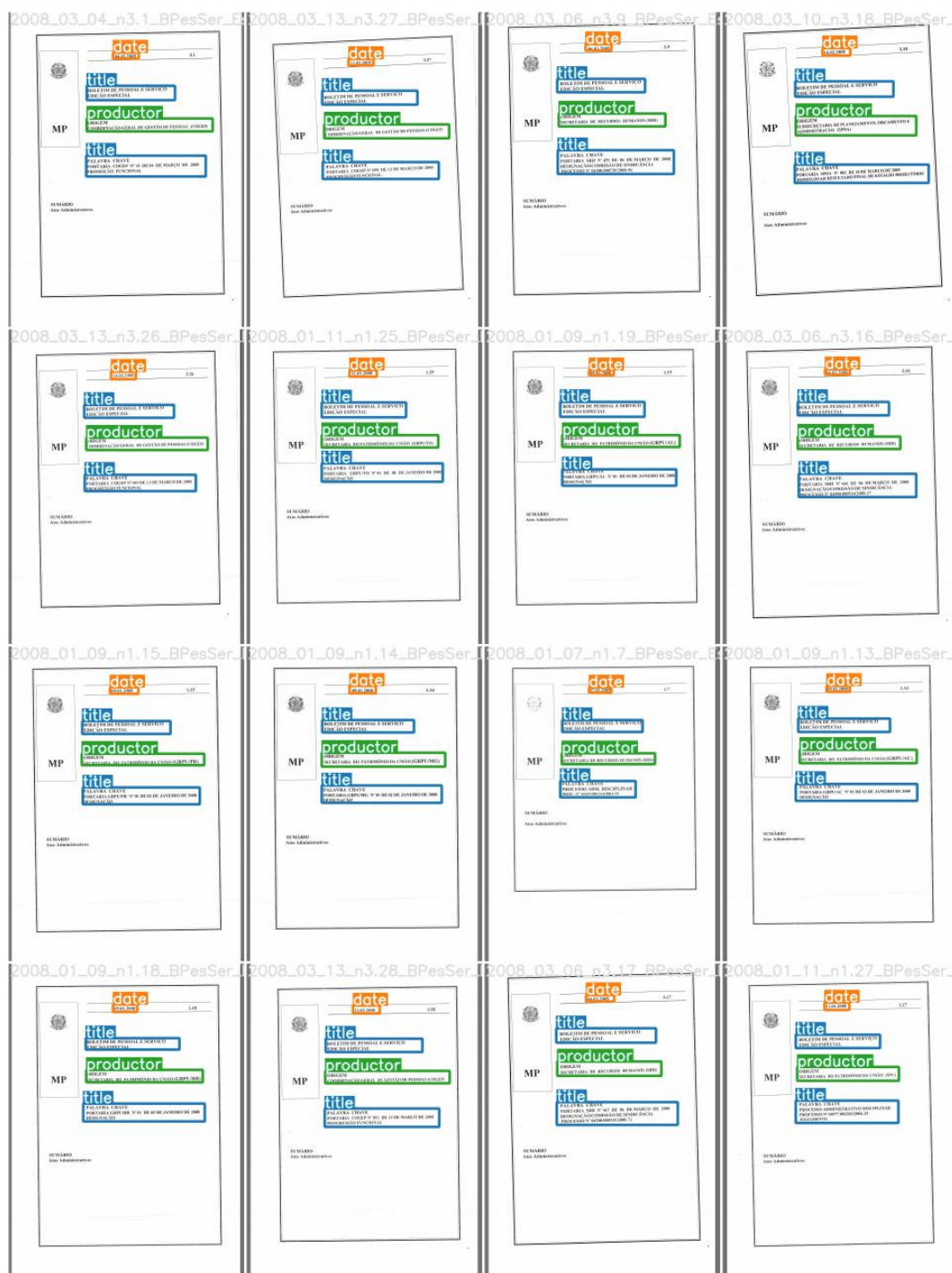


Figura 5.18: Marcação de metadados para treinamento.

Para a marcação dos *bouding boxes*, foi utilizada a ferramenta Makesense.ai [102], delimitando as 3 entidades conforme o exibido na Figura 5.18. Em seguida é gerado um arquivo TXT para cada imagem, contendo o número e os 4 pontos cartesianos respectivos para cada entidade.

O modelo é baseado na arquitetura YoloV5 da empresa Ultralytic, que, por sua vez, é baseada

na Yolov3 (you only look once) [82]. Trata-se de uma rede de detecção de objetos em tempo real que aplica uma rede neural, divide a imagem em regiões e prediz *bouding boxes* e probabilidades para cada região.

Para tal, foram utilizadas 100 imagens para treinamento com os seus respectivos *bouding boxes* e mais 20 imagens para validação durante 300 épocas. O treinamento resultou em uma média de precisão mAP de 0.71@0.95. O resultado pode ser observado na Figura 5.19.

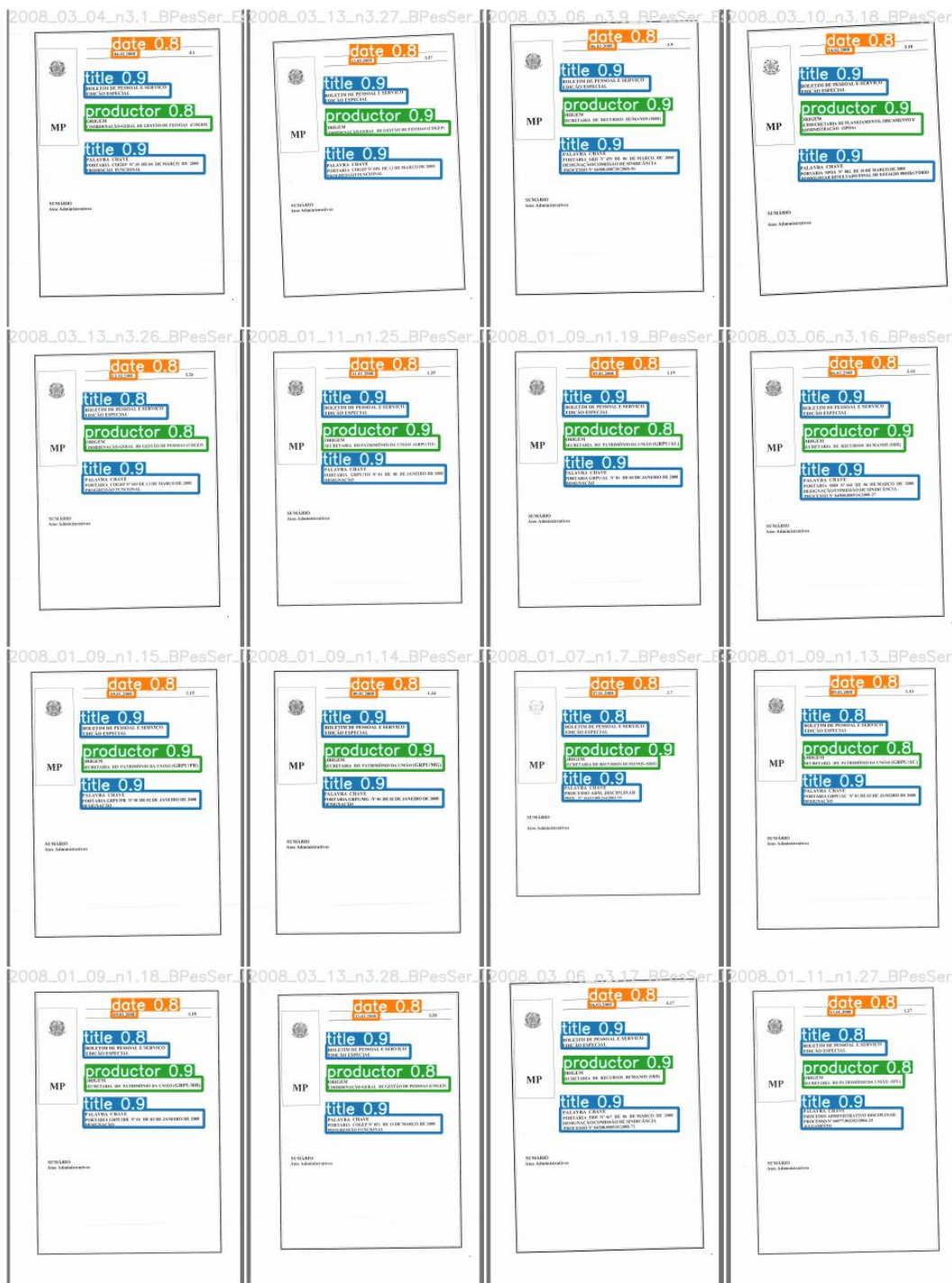


Figura 5.19: Modelo preditivo treinado.

Dessa forma, independente se o termo mudou a formatação, quantidade de espaços, ou localização na imagem, o detector de objetos irá encontrar o que foi treinado para encontrar. Uma vez que a região da entidade é localizada, computa-se o OCR específico daquela localização e obtém-se o valor textual referente à entidade em questão, sendo possível, então, conforme o modelo treinado, extrair os seguintes metadados:

- *data*: a data de publicação do documento;
- *nomeArquivo*: nome do documento concatenado com a versão e a data de publicação do documento;
- *nomeProdutores*: os órgãos contidos no documento.

Para processo de extração, são executados os seguintes etapas:

- Serviço de carga de imagem: as imagens dos documentos são carregadas, e a técnica de detecção de objetos é aplicada por meio de uma chamada de função.
- Conversão de imagem: as imagens são convertidas para *bits*.
- Inicialização do modelo: nessa etapa, são inicializados os pesos do modelo treinado.
- Inferências: inferência nas imagens de documentos recebidas.
- Carregamento de dicionários: são carregadas as listas de dicionários com *label* e sua localização.
- Serviço de saída de dados: o *webService* retorna os dados extraídos.

## 6 VALIDAÇÃO DA RESILIÊNCIA DO MODELO

Para a validação da proposta, foi necessário verificar a resiliência dos serviços de segurança para garantir a CoC documental. A partir do protótipo PoC funcional, buscou-se evidências da resiliência da arquitetura proposta a ameaças de segurança que pudessem comprometer a autenticidade e integridade de documentos arquivísticos digitais e metadados associados que fluem através da linha de produção documental, do scanner ao armazenamento no RDC-Arq e além, por meio da transferência para a Plataforma de Acesso.

### 6.1 RESILIÊNCIA A ATAQUES DE FALSIFICAÇÃO

Esta técnica de ataque tem o objetivo de falsificar o tráfego, por exemplo, copiando ou forjando pacotes para injetar tráfego malicioso.

Para testar a resiliência a ataques de falsificação, foram realizados ataques que forjaram as requisições HTTP, simulando um intruso na rede, principalmente no ponto mais vulnerável, o Gateway de API que opera a comunicação entre Dispositivos Inteligentes e as demais partes da estrutura.

As contramedidas do IoTSec2DCoC para impedir os invasores que realizam ataques de negação de serviço baseados em falsificação são impor várias barreiras as quais devem ser superadas simultaneamente, obrigando o invasor a gastar muitos recursos e tempo, o suficiente para ser detectado antes de conseguir.

As tentativas de falsificar o tráfego, por exemplo, copiando ou forjando pacotes para injetar tráfego malicioso, para serem bem-sucedidas, exigirão que o invasor execute ações ao mesmo tempo e de forma síncrona a fim de sintetizar uma assinatura válida no esquema de Serviços de Validação de Documentos e Assinatura Digital (Check-out e Check-in), no Serviço de Autoridade Certificadora CFSSL) e no Serviço de Autenticação (JWT).

Os ataques e as contramedidas aplicadas serão detalhados nos subtópicos: 6.1.1; 6.1.2; 6.1.3 e 6.1.4.

#### 6.1.1 Tentativa de Interceptar o Tráfego Entre os Dispositivos Inteligentes e o Middleware Edge

Começando com os Dispositivos Inteligentes, todos os componentes que oferecem serviços baseados em HTTP têm um certificado digital válido e configurado, solicitado ao CFSSL, definindo, então, esse certificado para o fornecimento de seus serviços via HTTP seguro sobre canais TLS (HTTPS). Isso é usado para autenticar a fonte original do fluxo e para garantir a confiden-

cialidade das transferências, uma vez que o cliente do serviço criptografa as mensagens com a chave pública do serviço. O protótipo da PoC usa o Wireshark [103] para verificar se um agente malicioso poderia interceptar as comunicações. Conforme mostrado na Figura 6.1, a ferramenta não consegue decifrar as mensagens, resultado já esperado, considerando a maturidade do TLS e a certificação das chaves públicas oferecida pela AC.

```

▶ Transmission Control Protocol, Src Port: 443, Dst Port: 47110, Seq: 1, Ack: 2, Len: 56
▼ Secure Sockets Layer
  ▼ TLSv1.2 Record Layer: Application Data Protocol: http-over-tls
    Content Type: Application Data (23)
    Version: TLS 1.2 (0x0303)
    Length: 51
    Encrypted Application Data: ef6ba57157ed01c9c83d23a9d31eca3afa12a63e121f7a50...

```

0000	40 a3 cc 01 07 b1 b4 2a 0e d1 4e f4 86 dd 60 0d	@.....* ..N....
0010	1d 39 00 58 06 74 28 00 03 f0 40 04 08 01 00 00	·9·X·t(· ..@.....
0020	00 00 00 00 20 0e 28 04 01 4c 65 e4 43 aa 19 8f	.....(· ..Le·C··
0030	b8 65 58 5e 8c 7d 01 bb b8 06 31 53 71 53 51 e4	·eX^·}·· ·1SqSQ·
0040	a9 f7 80 18 01 0e df 47 00 00 01 01 08 0a 56 d4	.....G .....V·
0050	fe de a3 2e a8 50 17 03 03 00 33 ef 6b a5 71 57	.....P·· ·3·k·qW
0060	ed 01 c9 c8 3d 23 a9 d3 1e ca 3a fa 12 a6 3e 12	.....=#·· ·:·>·
0070	1f 7a 50 48 93 84 c1 be 91 ff 4f c0 26 0f 1e 80	·zPH..... ·0·&··
0080	47 29 b2 5a c4 e7 e1 d4 77 84 59 e8 ae ce	G)·Z···· w·Y··

Figura 6.1: Intercepção malsucedida do conteúdo da mensagem.

### 6.1.2 Tentativa de Trocar Dados com o Middleware Edge sem Autenticação Válida

Os componentes IoT, Middleware Edge e Middleware Cloud recebem pacotes de documentos provenientes dos Dispositivos Inteligentes e do Serviço de Admissão de Arquivos Nato Digitais, exigindo que essas fontes apresentem um *token* JWT válido. O JWT é considerado inválido se estiver malformado ou se sua assinatura digital HMAC for inválida de acordo com a chave simétrica configurada.

A Figura 6.2 demonstra, em seu primeiro caso, um usuário legítimo (um Dispositivo Inteligente devidamente cadastrado) realizando uma chamada ao middleware e informando um JWT válido. No caso seguinte, demonstra um adversário tentando o envio de dados sem nenhuma autenticação, recebendo como retorno um erro HTTP 401 (*Unauthorized*). Já no próximo caso, o adversário informa um JWT, mas ele está malformado ou possui a assinatura HMAC inválida (não foi autenticado pela solução), recebendo como retorno um erro HTTP 403 (*Forbidden*).

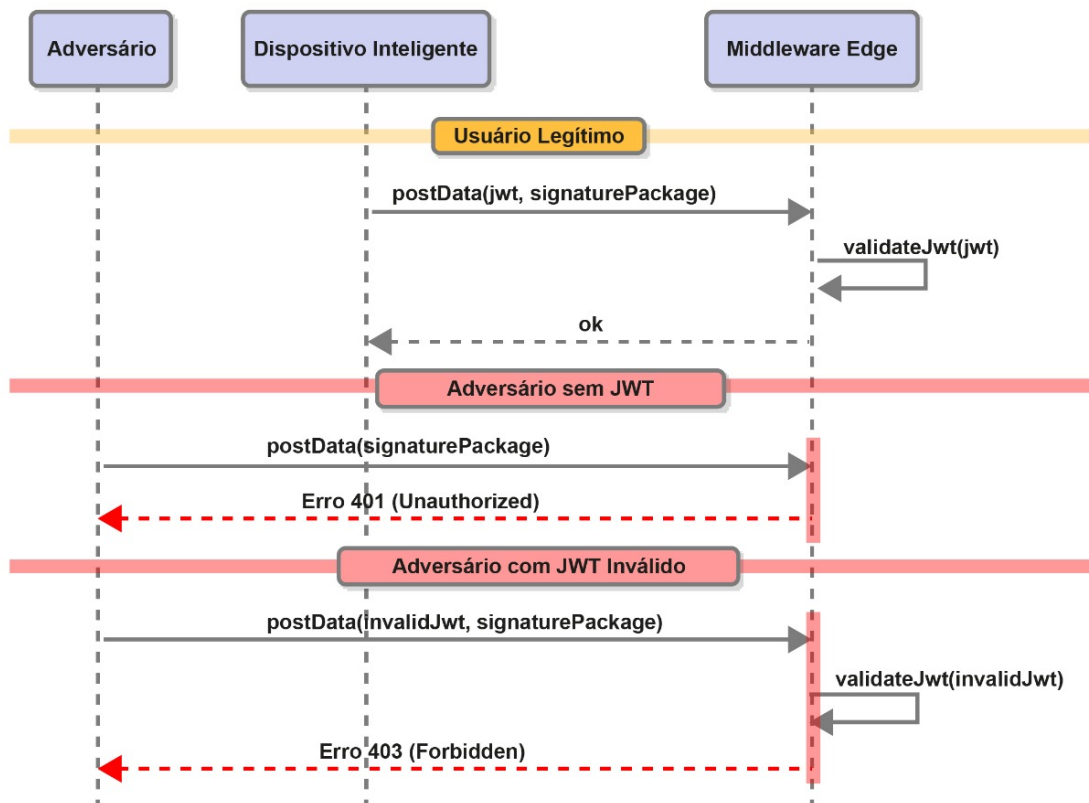


Figura 6.2: Diagrama de sequência com exemplificação do ataque sem autenticação válida.

O trecho de Código 6.1 do Middleware Edge realiza a validação do JWT em todas as requisições. Como se trata de uma aplicação Flask, é possível registrar uma função de *callback*, executada anteriormente a todas as requisições HTTP feitas ao middleware (`validate_auth_jwt`). Logo em seguida, é definida essa função, a qual, inicialmente, valida se a requisição não é interna (vinda do endereço 127.0.0.1), já que essas não devem ser validadas por serem requisições executadas pelos componentes internos do middleware. Logo após, ela faz a validação do JWT e informa o erro de cliente HTTP 401 (*Unauthorized*), em caso do JWT não ter sido informado, ou erro de cliente HTTP 403 (*Forbidden*), caso o JWT seja considerado inválido, ou seja, esteja malformado ou sua assinatura digital HMAC esteja inválida de acordo com a chave simétrica configurada (`JWT_KEY`).

Código 6.1: Validação do JWT no Middleware Edge.

```
# flask configuration
self.app = Flask(__name__)
self.app.before_request(validate_auth_jwt)
# ...
def validate_auth_jwt():
    if request.environ.get('HTTP_X_FORWARDED_FOR') is None:
        remote_addr = request.environ['REMOTE_ADDR']
```

```

else: remote_addr = request.environ['HTTP_X_FORWARDED_FOR']
if remote_addr == '127.0.0.1':
    return
elif 'Authorization' in request.headers and 'Bearer' in
request.headers['Authorization']:
    authorization = request.headers['Authorization'].
    replace('Bearer ', '')
    try:jwt.decode(authorization, JWT_KEY, algorithms=
        ['HS512'])
    except Exception as e:
        msg = 'Authorization problem. Check your
        Authorization Token'
        logger.exception(msg)
        abort(403, {'message': msg, 'details': ''})
else: msg = 'No Authorization Token was provided, access
    denied'
    logger.warn(msg)
    abort(401, {'message': msg, 'details': ''})

```

O protótipo da PoC permitiu simular a tentativa do invasor de forjar autenticação por meio do envio de solicitações HTTP com JWT inválido, conforme o Código 6.2.

Código 6.2: Requisição HTTP com JWT inválido no cabeçalho.

```

POST /d/data HTTP/1.1
User-Agent: Mozilla/4.0 (compatible; MSIE5.01; Windows NT)
Host: ec2-18-228-152-200.sa-east-1.compute.amazonaws.com: 8080
Authorization: Bearer eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9 \
.eyJzdWIiOiIxMjMONTY3ODkwIn0 \
.dozjgNryP4J3jVmNH10w5N_XgL0n3I9P1FUP0THsR8U
Content-Type: application/json
{# Correct body with middleware data}

```

Como resultado relatado, todas essas tentativas do invasor de forjar autenticação foram detectadas, e a transação recusada pelo Middleware Edge, conforme mostrado no print da tela com os comandos na Figura 6.3.

```

2020-03-26 15:15:19,960 - ERROR <PID 8516:MainProcess> gateway.listeners.http.validate_auth_jwt(): Authorization problem. Check your
Authorization Token
Traceback (most recent call last):
  File "/home/ubuntu/uiot/gateway/gateway/listeners/http.py", line 119, in validate_auth_jwt
    jwt.decode(authorization, JWT_KEY, algorithms=['HS512'])
  File "/home/ubuntu/uiot/gateway/venv/lib/python3.7/site-packages/jwt/api_jwt.py", line 92, in decode
    jwt, key=key, algorithms=algorithms, options=options, **kwargs
  File "/home/ubuntu/uiot/gateway/venv/lib/python3.7/site-packages/jwt/api_jws.py", line 156, in decode
    key, algorithms)
  File "/home/ubuntu/uiot/gateway/venv/lib/python3.7/site-packages/jwt/api_jws.py", line 223, in _verify_signature
    raise InvalidSignatureError('Signature verification failed')
jwt.exceptions.InvalidSignatureError: Signature verification failed
179.187.102.209 - - [26/Mar/2020 15:15:19] "GET / HTTP/1.1" 403 -

```

Figura 6.3: Falha de validação do JWT no Middleware Edge.



### 6.1.3 Tentativa de Enviar Pacotes de Documentos com Assinaturas Inválidas para o Middleware Edge

A verificação sistemática das transferências é apontada como um pilar da proposta, uma vez que garante a integridade e autenticidade dos documentos em todo o fluxo da CoC documental. Essa tarefa é executada pelos pares de dispositivos virtuais Check-out e Check-in. O componente Check-out realiza a assinatura digital quando um documento é alterado ou quando um novo documento é gerado, antes de transferir o documento para a próxima entidade na cadeia. Por outro lado, o componente Check-in valida a assinatura digital mediante o recebimento de documentos por cada componente. A Figura 6.4 demonstra, em seu primeiro caso, como um usuário legítimo (Dispositivo Inteligente devidamente cadastrado) realiza o envio de dados para o middleware, com uma autenticação JWT válida e uma assinatura digital do documento válida. No caso seguinte, é demonstrado um adversário que, de alguma forma, conseguiu se autenticar na solução e possui um JWT válido, mas envia o documento com uma assinatura digital inválida, resultando em um erro HTTP 400 (*Bad Request*) e no cancelamento do envio do documento para a próxima entidade da cadeia.

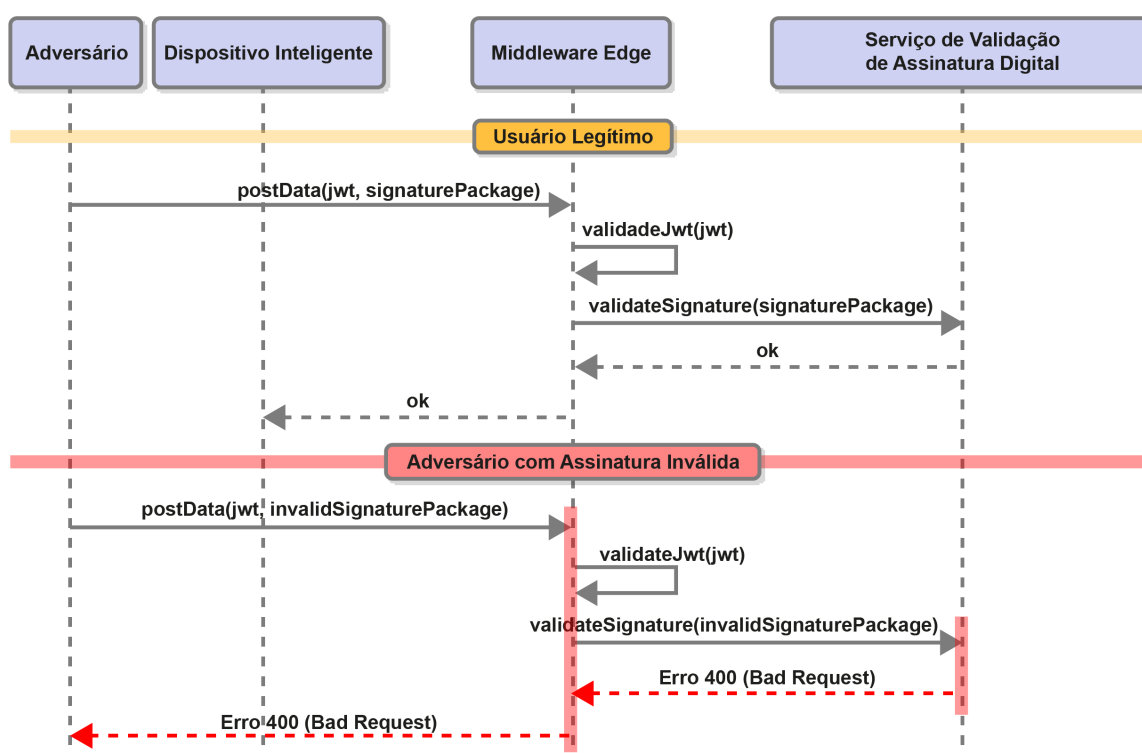


Figura 6.4: Diagrama de sequência com exemplificação do ataque sem assinatura digital de documentos válida.

Se um adversário tentar alterar um documento que está sendo transferido ou tentar injetar um documento ilegítimo no fluxo, esse documento forjado será rejeitado, porque contém uma assinatura digital inválida. O componente Check-in espera receber do componente Check-out anterior um arquivo compactado assinado digitalmente com todos os dados e metadados do documento.



A assinatura digital é feita no *hash* SHA512 do arquivo, conforme o trecho de código 6.3. Essa assinatura é verificada pelo componente Check-in usando a chave pública vinculada ao certificado digital X.509 do componente Check-out anterior, conforme o trecho de Código 6.4.

Código 6.3: Realização da assinatura digital no Componente de check-out.

```
def _sign(self, doc_package, priv_key_pem):
    private_key = serialization.load_pem_private_key(
        priv_key_pem,
        password=PRIV_KEY_PASSWD,
        backend=default_backend()
    )
    signature = private_key.sign(
        doc_package,
        ec.ECDSA(hashes.SHA512())
    )
    return signature
```

Código 6.4: Validação da assinatura digital chamado pelo Componente de check-in.

```
def _validate_signature(self, doc_package, signature, cert_pem):
    cert = x509.load_pem_x509_certificate(
        cert_pem,
        default_backend()
    )
    public_key = cert.public_key()
    public_key.verify(
        signature,
        doc_package,
        ec.ECDSA(hashes.SHA512())
    )
```

Para que um agente malicioso envie com êxito um documento ilegítimo para o Middleware Edge, esse invasor deve se passar por um Dispositivo Inteligente, uma tarefa que exigiria que ele se registrasse previamente no CFSSL. Portanto, o invasor deve contornar pelo menos dois serviços de segurança de forma síncrona, uma tarefa considerada difícil.

Para efeito de validação, os documentos falsos foram inseridos no fluxo do protótipo da PoC, como na requisição HTTP do Código 6.4.

Código 6.5: Requisição HTTP com assinatura digital inválida.

```
POST /d/data HTTP/1.1
User-Agent: Mozilla/4.0 (compatible; MSIE5.01; Windows NT)
Host: ec2-18-228-152-200.sa-east-1.compute.amazonaws.com:8080
Content-Type: application/json
{"clientTime": 1585220368,
 "tags": [
  "IoTSec2DCoC",
  "smart-device"
 ],
 "sensitive": 1,
 "chipset": "Broadcom BCM2837B0",
 "mac": "40:A3:CC:01:07:B1",
 "serviceNumber": 1,
 "value": ["<ZIP_AND_INVALID_SIGNATURE_BASE64>"]
}
```

Mesmo simulando uma assinatura digital, os documentos foram recusados pelo componente Check-in, conforme mostrado nos comandos da tela capturada na Figura 6.5.

```
26/03/2020 10:59:29 - app - INFO: Received signature_package
26/03/2020 10:59:29 - app - INFO: signature_package id=5ac18e2c-5047-4ff6-81ee-412e487e8309 data extracted with success
26/03/2020 10:59:29 - app - ERROR: doc_package id=5ac18e2c-5047-4ff6-81ee-412e487e8309 signature validation error:
Traceback (most recent call last):
  File "/home/ubuntu/iotsec2coc-checkout/iotsec2coc_checkout/app.py", line 104, in exec_validate_signature
    self._validate_signature(doc_package + b'aaa', signature, cert_pem)
  File "/home/ubuntu/iotsec2coc-checkout/iotsec2coc_checkout/app.py", line 55, in _validate_signature
    ec.ECDSA(hashes.SHA512())
  File "/home/matheus/.virtualenvs/iotsec2coc-checkout/lib/python3.6/site-packages/cryptography/hazmat/backends/openssl/ec.py", line 352, in
verify
    _ecdsa_sig.verify(self._backend, self._signature, data)
  File "/home/ubuntu/.virtualenvs/iotsec2coc-checkout/lib/python3.6/site-packages/cryptography/hazmat/backends/openssl/ec.py", line 101, in
_ecdsa_sig.verify
    raise InvalidSignature
cryptography.exceptions.InvalidSignature
26/03/2020 10:59:29 - _internal - INFO: 172.31.16.174 - - [26/Mar/2020 10:59:29] "POST /api/v1.0/signature/validate HTTP/1.1" 400
```

Figura 6.5: Falha de validação de assinatura digital solicitado pelo Componente de Check-in.

Além da validação da assinatura digital, o certificado digital permite solicitar ao CFSSL a validação de toda a cadeia de confiança desse certificado. Essa validação da cadeia de confiança também é executada por solicitação do componente Check-in, conforme o trecho de código 6.6.

Código 6.6: Validação do certificado digital do assinante chamado pelo Componente de check-in.

```
def _validate_chain_of_trust(self, cert_pem, trusted_certs_pem)
: certificate = crypto.load_certificate(crypto.FILETYPE_PEM,
cert_pem)
    store = crypto.X509Store()
    for trusted_cert_pem in trusted_certs_pem:
        trusted_cert = crypto.load_certificate(crypto.FILETYPE
_PEM, trusted_cert_pem)
        store.add_cert(trusted_cert)
    store_ctx = crypto.X509StoreContext(store, certificate)
    store_ctx.verify_certificate()
```

Se o certificado digital utilizado para assinar o arquivo compactado não for confiável, ou seja, não for emitido pela raiz ou por uma autoridade de certificação intermediária, é detectada uma recusa de validação, conforme mostrado nos comandos da tela capturada na Figura 6.6.

```
26/03/2020 11:41:51 - app - ERROR: doc_package_id=5ac18e2c-5047-4ff6-81ee-412e487e8309 certificate chain of trust validation error:
Traceback (most recent call last):
  File "/home/ubuntu/iotsec2coc-checkout/iotsec2coc_checkout/app.py", line 84, in exec_validate_cert_chain
    self._validate_chain_of_trust(cert_pem, trusted_certs_pem)
  File "/home/ubuntu/iotsec2coc-checkout/iotsec2coc_checkout/app.py", line 43, in _validate_chain_of_trust
    store_ctx.verify_certificate()
  File "/home/ubuntu/.virtualenvs/iotsec2coc-checkout/lib/python3.6/site-packages/OpenSSL/crypto.py", line 1766, in verify_certificate
    raise self.exception_from_context()
OpenSSL.crypto.X509StoreContextError: [20, 0, 'unable to get local issuer certificate']
26/03/2020 11:41:51 - _internal - INFO: 172.31.16.174 - - [26/Mar/2020 11:41:51] "POST /api/v1.0/signature/validate HTTP/1.1" 400
```

Figura 6.6: Falha de validação do certificado digital solicitado pelo Componente de Check-in.

### 6.1.4 Contramedidas Contra Vazamento de Credenciais e Chaves de Criptografia

É importante, ao usar autenticação JWT, assinaturas digitais e TLS, levar em consideração a vulnerabilidade relacionada ao possível vazamento de credenciais (no caso de autenticação) e chaves privadas (em assinaturas e túneis seguros), pois esse risco surge com descuidos dos portadores de credenciais ou situações de engenharia social. Um ataque a partir de uma credencial de Dispositivo Inteligente vazada permite que o invasor se autentique na solução e tente o envio de dados ilegítimos ao Middleware Edge.

A primeira contramedida nesse caso se dá por meio da validação da assinatura digital, pois o invasor não conseguirá enviar arquivos devido à sua assinatura digital inválida (já que ele não possuirá um certificado digital e chave privada emitidos pelo CFSSL). A segunda contramedida acionada por esse evento de segurança é a identificação do vazamento e o consequente bloqueio da conta do dispositivo no Serviço de Autenticação. Após uma investigação forense do caso, novas credenciais podem ser geradas e configuradas para o Dispositivo Inteligente afetado.

Em uma nova tentativa de ataque sem um JWT válido, na qual o invasor tente ler conversas, assinar documentos, tentando enviá-los para o Middleware Edge, ele não teria sucesso, pois não seria autenticado com um JWT válido.

Como contra medida nesse caso, a detecção desse vazamento crítico acionará a revogação do certificado digital correspondente pelo CFSSL, informado na Lista de Revogação de Certificados - CRL, a qual permite ao próprio dispositivo identificar que o certificado é inválido, solicitando um novo.

## 6.2 RESILIÊNCIA DO MODELO A ATAQUES CONTRA SENHAS

Um dos mais comuns ataques contra a segurança de sistemas é o ataque de força bruta, no qual a estratégia do atacante consiste em testar combinações de chaves criptográficas ou de senhas de forma exaustiva para quebrá-las. Uma das formas de executar um ataque de força bruta contra senhas é fazendo uso de um dicionário de senhas (uma lista de senhas comuns ou senhas baseadas

em engenharia social contra um indivíduo), o que pode ser eficaz em casos de senhas fracas.

Na solução, o Serviço de Autenticação pode ser um alvo de um ataque de força bruta com dicionários. A contramedida implementada utiliza uma estratégia de mitigação desse ataque por meio de uma política de bloqueio de conta, um método feito para se evitar a pressuposição de senhas, bloqueando a conta da identidade (pessoa ou dispositivo) que possua múltiplas tentativas de autenticação falhas. A configuração da política de bloqueio de conta é feita no Serviço de Autenticação da forma a seguir:

- Bloqueio de conta habilitado (e.g. Sim ou não).
- Duração do bloqueio de conta em segundos (e.g. 3600s ou 0s para somente ser desbloqueado com ação do administrador).
- Quantidade máxima de tentativas falhas por identidade antes do bloqueio (e.g. 3).
- Quantidade de tempo máximo entre as tentativas falhas de autenticação antes do bloqueio da senha (e.g. 60s).

A Figura 6.7 ilustra que, após 3 tentativas falhas de autenticação dentro do tempo máximo entre tentativas falhas configuradas (60s), o serviço indica que a conta foi desabilitada (código HTTP 422, Entidade não processável).

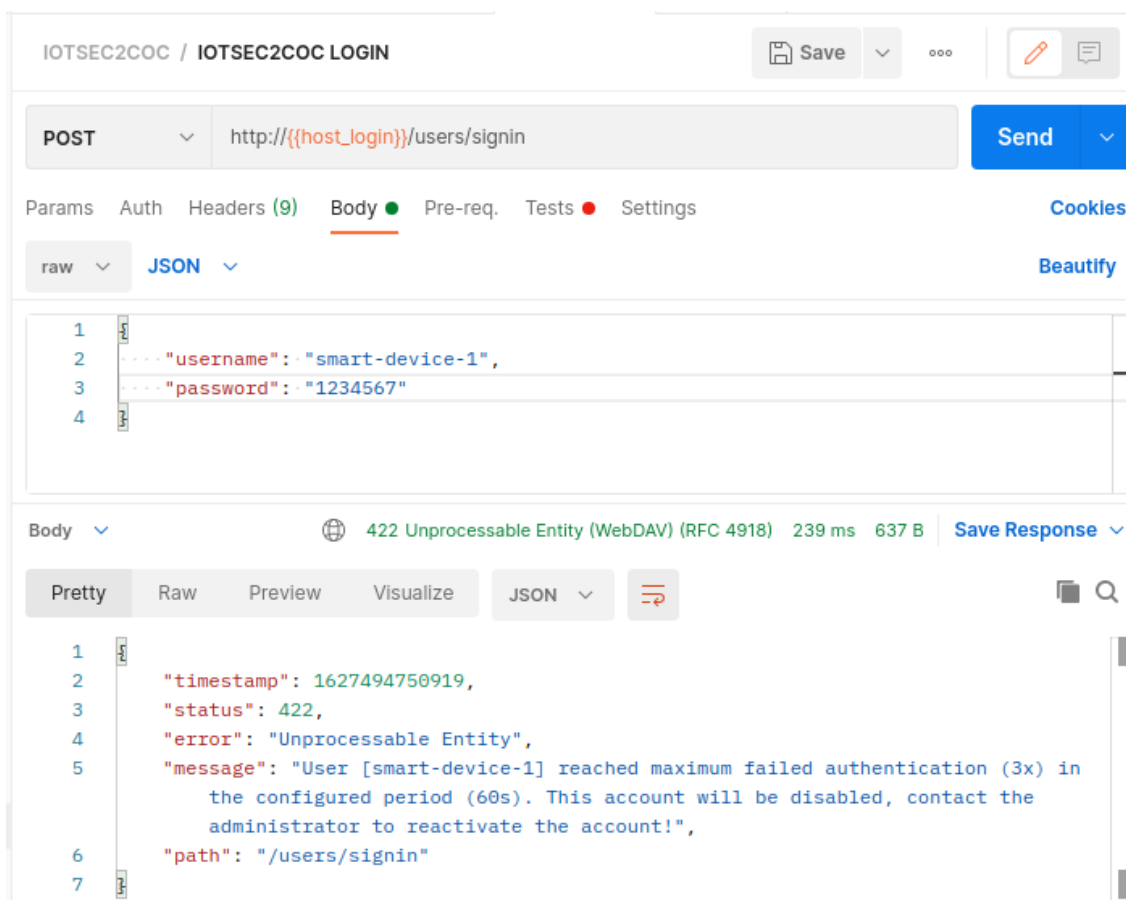


Figura 6.7: Número máximo configurado de autenticações falhas atingido e desabilitação da conta.

Por fim, na Figura 6.8, é possível observar que, em uma próxima tentativa de autenticação, o usuário recebe uma mensagem indicando que a sua conta está desabilitada (código HTTP 423, Trancado) e que ele deve entrar em contato com o administrador, pois no caso a duração do bloqueio foi configurado como 0s, ou seja, por tempo indeterminado.

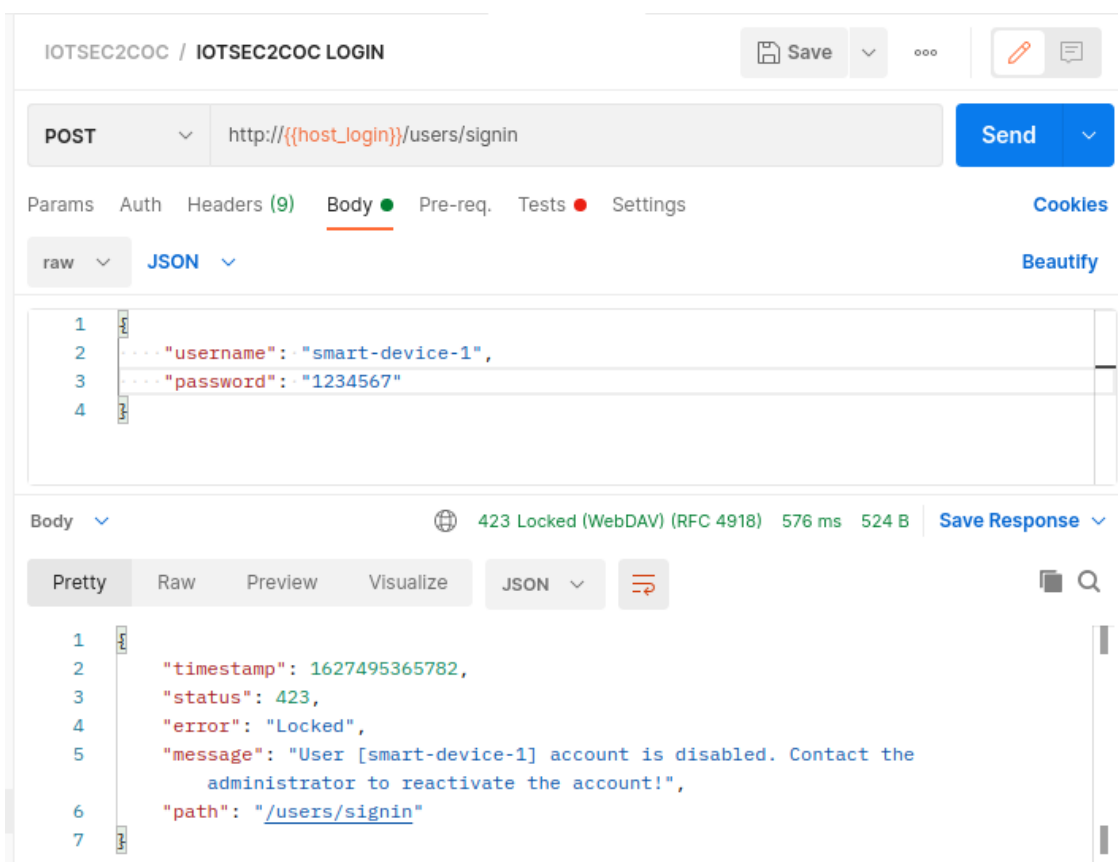


Figura 6.8: Tentativa de autenticação com a conta desabilitada.

### 6.3 CAPACIDADE EFETIVA DE MONITORAMENTO E RASTREABILIDADE

A solução proposta disponibiliza uma forma de acessar uma trilha de auditoria sobre o processamento de um determinado documento digital. Durante todo o fluxo da cadeia de custódia os eventos são registrados em logs específicos para esse objetivo, chamados de *checkpoints*. Por meio desses *checkpoints*, pode-se rastrear o momento em que um determinado evento ocorreu com o documento, em qual etapa do fluxo (componente) ocorreu e sobre qual documento se trata. Essa auditoria pode ser feita de forma centralizada e de fácil acesso e análise, graças a solução de logs centralizados. Tais registros também facilitam descobertas de anomalias e de possíveis ataques, já que o principal registro de log tem relação com a realização e validação de assinaturas digitais em cada etapa do processo, onde, em caso de falha, ela poderá ser facilmente indentificada. A Figura 6.9 possui um exemplo de trilha de log de processamento de um documento específico, podendo ser visualizado por meio do Serviço de Monitoramento com o uso da ferramenta Kibana.

Discover / TRILHA DE AUDITORIA (CHECKPOINT) ✓ Options New Save Open Share Inspect

message:checkpoint KQL Last 30 days Show dates Refresh

39 hits Reset search Show chart

>	Jun 9, 2021 @ 09:31:44.205	[CHECKOUT] [CHECKPOINT] Signed Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=uiot-edge-1
>	Jun 9, 2021 @ 09:31:38.756	[CHECKIN] [CHECKPOINT] Validated Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=uiot-edge-1
>	Jun 9, 2021 @ 09:31:36.753	[PKI] [CHECKPOINT] Got Private Key and Certificate from PKI for Module=smart-device-1
>	Jun 9, 2021 @ 09:31:36.108	[CHECKOUT] [CHECKPOINT] Signed Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=smart-device-1
>	Jun 9, 2021 @ 09:31:35.998	[LOGIN] [CHECKPOINT] Successfully logged-in Module=smart-device-1
>	Jun 9, 2021 @ 09:31:34.751	[SMART DEVICE] [CHECKPOINT] IA processing finished finished Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=smart-device-1
>	Jun 9, 2021 @ 09:31:19.964	[SMART DEVICE] [CHECKPOINT] OCR finished Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=smart-device-1
>	Jun 9, 2021 @ 09:31:09.572	[SMART DEVICE] [CHECKPOINT] Scanning finished Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=smart-device-1
>	Jun 9, 2021 @ 09:31:05.482	[SMART DEVICE] [CHECKPOINT] Starting processing Package_ID=2ab7ecdd-b8f9-4a16-883b-38ef56ca9104 Module=smart-device-1

Figura 6.9: Log de auditoria consultada na interface web Kibana.

## 7 DISCUSSÃO DOS RESULTADOS

Este capítulo apresenta e discute os principais resultados desta tese, os quais foram extraídos a partir da prova de conceito e validação da solução IoTSec2DCoC.

### 7.1 ASPECTOS DO MODELO E SEU RELACIONAMENTO COM OS COMPONENTES DA SOLUÇÃO

Como o IoTSec2DCoC propõe uma medida de segurança totalmente distribuída, a Tabela 7.1 ilustra a contribuição de cada componente para responder aos requisitos de segurança e integração, e seus detalhes operacionais são descritos nas seções a seguir.

Tabela 7.1: Contribuições de cada componente para os serviços de segurança do IoTSec2DCoC.

Parâmetros	Componentes IoTSec2DCoC										
	DI	API	MID	EMP	AC	AUT	ASS	CB	RS	SR	SM
<b>Reforçando propriedades de segurança</b>											
Confidencialidade					✓						
Integridade					✓		✓				
Autenticidade	✓	✓			✓	✓	✓				
<b>Combatendo ameaças à segurança</b>											
Perturbação de operação normal									✓		
Ataque de sniffing					✓						
Ataque de falsificação de assinatura digital					✓		✓				
Ataque de falsificação de autenticação							✓				
Ataque de autenticação de força bruta							✓				
<b>Reforçando as propriedades da solução e de integração</b>											
Monitoramento e rastreabilidade										✓	✓
Integração de componentes				✓					✓		
Escaneamento de documentos	✓				✓						
Classificação e extração de metadados	✓										

**Legenda:** DI - Dispositivo Inteligente; API - Gateway de API; MID - Middleware IoT; EMP – Empacotador; AC - Autoridade Certificadora; AUT - Serviço de Autenticação; ASS - Assinatura Digital de Check-out e Check-in; CB - Circuit Breaker, RS - Registro de Serviços; SR - Serviço de Rastreamento; SM - Serviço de Monitoramento.

#### 7.1.1 Componentes e o reforço às propriedades de segurança da informação

Com as propriedades da arquitetura IoTSec2DCoC destinadas a incrementar a segurança, foi possível verificar, nos resultados da prova de conceito, que os componentes da solução apontados na matriz da Tabela 7.1 demonstraram que o componente GC promove a **confidencialidade**, uma vez que emite certificados digitais (padrão X.509) e chaves privadas que permitem configurar o uso do TLS (v1.2 ou 1.3), garantindo um túnel seguro de comunicação.

Essas mesmas chaves permitem a realização de assinaturas digitais (ECDSA-521 para realiza-

ção/validação de assinaturas em cima de *hashes* SHA-521) juntamente com o componente ASS, o qual realiza e valida as assinaturas digitais por meio dos componentes Check-out e Check-in (respectivamente), garantindo assim a **integridade** dos documentos.

Observa-se também que a orquestração dos componentes garante **autenticidade**: a AC emite certificados digitais e chaves privadas que servirão para realização de assinaturas digitais e para validação da confiança do certificado; o DI deve sempre se autenticar na solução adquirindo um JWT; a API sendo ponto de acesso ao IoT utilizado pelos DI's e Sistemas de Negócio, podendo realizar autenticação (validação do JWT); o serviço de autenticação (AUT), gerando os JWTs para os DI's e Sistemas de Negócio; e ASS por meio do Check-out e Check-in e seus processos de assinatura digital que garantem a autenticidade dos emissores dos documentos, alinhado com a validação da confiança do certificado.

### 7.1.2 Componentes e o combate às ameaças de segurança da informação

As medidas e contramedidas implementadas na solução proposta IoTSec2DCoC foram eficazes no ambiente controlado da prova de conceito, onde foi possível aferir que os componentes da solução apontados na matriz da Tabela 7.1 demonstraram que a AC contribui efetivamente na segurança, pois os certificados digitais e chaves privadas emitidos não configuram os túneis TLS nas comunicações, não possibilitando a leitura em claro das informações trafegadas (**ataque de sniffing**).

A AC também emite certificados digitais e chaves privadas que servirão para a realização de assinaturas digitais e para validação da confiança do certificado, não possibilitando **ataques de falsificação das assinaturas**; e o ASS (Check-out e Check-in) realiza os processos citados na seção 7.1.1. Quanto aos **ataques de falsificação de autenticação**, a autenticação do DI (ou sistemas de negócio), previamente cadastrados a partir da geração, um *token* JWT garante autenticidade. O serviço AUT oferece mecanismos que previnem o **ataque de força bruta** para autenticação, como o limite de autenticações com falha e respectivo bloqueio de senha.

O componente CB (ferramenta e padrão de microsserviços), baseado na ferramenta Hystrix do Spring Cloud, previsto na arquitetura da solução, não foi implementado no protótipo da PoC. Observa-se claramente que, por suas funcionalidades de acordo com arranjo proposto, ele teria mecanismos para minimizar indisponibilidades que podem degradar a solução como um todo, protegendo o IoTSec2DCoC de **perturbações do funcionamento normal**.

### 7.1.3 Componentes e o reforço às propriedades da solução e de integração

As demais funcionalidades da solução proposta IoTSec2DCoC e suas propriedades de integração demonstraram a eficiência necessária desejada no ambiente controlado da prova de conceito onde, de acordo com a matriz da Tabela 7.1, comprovaram que o serviço SR oferece o rastreamento no nível de requisições às APIs dos microsserviços, indicando qual serviço foi consumido



ou requisitado, inclusive indicando os erros desencadeados. Já o serviço SM fornece a centralização de logs, possibilitando o rastreamento dos eventos relacionados com logs dos componentes do sistema. Esse dois serviços possibilitam o **monitoramento e rastreabilidade**

Os middlewares IoT promovem a **integração de componentes**, DIs e Sistemas de Negócio, os quais são clientes em camadas Edge e Cloud, fornecendo a mediação e serviços a nível de IoT para todos os componentes da solução. O RS fornece a centralização do registro dos serviços da solução, podendo oferecer nomes (DNS) para facilitar a comunicação entre os serviços.

OS DIs realizam a **digitalização de documentos** (via scanner). Em seguida, ocorre a **classificação e extração de metadados** com a execução dos algoritmos do componente IA. Por fim, o componente EMP possibilita a revisão da digitalização e dos metadados, viabilizando, caso necessário, complementá-los ou corrigi-los.

## 7.2 RESULTADOS DA CLASSIFICAÇÃO AUTOMÁTICA DOS DOCUMENTOS E EXTRAÇÃO DE METADADOS

Para auxiliar a prova de conceito da solução proposta a partir da amostra descrita na seção 5.1.4, foi implementado um modelo simples de classificação automática de documentos e extração de metadados para que o fluxo do processo fosse mais ágil, conforme já detalhado na seção 5.5.4.

Cabe ressaltar que essa implementação foi realizada de forma a apenas dar agilidade ao processo, ou seja, utilizado como ferramenta de apoio, não sendo parte fundamental deste trabalho. No entanto, os resultados alcançados com a aplicação das técnicas de IA foram relevantes conforme pode ser observado na Figura 7.1. Nela é possível observar que o modelo de detecção de objetos para extração posterior do metadado atingiu um mAP com um IoU de 0.95 de 0.74.

Em detecção de imagens, uma métrica muito importante é a intersecção sobre união, do inglês *Intersection over Union* (IoU). Ela mede a quantidade da área predita pelo modelo sobre a área real do objeto, ou seja, mede o quanto os limites do *bounding box* predito sobrepõe os limites reais do objeto. Desse modo, utiliza-se um *threshold* de IoU para classificar se a predição está correta ou não.

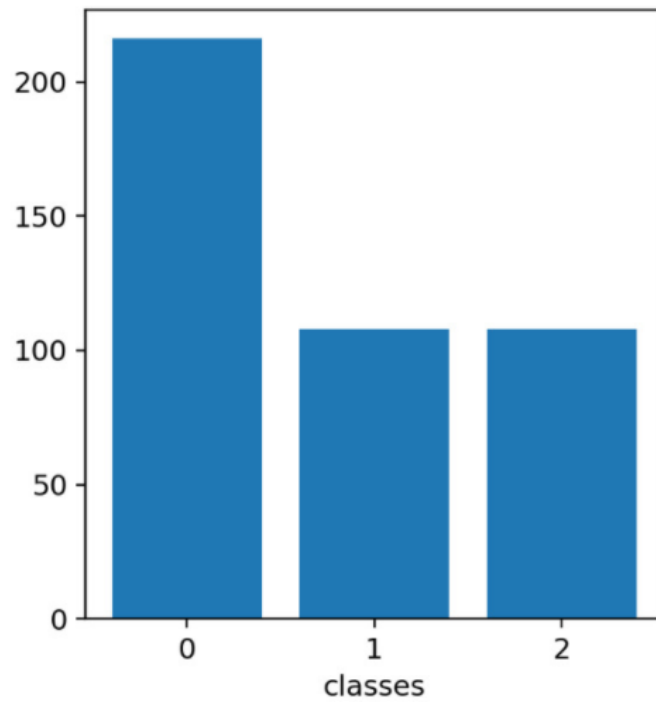


Figura 7.1: Distribuição de classes ao longo do dataset de treinamento de detecção de objetos.

A Figura 7.1 apresenta a distribuição de classes no *dataset* de treinamento utilizado para a detecção de objetos na imagem. Nela observa-se que a classe 0, a qual representa a classe Título, possui o dobro de entradas das outras classes, visto que, em cada imagem, havia exatamente duas regiões em que o Título estava presente: uma região com a Data, representada pela classe 1; e uma região para os Produtores, representada pela classe 2, de modo que há 100 entradas dessas classes.

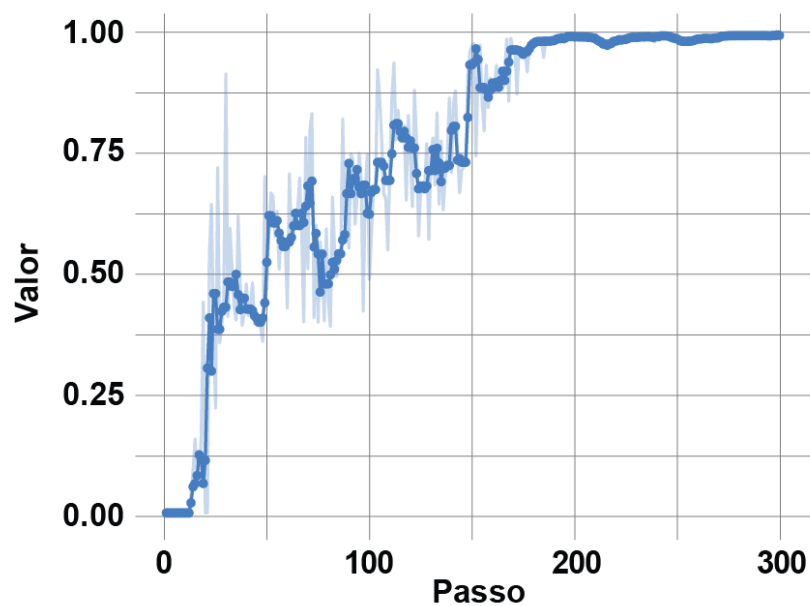


Figura 7.2: Análise de precisão ao longo das 300 épocas de treinamento.

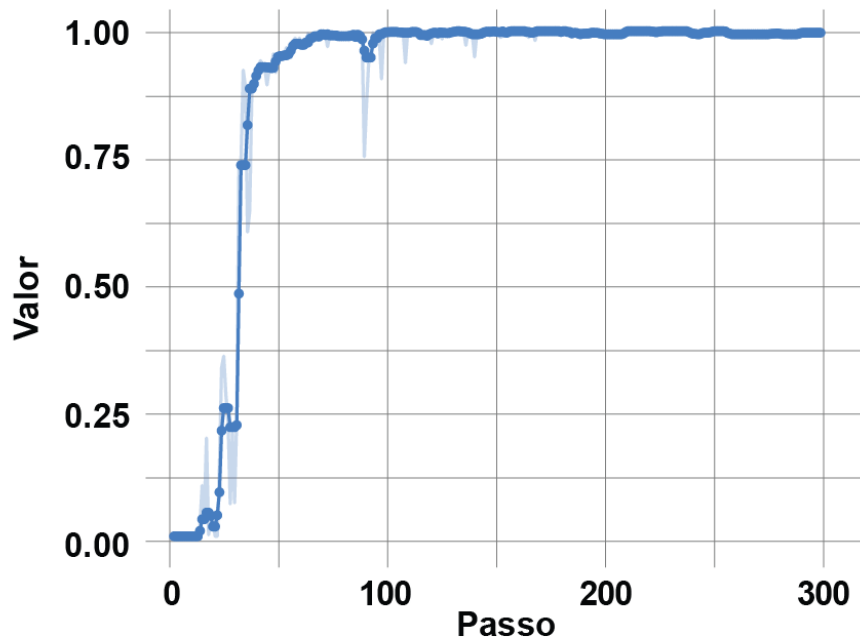


Figura 7.3: Análise da revocação ao longo das 300 épocas de treinamento para o sistema de detecção de objetos.

As Figuras 7.2 e 7.3 representam as métricas de precisão e revocação durante o treino do modelo ao longo das 300 épocas. Como os objetos a serem encontrados estão dispostos de forma padronizada e sem grandes variações, o modelo conseguiu distinguir sem erros a localização de cada objeto, tendo ambas precisão e revocação com 99,9%.

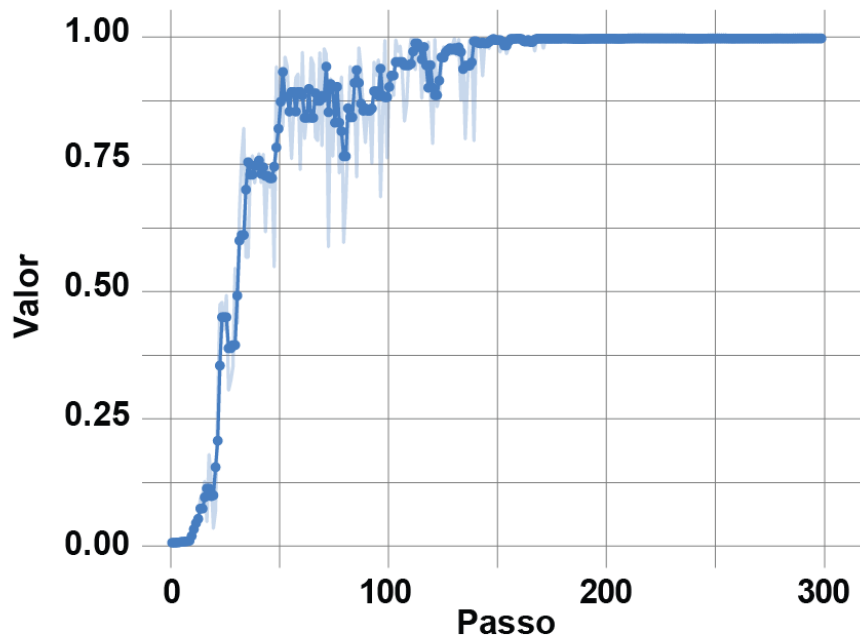


Figura 7.4: Valor médio de precisão com um threshold de 0.5 ao longo das 300 épocas de treinamento.

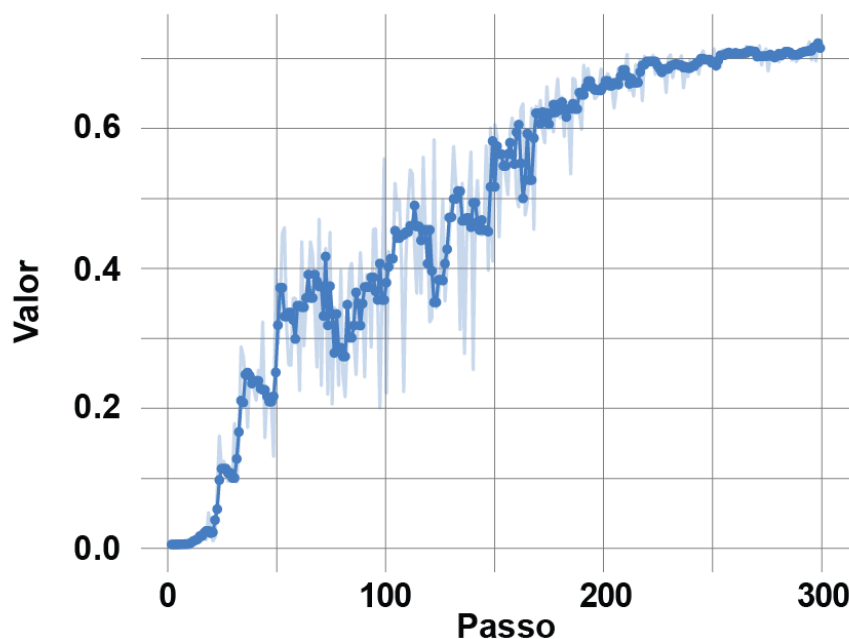


Figura 7.5: Valor médio de precisão com um threshold de 0.95 ao longo das 300 épocas de treinamento.

As Figuras 7.4 e 7.5 representam as métricas de mAP com dois *thresholds* diferentes para o IoU, em 0.5 e 0.95, ou seja, quando a área predita sobrepõe a área real em mais de 50% e em mais de 95%, respectivamente. Ao passo que as métricas de precisão e revocação foram próximas de 100%, era de esperar que a área sobreposta em um IoU 0.5 obtivesse uma mesma assertividade, representada pela Figura 7.4. Já quando o IoU sobe para 0.95, a precisão média cai para 0.74, porém não impacta no objetivo final da solução.

Uma vez que os objetos são localizados nas imagens, o processo de OCR fica mais assertivo, visto que não precisará realizar sobre toda imagem, e sim apenas em uma área específica, já sabendo o tipo de metadado localizado.

### 7.3 RESULTADOS DA RESILIÊNCIA DO MODELO A PERTURBAÇÃO DO FUNCIONAMENTO NORMAL

Uma solicitação básica para a PoC IoTSec2DCoC é demonstrar a plasticidade e flexibilidade da proposta quanto à sua configuração, bem como sua estabilidade para operar todo o fluxo documental, permitindo a instalação, configuração, implantação e o gerenciamento automatizado de todos os componentes e dispositivos conectados à rede e resolvendo problemas de desempenho. A escolha da PoC para validar a proposta é a *suite* Spring Cloud, a qual possui uma importante característica de auto-registro dos componentes, ou seja, cada dispositivo integrado na rede envia sua identidade para que o serviço de descoberta possa conhecê-lo, identificando rapidamente o que entrou ou saiu da rede. Esse recurso garante que os dispositivos sejam sempre verificados

para fornecerem seus serviços.

Essa escolha da PoC também permite a utilização de outros módulos, como Spring Cloud Zull, o qual atua como o Gateway de API e balanceador de carga, decidindo a rota de cada requisição que trafega pela rede, considerando a carga atual de recursos. Outro módulo é o Spring Cloud Hystrix, que atua como Circuit Breaker, identificando falhas e sobrecargas na rede, interrompendo as solicitações de serviço até detectar a recuperação de componentes e serviços. Por último, existe o módulo de gestão de logs Zipkin, operando para coletar todos os logs dos outros serviços que são parte da rede, permitindo, assim, o alarme e a análise de erros de execução para procurar rapidamente as correspondentes medidas de reparo.

Não obstante a integração dos serviços obtida com o Spring Cloud, com sua medida de contingência para manter os recursos mínimos funcionando e rotas suficientes para manter o fluxo completo do processo, há um ponto de atenção quanto ao gateway principal do Spring Cloud, sendo a única porta para processar todos os pedidos. É fundamental manter o monitoramento desse serviço para que, caso haja algum problema, a administração da linha de produção o resolva rapidamente.

Este conjunto de componentes estão previstos no modelo arquitetural da solução, mas não foram implementados na POC.

#### **7.4 RESULTADOS DA RESILIÊNCIA A ATAQUES DE FALSIFICAÇÃO**

As tentativas de falsificar o tráfego, copiando ou forjando pacotes para injetar tráfego malicioso, não foram bem-sucedidas, pois o invasor não conseguiu de forma síncrona executar ações para sintetizar uma assinatura válida no esquema de Serviços de Validação de Documentos e Assinatura digital (Check-out e Check-in), no Serviço de Autoridade Certificadora (CFSSL) e no Serviço de Autenticação (JWT).

A verificação sistemática das transferências é apontada como um pilar da proposta, uma vez que garante a integridade e autenticidade dos documentos em todo o fluxo da CoC documental. Essa tarefa é executada pelos pares de dispositivos virtuais Check-out e Check-in. O componente Check-out realiza a assinatura digital quando um documento é alterado ou quando um novo documento é gerado, antes de transferir o documento para a próxima entidade na cadeia. Por outro lado, o componente Check-in valida a assinatura digital mediante o recebimento de documentos por cada componente.

No ambiente controlado do laboratório da PoC, os resultados mostram as contramedidas sem falhas, embora seja aconselhável que os testes sejam realizados em uma configuração de estudo de caso mais exposta a ataques externos. Embora esse estudo de caso seja considerado para estudos posteriores, a estratégia adotada parece ser capaz de conter ataques o suficiente para desencadear outras medidas defensivas.

## **7.5 RESULTADOS DAS CONTRAMEDIDAS CONTRA VAZAMENTO DE CREDENCIAIS E CHAVES DE CRIPTOGRAFIA**

É importante, ao usar autenticação JWT, assinaturas digitais e TLS, levar em consideração a vulnerabilidade relacionada ao possível vazamento de credenciais (no caso de autenticação) e chaves privadas (em assinaturas e túneis seguros), mesmo considerando que esse risco surja apenas em circunstâncias incomuns. Uma credencial de Dispositivo Inteligente vazada permite que o invasor se autentique na solução e tente o envio de dados ilegítimos ao Middleware Edge, mas o invasor não conseguirá enviar arquivos devido à sua assinatura digital inválida (já que ele não possuirá um certificado digital e chave privada emitidos pelo CFSSL). Esse evento aciona a identificação do vazamento e o consequente bloqueio da conta do dispositivo no Serviço de Autenticação. Após uma investigação forense do caso, novas credenciais podem ser geradas e configuradas para o Dispositivo Inteligente afetado. O mesmo é verdadeiro para o caso em que o invasor tente ler conversas e assinar documentos, tentando enviá-los para o Middleware Edge: o invasor não teria sucesso, porque não seria autenticado com um JWT válido. Nesse caso, a detecção desse vazamento crítico acionará a revogação do certificado digital correspondente pelo CFSSL, informado na Lista de Revogação de Certificados - CRL, a qual permite ao próprio dispositivo identificar que o certificado é inválido, solicitando um novo.

## **7.6 RESULTADOS DA RESILIÊNCIA DO MODELO A ATAQUES CONTRA SENHAS**

O Serviço de Autenticação demonstrou ser eficiente a ataques de força bruta com dicionários, utilizado como contramedida uma estratégia de mitigação desse ataque através de uma política de bloqueio de conta, um método feito para se evitar a pressuposição de senhas, bloqueando a conta da identidade (pessoa ou dispositivo) que possua múltiplas tentativas de autenticação falhas.

Com essa configuração de política de bloqueio de conta, o ataque de força bruta é minimizado, pois, para se testar uma grande combinação de senhas (um dicionário de senhas pode ter milhões de registros), seria necessário um espaço de tempo inviável para o atacante. Supondo um dicionário de senhas de 1 milhão de combinações, com uma política de bloqueio de conta por 1 hora e com, no máximo, 3 tentativas falhas para o bloqueio, o atacante levaria até 38 anos para testar todas as possibilidades do dicionário.

## **7.7 RESULTADOS DA RESILIÊNCIA DO MODELO A ATAQUES DOS**

Ataques de negação de serviço (DoS) tem como objetivo ferir o princípio de disponibilidade de sistemas computacionais, em que o atacante, de alguma forma, tenta sobrecarregar os recursos do sistema para ele ficar indisponível aos usuários legítimos e autorizados. Existem inúmeras

técnicas de ataque DoS (ICMP Flood, TCP SYN Flood, Ataques Distribuídos – DDoS, etc.) que podem ser combinadas de inúmeras formas para atingir o alvo. Por conta do potencial dano que esses casos podem causar e por eles poderem estar atrelados a níveis mais baixos de rede, as soluções de proteção para esses tipos de ataque raramente se encontram no nível de Sistema da solução, ou seja, os módulos do IoTSec2DCoC como software não irão conseguir tratá-los sozinhos, sendo necessário o auxílio de softwares (e até mesmo hardware) externos à solução.

Ataques de DoS podem ser tratadas de forma mais simples com uso de *firewalls* tradicionais, caso se conheça a origem do ataque e seja criada uma regra de rejeição de mensagens oriundas dessa origem, mas não é muito eficaz. Existem empresas que fornecem soluções de CDN – Content Delivery Network, sendo um conjunto de servidores distribuídos pela internet que realizam *caching*, proteção contra DoS e, até mesmo, recursos mais avançados de segurança, como WAF's – Web Application Firewalls. Esse tipo de conjunto de Hardware e Software pode ser contratado pelas organizações que optarem pela utilização do IoTSec2DCoC para fornecer maior segurança à solução.

## 7.8 RESULTADOS DO MONITORAMENTO E RASTREABILIDADE

A partir dos logs armazenados no Serviço de Monitoramento, gerados na prova de conceito conforme a subseção 6.3, foi possível analisar os tempos decorridos sequenciados na cadeia de produção documental, proposta neste trabalho, incluindo medidas de tempo para a sequência de ações e transferências do Dispositivo Inteligente para os Middlewares IoT, o Empacotador, o RDC-Arq e, finalmente, a Plataforma de Acesso, conforme ilustrado no diagrama da Figura 7.6.

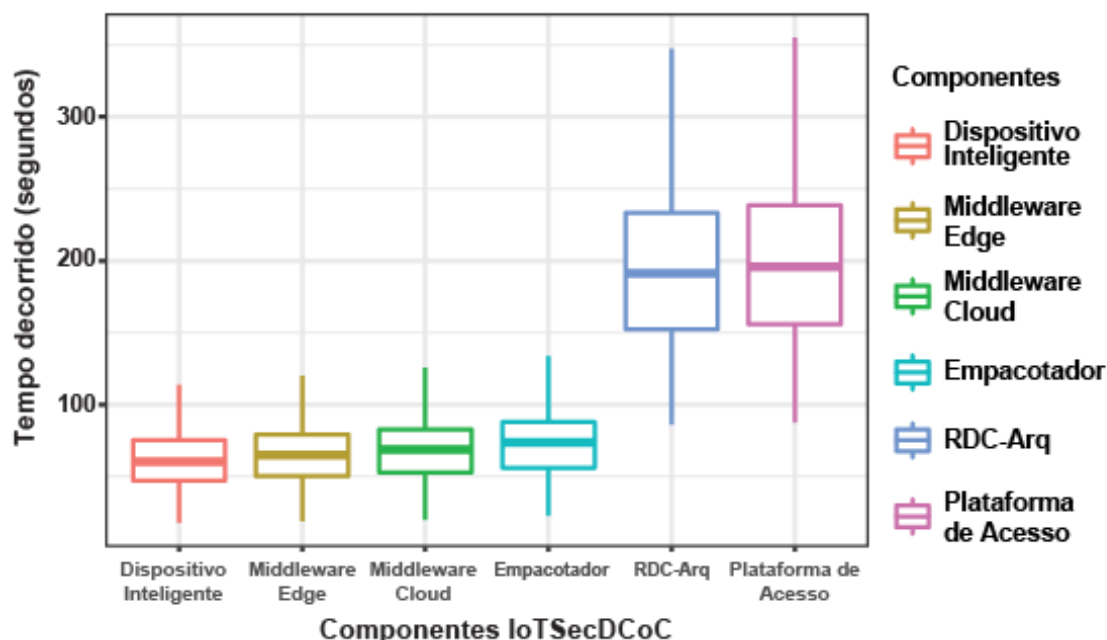


Figura 7.6: Tratamentos e transferências sequenciados de componentes IoTSec2DCoC.

Os dados da Figura 7.6 são coletados com a ferramenta Zipkin, rastreador de log centralizado do Spring Cloud. Um subconjunto de documentos foi processado em uma configuração controlada, e os resultados coletados foram condensados na Figura 7.6, a qual mostra *boxplots* para os seis componentes envolvidos na produção desses documentos. O eixo *X* representa, da esquerda para a direita, a sequência de tratamentos, desde a entrada na cadeia de produção IoTSec2DCoC até o final dessa cadeia. O eixo *Y* é o tempo decorrido, tanto em cada unidade de tratamento quanto na transferência para a próxima. Como há variação do tempo decorrido em função do tamanho do documento e da carga do sistema operacional e da rede no momento do processamento de cada documento, a escolha dos *boxplots* permite representar os valores médios e sua variância. Portanto, a barra horizontal no meio de cada caixa é o valor médio para os documentos definidos, enquanto os limites superior e inferior da caixa são o terceiro e primeiro quartis, respectivamente.

Sendo assim, há uma ideia do tempo decorrido acumulado para toda a cadeia produtiva. Por exemplo, o tempo de processamento do Middleware Edge no gráfico refere-se a um intervalo de tempo após o Dispositivo Inteligente ter transferido o documento somado ao tempo de processamento do próprio middleware.

Além disso, para cada variável, a diferença entre o terceiro e o primeiro quartil é calculada internamente, dando o chamado intervalo interquartil. A partir desse intervalo, o Limite Teórico Superior (LTS) e o Limite Teórico Inferior (LTI) são calculados da seguinte forma:  $LTI = PrimeiroQuartil - (1.5 * IntervaloInterquartil)$  e  $LTS = TerceiroQuartil + (1.5 * IntervaloInterquartil)$ . Assim, as áreas abaixo e acima das linhas médias dentro das caixas referem-se aos valores que estão dentro desses limites superior e inferior, mas fora do intervalo interquartil. Se houver uma observação maior que o LTS ou menor que o LTI, essa informação torna-se um ponto isolado no gráfico, indicando a possibilidade da existência de um *outlier*, ou seja, um evento de atraso anormal na linha de produção.

As informações apresentadas na Figura 7.6 constituem de fato um instrumento de painel de controle para monitorar a cadeia de controladores, seguindo o paradigma IIoT para a linha de produção documental.



## 8 CONCLUSÃO

Esta tese busca contribuir para a confiança nos documentos digitais, tão necessária para que pessoas se beneficiem de um ambiente sustentável no que diz respeito à produção e preservação de documentos. Para construir a confiança na digitalização de documentos e processos, foi proposto um conjunto totalmente distribuído de medidas de segurança para apoiar a CoC de documentos digitais, conforme recomendado pelo modelo de referência OAIS, um padrão ISO amplamente aceito pela comunidade internacional. Assim, a proposta do Modelo de Segurança da Internet das Coisas para a CoC Digital (IoTSec2DCoC) visa aumentar a segurança, integrando medidas orientadas para garantir a CoC, distribuídas desde o produtor do documento até seu consumidor.

Para tanto, o modelo proposto na tese unifica diversas áreas de conhecimento originalmente não relacionadas, como a cadeia de custódia de documentos arquivísticos digitais, arquitetura de uma solução IoT para prover segurança em um ambiente ciberfísico de produção e em suas contrapartes lógicas, aprendizado de máquina aplicada à extração de dados de uma fonte não estruturada e a elaboração de um modelo adversarial em seu domínio, de modo que a proposta inclua medidas de segurança e serviços operando efetivamente para garantir a cadeia de custódia de documentos arquivísticos.

Para verificar o atendimento às hipóteses apresentadas na subseção (3.2), são apresentadas a seguir conclusões específicas sobre cada hipótese.

Quanto à hipótese 1, foram implementadas medidas de segurança no modelo que incluem identificação, autenticação, proteção e monitoramento dos fluxos de comunicação e armazenamento. Todos os componentes são integrados a uma infraestrutura de rede IoT hierárquica, alinhada a uma arquitetura de software de microsserviços para melhor integração dos componentes de segurança da solução, utilizando técnicas de criptografia e assinatura digital para garantir a confidencialidade, autenticidade e integridade dos pacotes de informações transferidos. A validação utilizando a PoC levou a resultados que comprovam o atendimento da hipótese 1.

Quanto à hipótese 2, a estrutura de tratamento de documentos baseada em IIoT e suas medidas de segurança totalmente distribuídas apresentam um conjunto de recursos para oferecer suporte a várias classes de confiança de usuários em cidades sustentáveis em relação aos documentos produzidos, protegidos e preservados. Esses serviços permanecem explicitamente ativos e visíveis para o usuário dos documentos da CoC, demonstrando, assim, a efetividade da hipótese 2. Vale observar que, como a confiança computacional e a segurança da informação contribuem mutuamente para construir uma à outra, esta tese contribui efetivamente para construir a confiança em documentos digitais, garantindo a cadeia de custódia (CoC) dos documentos digitais produzidos e preservados.

Com relação às hipóteses 3 e 4, esta tese considera o paradigma IIoT e o modelo da CoC para conceber um serviço de segurança baseado em IIoT para a cadeia de custódia documental,

estruturando, assim, uma cadeia de produção segura que apresenta a plasticidade adequada para cidades inteligentes e sustentáveis.

A digitalização de documentos e processos é considerada um fator essencial para cidades sustentáveis. Embora contribua para sociedades e ambientes urbanos favoráveis ao clima, também apoia a confiança nas interações socioeconômicas. Como esta tese compreende uma estrutura de digitalização de documentos baseada em IIoT e medidas de segurança para documentos produzidos por essas interações, entende-se que as contribuições para cidades sustentáveis vêm principalmente do princípio de que os serviços de segurança totalmente distribuídos propostos, apoiados por uma instância adaptável de IIoT, constituem um *framework* que sustenta a confiança do usuário em relação aos documentos protegidos e preservados, pois esses serviços permanecem explicitamente ativos e visíveis para o usuário da CoC de documentos arquivísticos digitais.

A desmaterialização de coleções de documentos impressos a partir de seu suporte analógico e a substituição de documentos impressos por digitais têm efeito direto na pegada de carbono ao permitir linhas de produção documental que mostram melhor desempenho ambiental em relação a emissões diretas e indiretas, sequestro de carbono em florestas, valor de bioenergia, emissões evitadas, bem como custo de armazenamento, energia e ocupação do espaço urbano.

Desse modo, comprova-se o atendimento às hipóteses 3 e 4 da tese.

Quanto às hipóteses 5 e 6 adotadas na tese, a adoção dessa hipóteses se reflete na concepção e na experimentação da proposta.

O IoTSec2DCoC proposto compreende um conjunto de recursos e meios para garantir as propriedades de segurança exigidas para documentos digitalizados, sendo um conceito abrangente que permite integrar as medidas de segurança em todo o fluxo de informação produtor-consumidor. Como o usuário é informado explicitamente sobre essas medidas, ele pode confiar que as proteções disponíveis foram implementadas para seus documentos digitais. Considerando os conceitos de confiança discutidos na seção 2.6.1, vale destacar as classes de confiança específicas que o usuário pode perceber a partir dos recursos IoTSec2DCoC projetados: confiança na identificação de dispositivos e operadores CoC; confiança na autenticação da fonte do documento; confiança na integridade do documento do produtor ao consumidor; confiança na Disponibilidade do fluxo de CoC; confiança em múltiplas medidas para monitorar fluxos de comunicação e armazenamento; confiança na preservação e disponibilidade de documentos; confiança em que as partes da CoC não negarão a produção de documentos; confiança em que documentos falsos não serão atribuídos ao CoC; confiança nos documentos pode ser verificada para fins forenses. Portanto, em resumo, uma soma de recursos para apoiar a confiança do usuário nos documentos produzidos, protegidos e preservados.

Para validar a proposta, um protótipo PoC desenvolvido foi submetido a testes de resiliência em cenários de ataques e com interrupções de operação normal. Os testes permitiram a avaliação do equilíbrio e da estabilidade dos serviços, da latência de comunicação com instâncias do Middleware Edge, bem como da rastreabilidade e recuperação em caso de ataque bem-sucedido a um ou mais componentes. Assim, foram verificadas as hipóteses 5 e 6 adotadas na tese.

A validação experimental também indica que, de forma geral, as hipóteses (3.2) obtiveram respostas de sucesso, pois o modelo proposto, considerando as limitações deste estudo, pode aumentar a segurança das linhas de produção documental que adotam o modelo de referência OAIS, promovendo um fluxo documental protegido, ininterrupto e rastreável - do produtor ao consumidor de informações.

Todos os componentes foram incorporados à CoC em estrita conformidade com a ISO 14721:2003 [23], respeitando outros padrões relativos ao pacote de submissão dos metadados [28] e os atributos e responsabilidades pertinentes a um repositório seguro [31]. Esses critérios de conformidade demonstram que, com base nos resultados apresentados, a solução IoTSec2DCoC está pronta para ser dimensionada para cenários reais.

## 8.1 TRABALHOS FUTUROS

Como trabalho futuro, a possibilidade de implantar agentes IoT em repositórios redundantes (Archivematica, RODA, LOCKSS) pode ser investigada para estender a abordagem de monitoramento IoT, nos aspectos de plasticidade e flexibilidade, para linhas documentais de produção e disseminação específicas.

Uma tecnologia promissora para preservação de documentos é o Blockchain, o qual também seria interessante para registrar transferências de documentos em uma cadeia de produção. Apesar dessa vantagem, o Blockchain não foi projetado para interferir em uma linha de produção, função que requer que a IoT traga a plasticidade e flexibilidade de implantação em linhas de produção documentais variáveis, permitindo que a medida de segurança totalmente distribuída estenda seu controle a cada componente que contribui para a CoC dos documentos produzidos e transferidos. Assim, para cumprir suas possibilidades, questões interessantes relacionadas devem ser abordadas. Como um sistema de software de Blockchain puro dificilmente é capaz de agir no mundo físico, a IoT é necessária para processamento de documentos, requerendo medidas para implementar e integrar instâncias de Blockchain para linhas independentes de produção de documentos.

Como nas organizações reais o cenário é possivelmente a necessidade de múltiplas linhas de produção de documentos, fisicamente e logicamente separadas, é necessário pesquisar, implementar e integrar estratégias para coordenar múltiplas cadeias de custódia.

Seria relevante a análise conceitual e a realização de testes de validação estática e dinâmica de soluções que adotam a CoC documental convencional.

Outra boa oportunidade de trabalhos futuros seria realizar experimentos com os componentes e serviços destinados a prover resiliência a perturbação do funcionamento normal.

Finalmente, vale a pena considerar que, no ambiente de laboratório PoC controlado, os resultados mostram um *framework* altamente resiliente sem falhas, embora seja aconselhável que os

testes sejam realizados em uma configuração de estudo de caso mais exposta a ataques externos. Embora este estudo de caso seja considerado para um estudo mais aprofundado, a estratégia adotada mostra-se capaz de prevenir e conter ataques, desencadeando medidas de defesa abrangentes.

## REFERÊNCIAS BIBLIOGRÁFICAS

- 1 DONALDSON, D. R. Trust in Archives–Trust in Digital Archival Content Framework. *Archivaria*, Association of Canadian Archivists, v. 88, p. 50–83. Disponível em: <https://www.muse.jhu.edu/article/740193>, 2019.
- 2 DE OLIVEIRA ALBUQUERQUE, R.; VILLALBA, L. J. G.; OROZCO, A. L. S.; DE SOUSA JÚNIOR, R. T.; KIM, T.-H. Leveraging information security and computational trust for cybersecurity. *Journal of Supercomputing*, 2015. ISSN 15730484 09208542. Disponível em: <https://doi.org/10.1007/s11227-015-1543-4>.
- 3 ISO-14721. Iso 14721:2012: Space data and information transfer systems: Open archival information system – reference model. *International Organization for Standardization.*, p. 1–156, 2003.
- 4 CONARQ, A. N. *Diretrizes para a Presunção de Autenticidade de Documentos Arquivísticos Digitais*. 2012. Disponível em: [http://conarq.gov.br/images/publicacoes\\_textos/conarq\\_presuncao\\_autenticidade\\_completa.pdf](http://conarq.gov.br/images/publicacoes_textos/conarq_presuncao_autenticidade_completa.pdf).
- 5 SMITH, M.; BRONNER, W.; SHIMOMURA, E.; LEVINE, B.; FROEDE, R. *Quality assurance in drug testing laboratories*. September 1990. 503–516 p. Disponível em: <http://europepmc.org/abstract/MED/2253447>.ISSN:0272-2712.
- 6 ISO. ISO 16363: 2012 space data and information transfer systems – audit and certification of trustworthy digital repositories. *International Organization for Standardization.*, p. 1–70, 2012.
- 7 KAY, A. Tesseract: an open-source optical character recognition engine. *Linux Journal*, Belltown Media, v. 2007, n. 159, p. 2. ISSN: 1075–3583, 2007.
- 8 DA SILVA, D. A.; DE SOUSA, R. T.; DE OLIVEIRA ALBUQUERQUE, R.; SANDOVAL OROZCO, A. L.; GARCÍA VILLALBA, L. J. Iot-based security service for the documentary chain of custody. *Sustainable Cities and Society*, v. 71, p. 102940, 2021. ISSN 2210-6707. Disponível em: <https://doi.org/10.1016/j.scs.2021.102940>.
- 9 SILVA, D. A. d.; JÚNIOR, R. T. d. S. Iot-based security service for the documentary chain of custody (iotsec2coc):how to protect digital documents in the iot flow from information producer to consumer? *Brasil-Italy webinar, Flexible and Autonomous Manufacturing Systems for Custom-Designed Products (FASTEN)*, 2020, 2020.
- 10 SILVA, D. A. da; TORRES, J. A. S.; PINHEIRO, A.; FILHO, F. L. de C.; MENDONÇA, F. L.; PRACIANO, B. J.; KFOURI, G. de O.; SOUSA, R. T. de. Inference of driver behavior using correlated iot data from the vehicle telemetry and the driver mobile phone. In: IEEE. *2019 Federated Conference on Computer Science and Information Systems (FedCSIS)*. [S.l.], 2019. p. 487–491.
- 11 SILVA, D. A. d.; MACHADO, P. L.; COELHO, V. C. G.; BARBOSA, R. V.; MENDONÇA, F. L. L. d.; SANTOS, D. P. d.; JÚNIOR, R. T. d. S. Produção de indicadores de empregabilidade com base em técnicas de mineração de big data e business intelligence. *Inclusão Social*, v. 12, n. 2, jun. 2019. Disponível em: <http://revista.ibict.br/inclusao/article/view/4670>.
- 12 CONARQ. *e-ARQ. Modelo de requisitos para sistemas informatizados de gestão arquivística de documentos*. [S.l.]: Arquivo Nacional, 2011.
- 13 CONARQ. *Glossário de Documentos Arquivísticos Digitais*. setembro 2016. Disponível em: [http://conarq.arquivonacional.gov.br/images/ctde/Glossario/2016-CTDE-Glossario\\_V7\\_public.pdf](http://conarq.arquivonacional.gov.br/images/ctde/Glossario/2016-CTDE-Glossario_V7_public.pdf).

- 14 ROUSSEAU, J.-Y.; COUTURE, C.; ARÈS, F. *Os fundamentos da disciplina arquivística*. [S.l.]: Publicações Dom Quixote Lisboa, 1998.
- 15 BRASIL, P. d. R. *Dispõe sobre a política nacional de arquivos públicos e privados e dá outras providências*. jan 1991. Disponível em: [http://www.planalto.gov.br/ccivil\\_03/Leis/L8159.htm](http://www.planalto.gov.br/ccivil_03/Leis/L8159.htm).
- 16 BRASIL, P. d. R. *Dispõe sobre o uso do meio eletrônico para a realização do processo administrativo no âmbito dos órgãos e das entidades da administração pública federal direta, autárquica e fundacional*. Brasília, DF, 2015. Disponível em: [http://www.planalto.gov.br/ccivil\\_03/\\_Ato2015-2018/2015/Decreto/D8539.htm](http://www.planalto.gov.br/ccivil_03/_Ato2015-2018/2015/Decreto/D8539.htm).
- 17 BRASIL, P. d. R. *Autoriza o armazenamento, em meio eletrônico, óptico ou equivalente, de documentos públicos ou privados, compostos por dados ou por imagens, observado o disposto nesta Lei, nas legislações específicas e no regulamento*. Brasília, DF, 2020a.
- 18 BRASIL, P. d. R. *Regulamenta o disposto no inciso X do caput do art. 3º da Lei nº 13.874, de 20 de setembro de 2019, e no art. 2º-A da Lei nº 12.682, de 9 de julho de 2012, para estabelecer a técnica e os requisitos para a digitalização de documentos públicos ou privados, a fim de que os documentos digitalizados produzam os mesmos efeitos legais dos documentos originais*. Brasília, DF, 2020b. Disponível em: <https://www.in.gov.br/en/web/dou/-/decreto-n-10.278-de-18-de-marco-de-2020-248810105>.
- 19 FORCE, A. T. Requirements for assessing and maintaining the authenticity of electronic records. *InterPARES Project, Vancouver, Canada*. Disponível em: [http://www.interpares.org/display\\_file.cfm?doc=ip1\\_authenticity\\_requirements.pdf](http://www.interpares.org/display_file.cfm?doc=ip1_authenticity_requirements.pdf), 2002.
- 20 SANTOS, C. A. C. M. dos; LUZ, C. dos S.; AGUIAR, F. L. Introdução à organização de arquivos: conceitos arquivísticos para bibliotecários. 2014.
- 21 FLORES, D.; ROCCO, B. C. de B.; SANTOS, H. M. dos. Cadeia de custódia para documentos arquivísticos digitais. *Acervo*, v. 29, n. 2 jul-dez, p. 117–132, 2016.
- 22 CONARQ, A. N. *Resolução nº 39, de 29 de abril de 2014*. 2014. Disponível em: <http://www.conarq.gov.br/index.php/resolucoes-do-conarq/281-resolucao-n-39,-de-29-de-abril-de-2014>.
- 23 CCSDS. *Reference Model for an Open Archival Information System (OAIS) - CCSDS 650.0-M-2. Magenta Book also available as ISO 16363:2011*. Disponível em: <https://public.ccsds.org/pubs/650x0m2.pdf>; CCSDS Secretariat, 2012. 135 p.
- 24 CCSDS. Reference model for an open archival information system (oais) - ccsds 650.0-r-1. red book. <http://ftp.ccsds.org/ccsds/documents/pdf/CCSDS-650.0-R-1%20.pdf>, Issue 1, 1999.
- 25 ISO-14721B. Iso 14721:2012: Space data and information transfer systems: Open archival information system – reference model. *International Organization for Standardization*., p. 1–156, 2012.
- 26 NBR-15472. Abnt nbr 15472: Sistemas espaciais de dados e informações - modelo de referência para um sistema aberto de arquivamento de informação (saai). *Associação Brasileira de Normas Técnicas*, 2007.
- 27 ROCHA, C. L. Repositórios para a preservação de documentos arquivísticos digitais. *Acervo*, v. 28, n. 2 jul-dez, p. 180–191. ISSN: 2237–8723, 2015.
- 28 CDS. *ISAD(G): General International Standard Archival Description - Second edition*. 2. ed. <https://www.ica.org/en/isadg-general-international-standard-archival-description-second-edition>; ICA, International Council on Archives, 2000. (1, v. 2). ISBN 0-9696035-5-X.

- 29 NOBRADE, C. N. d. A. *Nobrade: norma brasileira de descrição arquivística*. [S.l.]: Conselho nacional de arquivos, 2006.
- 30 CONARQ, A. N. *Diretrizes para a implementação de repositórios arquivísticos digitais confiáveis*. 2015. Disponível em: [http://www.conarq.gov.br/images/publicacoes\\_textos/diretrizes\\_rdc\\_arq.pdf](http://www.conarq.gov.br/images/publicacoes_textos/diretrizes_rdc_arq.pdf).
- 31 RLG-OCLC; OTHERS. *Trusted Digital Repositories: Attributes and Responsibilities*. An RLG-OCLC Report. *Mountain View, CA.: RLG*). Disponível em: <http://www.rlg.org/legacy/longterm/repositories.pdf>, 2002.
- 32 ISO-16363. *Space Data and Information Transfer Systems-Audit and Certification of Trustworthy Digital Repositories: ISO 16363*. [S.l.]: ISO, 2012.
- 33 STALLINGS, W. *Cryptography and network security: principles and practice*. [S.l.]: Pearson Upper Saddle River, 2019. ISBN 0135764181.
- 34 OWASP, T. *Top 10-2017 the ten most critical web application security risks*. Disponível em: [owasp.org/images/7/72/OWASP\\_Top\\_10-2017\\_%28en%29.pdf](http://www.owasp.org/images/7/72/OWASP_Top_10-2017_%28en%29.pdf), v. 29, 2017.
- 35 CERT.BR, C. d. E. R. e. T. d. I. d. S. n. B. *Cartilha de Segurança para Internet, Versão 4.0*. [S.l.]: Comitê Gestor da Internet no Brasil, 2012.
- 36 SLOAN, J. D. *High Performance Linux Clusters with OSCAR, Rocks, OpenMosix, and MPI: A Comprehensive Getting-Started Guide*. [S.l.]: "O'Reilly Media, Inc.", 2004. ISBN 9780596005702.
- 37 COULOURIS, G. F.; DOLLIMORE, J.; KINDBERG, T. *Distributed systems: concepts and design*. [S.l.]: pearson education, 2012. ISBN 9780133001372.
- 38 BUYYA, R.; BROBERG, J.; GOSCINSKI, A. M. *Cloud computing: Principles and paradigms*. [S.l.]: Wiley Publishing, 2011. ISBN 9780470887998.
- 39 MELL, P.; GRANCE, T. *The nist definition of cloud computing*. *Computer Security Division, Information Technology Laboratory, National*, n. 800-145, 2011. Disponível em: <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>.
- 40 SHI, W.; CAO, J.; ZHANG, Q.; LI, Y.; XU, L. *Edge computing: Vision and challenges*. *IEEE Internet of Things Journal*, v. 3, n. 5, p. 637–646, 2016. Disponível em: <https://doi.org/10.1109/JIOT.2016.2579198>.
- 41 NEWMAN, S. *Building microservices: designing fine-grained systems*. [S.l.]: "O'Reilly Media, Inc.", 2015. ISBN 9781491950357.
- 42 MARTINS, L. M. C. e.; FILHO, F. L. d. C.; JÚNIOR, R. T. d. S.; GIOZZA, W. F.; COSTA, J. a. P. C. da. *Increasing the dependability of iot middleware with cloud computing and microservices*. In: . *New York, NY, USA: Association for Computing Machinery*, 2017. (UCC '17 Companion), p. 203–208. ISBN 9781450351959. Disponível em: <https://doi.org/10.1145/3147234.3148092>.
- 43 IOT-GSI/ITU. *IoT-GSI/ITU. Overview of the Internet of things*. 2012. Disponível em: <http://handle.itu.int/11.1002/1000/11559>.
- 44 FERREIRA, H. G. C.; CANEDO, E. D.; SOUSA, R. T. de. *A ubiquitous communication architecture integrating transparent upnp and rest apis*. *International Journal of Embedded Systems*, v. 6, n. 2-3, p. 188–197, 2014. PMID: 63816. Disponível em: <https://doi.org/10.1504/IJES.2014.063816>.
- 45 SILVA, C. C. d. M.; CALDAS, F. L. d.; MACHADO, F. D.; MENDONÇA, F. L.; JÚNIOR, R. T. de S. *Proposta de auto-registro de serviços pelos dispositivos em ambientes de iot*. *34º Simpósio Brasileiro de Telecomunicações e Processamento de Sinais*, 2016.

- 46 FERREIRA, H. G. C.; de Sousa Junior, R. T. Security analysis of a proposed internet of things middleware. *Cluster Computing*, Springer, v. 20, n. 1, p. 651–660, 2017. Disponível em: <https://doi.org/10.1007/s10586-017-0729-3>.
- 47 KHAN, W. Z.; REHMAN, M. H.; ZANGOTI, H. M.; AFZAL, M. K.; ARMI, N.; SALAH, K. Industrial internet of things: Recent advances, enabling technologies and open challenges. *Computers & Electrical Engineering*, Elsevier, v. 81, p. 106522, 2020. Disponível em: <https://doi.org/10.1016/j.compeleceng.2019.106522>.
- 48 ESPOSITO, C.; CASTIGLIONE, A.; PALMIERI, F.; SANTIS, A. D. Integrity for an Event Notification Within the Industrial Internet of Things by Using Group Signatures. *IEEE Transactions on Industrial Informatics*, v. 14, n. 8, p. 3669–3678, 2018. Disponível em: <https://doi.org/10.1109/TII.2018.2791956>.
- 49 CRUZ, R. H. Proposição de um modelo e sistema de gerenciamento de dados distribuídos para internet das coisas–gddiot. 2017.
- 50 CRUZ, R. H.; JUNIOR, R. de S.; HOLANDA, M. de; ALBUQUERQUE, R. de O.; VILLALBA, L. G.; KIM, T.-H. Distributed data service for data management in internet of things middleware. *Sensors*, Multidisciplinary Digital Publishing Institute, v. 17, n. 5, p. 977, 2017. Disponível em: <https://doi.org/10.3390/s17050977>.
- 51 FERREIRA, H. G. C.; CANEDO, E. D.; SOUSA, R. T. de. A ubiquitous communication architecture integrating transparent UPnP and REST APIs. *International Journal of Embedded Systems*, Inderscience Publishers Ltd, v. 6, n. 2-3, p. 188–197, 2014. Disponível em: <https://doi.org/10.1504/IJES.2014.063816>.
- 52 MENEZES, J.; COSTA, P.; CUNHA, D.; FILHO, F.; MARTINS, L.; MENDONÇA, F. L. *Desenvolvimento de modelo hierárquico de middlewares com aplicação de fog computing para redes IoT*. 12 2019. 155-162 p. Disponível em: [https://doi.org/10.33965/ciaca2019\\_201914L020](https://doi.org/10.33965/ciaca2019_201914L020).
- 53 TONG, L.; LI, Y.; GAO, W. A hierarchical edge cloud architecture for mobile computing. In: IEEE. *IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on Computer Communications*. [S.l.], 2016. p. 1–9.
- 54 FILHO, F. L. d. C.; MARTINS, L. M. C. e.; ARAÚJO, I. P.; MENDONÇA, F. L. L. d.; COSTA, J. P. C. L. da; JÚNIOR, R. T. d. S. Design and evaluation of a semantic gateway prototype for iot networks. In: . New York, NY, USA: Association for Computing Machinery, 2017. (UCC '17 Companion), p. 195–201. ISBN 9781450351959. Disponível em: <https://doi.org/10.1145/3147234.3148091>.
- 55 RUSSELL, S. J.; NORVIG, P. *Artificial intelligence: a modern approach*. [S.l.]: Malaysia; Pearson Education Limited, 2020. ISBN 9780136958420.
- 56 NILSSON, N. J. *Principles of artificial intelligence*. [S.l.]: Morgan Kaufmann, 2014. ISBN 1493306065.
- 57 HINTON, G.; VINYALS, O.; DEAN, J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- 58 GARDNER, M. W.; DORLING, S. Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. *Atmospheric environment*, Elsevier, v. 32, n. 14-15, p. 2627–2636, 1998.
- 59 LECUN, Y.; BENGIO, Y.; HINTON, G. Deep learning. *Nature*, v. 521, p. 436–44, 05 2015. Disponível em: <https://doi.org/10.1038/nature14539>.



- 60 REDMON, J.; DIVVALA, S.; GIRSHICK, R.; FARHADI, A. You only look once: Unified, real-time object detection. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. [s.n.], 2016. p. 779–788. Disponível em: <https://doi.org/10.1109/CVPR.2016.91>.
- 61 REDMON, J.; FARHADI, A. Yolo9000: Better, faster, stronger. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. [s.n.], 2017. p. 6517–6525. Disponível em: <https://doi.org/10.1109/CVPR.2017.690>.
- 62 JU, M.; LUO, H.; WANG, Z.; HUI, B.; CHANG, Z. The application of improved yolo v3 in multi-scale target detection. *Applied Sciences*, v. 9, n. 18, 2019. ISSN 2076-3417. Disponível em: <https://doi.org/10.3390/app9183775>.
- 63 HUANG, Y.; YAN, Q.; LI, Y.; CHEN, Y.; WANG, X.; GAO, L.; TANG, Z. A yolo-based table detection method. In: *2019 International Conference on Document Analysis and Recognition (ICDAR)*. [s.n.], 2019. p. 813–818. Disponível em: <https://doi.org/10.1109/ICDAR.2019.00135>.
- 64 ADNANE, A.; BIDAN, C.; de Sousa Júnior, R. T. Trust-based security for the olsr routing protocol. *Computer Communications*, v. 36, n. 10, p. 1159–1171, 2013. ISSN 0140-3664. Disponível em: <https://doi.org/10.1016/j.comcom.2013.04.003>.
- 65 YAHALOM, R.; KLEIN, B.; BETH, T. Trust relationships in secure systems-a distributed authentication perspective. *Proceedings 1993 IEEE Computer Society Symposium on Research in Security and Privacy*, p. 150–164, 1993. Disponível em: <https://doi.org/10.1109/RISP.1993.287635>.
- 66 STOYANOVA, M.; NIKOLOUDAKIS, Y.; PANAGIOTAKIS, S.; PALLIS, E.; MARKAKIS, E. K. A survey on the internet of things (iot) forensics: Challenges, approaches and open issues. *IEEE Communications Surveys & Tutorials*, IEEE, 2020. Disponível em: <https://doi.org/10.1109/COMST.2019.2962586>.
- 67 AAZAM, M.; ZEADALLY, S.; HARRAS, K. A. Deploying fog computing in industrial internet of things and industry 4.0. *IEEE Transactions on Industrial Informatics*, IEEE, v. 14, n. 10, p. 4674–4682, 2018. Disponível em: <https://doi.org/10.1109/TII.2018.2855198>.
- 68 MATTHYSSENS, P. Reconceptualizing value innovation for Industry 4.0 and the Industrial Internet of Things. *Journal of Business & Industrial Marketing*, Emerald Publishing Limited, 2019. Disponível em: <https://doi.org/10.1108/JBIM-11-2018-0348>.
- 69 BOYES, H.; HALLAQ, B.; CUNNINGHAM, J.; WATSON, T. The industrial internet of things (IIoT): An analysis framework. *Computers in industry*, Elsevier, v. 101, p. 1–12, 2018. Disponível em: <https://doi.org/10.1016/j.compind.2018.04.015>.
- 70 GARRIDO-HIDALGO, C.; OLIVARES, T.; RAMIREZ, F. J.; RODA-SANCHEZ, L. An end-to-end Internet of Things solution for reverse supply chain management in industry 4.0. *Computers in Industry*, Elsevier, v. 112, p. 103127, 2019. Disponível em: <https://doi.org/10.1016/j.compind.2019.103127>.
- 71 HANSCH, G.; SCHNEIDER, P.; FISCHER, K.; BÖTTINGER, K. A unified architecture for industrial iot security requirements in open platform communications. *2019 24th IEEE International Conference on Emerging Technologies and Factory Automation (ETFA)*, p. 325–332, 2019. Disponível em: <https://doi.org/10.1109/ETFA.2019.8869524>.
- 72 RICHARDSON, C. *Microservices Patterns: With examples in Java*. Manning Publications, 2018. ISBN 9781617294549. Disponível em: <https://books.google.com.br/books?id=UeK1swEACAAJ>.
- 73 SUBRAMANIAN, H.; RAJ, P. *Hands-On RESTful API Design Patterns and Best Practices: Design, develop, and deploy highly adaptable, scalable, and secure RESTful web APIs*. [S.l.]: Packt Publishing Ltd, 2019.

- 74 PIVOTAL, S. Spring. *Spring Framework*: <https://spring.io/projects/spring-framework>, 2002.
- 75 PALLETS. Flask: web application framework Electronic resource. *Disponível em*: <https://flask.palletsprojects.com/en/1.1.x/>, 2019.
- 76 SHEFFER, Y.; HARDT, D.; JONES, M. JSON Web Token Best Current Practice. *Disponível em*: <https://www.rfc-editor.org/info/rfc8725>, 2020. *Disponível em*: <https://doi.org/10.17487/RFC8725>.
- 77 CORPORATION, A. *TIFF, Revision 6.0*. [S.l.]: <https://www.itu.int/itudoc/itu-t/com16/tiff-fx/docs/tiff6.pdf>, 1992.
- 78 INC, D. Overview of docker compose. *Disponível em*: <https://docs.docker.com/compose/>, 2013.
- 79 MOSBERGER-TANG, D.; BECK, A. Sane - scanner access now easy. *Disponível em*: <http://www.sane-project.org/>.
- 80 BRADSKI, G.; KAEHLER, A. *Learning OpenCV: Computer Vision with the OpenCV Library*. [S.l.]: O'Reilly, 2008. ISBN 0596516134.
- 81 CHOLLET, F.; OTHERS. *Keras: The Python Deep Learning library*. 2018. ascl:1806.022 p. Provided by the SAO/NASA Astrophysics Data System. *Disponível em*: <https://ui.adsabs.harvard.edu/abs/2018ascl.soft06022C>.
- 82 REDMON, J.; FARHADI, A. Yolov3: An incremental improvement. *arXiv*, 2018.
- 83 PYTHON. *Cryptography*. 2019. *Disponível em*: <https://pypi.org/project/cryptography/>.
- 84 SHARMA, S. *Mastering Microservices with Java: Build Enterprise Microservices with Spring Boot 2.0, Spring Cloud, and Angular*. [S.l.]: Packt Publishing Ltd, 2019.
- 85 RESCORLA, E. *The Transport Layer Security (TLS) Protocol Version 1.3*. 2018. *Disponível em*: "<https://tools.ietf.org/html/rfc8446>".
- 86 ASTHANA ABHINAV; SOBTI, A.; KANE, A. Cfssl: Cloudflare's pki and tls toolkit. <https://github.com/cloudflare/cfssl>, 2015.
- 87 ITI. Infraestrutura de chaves públicas brasileira – ICP-Brasil. *Disponível em*: <https://antigo.iti.gov.br/icp-brasil>, 2017.
- 88 LI, S.; WANG, N.; DU, X.; LIU, A. Internet Web Trust System Based on Smart Contract. *International Conference of Pioneering Computer Scientists, Engineers and Educators*, p. 295–311, 2019. *Disponível em*: [https://doi.org/10.1007/978-981-15-0118-0\\_23](https://doi.org/10.1007/978-981-15-0118-0_23).
- 89 JOHNSON, D.; MENEZES, A.; VANSTONE, S. A. The Elliptic Curve Digital Signature Algorithm (ECDSA). *Int. J. Inf. Sec.*, <http://dblp.uni-trier.de/db/journals/ijisec/ijisec1.html#JohnsonMV01>, v. 1, 2001. *Disponível em*: <https://doi.org/10.1007/s102070100002>.
- 90 TONG, Z. Elasticsearch index. *Disponível em*: <https://www.elastic.co/pt/what-is/elk-stack>.
- 91 KEMENADE, J. V. *The CERN Digital Memory Platform: Building a CERN scale OAIS compliant Archival Service*. Dissertação (Mestrado) — Vrije Universiteit Amsterdam and Universiteit van Amsterdam, <https://cds.cern.ch/record/2728246/files/CERN-THESIS-2020-092.pdf>. CERN-THESIS-2020-092, 2020.
- 92 ARTEFACTUAL. Manual normalization. *Disponível em*: <https://www.ar-chivematica.org/en/docs/archivematica-1.12/user-manual/ingest/manual-normalization/>, 2020.

- 93 VOINOV, N.; DROBINTSEV, P.; KOTLYAROV, V.; NIKIFOROV, I. Distributed OAIS-based digital preservation system with HDFS technology. *2017 20th Conference of Open Innovations Association (FRUCT)*, p. 491–497, 2017. Disponível em: <https://doi.org/10.23919/FRUCT.2017.8071353>.
- 94 ROSENTHAL, D.; REICH, V. Lockss, a permanent web publishing and access system: Brief introduction and status report. *Serials*, UKSG in association with Ubiquity Press, v. 14, n. 3, 2003. Disponível em: <https://doi.org/10.1629/14239>.
- 95 ARTEFACTUAL. *Archival storage*. 2009. Disponível em: <https://www.archivemata.org/en/docs/archivemata-1.8/user-manual/archival-storage/archival-storage/#aip-encryption>.
- 96 FEDERATION, D. L. *Metadata Encoding and Transmission Standard: Primer and Reference Manual*. 2007. Disponível em: <http://www.loc.gov/standards/mets/mets-schemadocs.html>.
- 97 CAPLAN, P. Understanding premis. In: LIBRARY OF CONGRESS WASHINGTON DC, USA. 2017. Disponível em: <https://www.loc.gov/standards/premis/understanding-premis-rev2017.pdf>.
- 98 ISO. ISO 15836-1:2017 information and documentation – the dublin core metadata element set – part 1: Core elements. *International Organization for Standardization.*, p. 1–7, 2017.
- 99 FANG, L.; BITAR, N.; ZHANG, R.; DAIKOKU, M.; PAN, P. *The BagIt File Packaging Format*. ISSN: 2070-1721. 2013. Disponível em: <https://datatracker.ietf.org/doc/html/rfc8493>.
- 100 CALLAS, J.; DONNERHACKE, L.; FINNEY, H.; SHAW, D.; THAYER, R. *OpenPGP Message Format (RFC 4880)*. 2007. Disponível em: <https://tools.ietf.org/html/rfc4880>.
- 101 SHINTAKU, M.; ABREU, J. P. L. de; Santarem Segundo, J. E.; CASTRO, P. d. P. *Guia de usuário do AtoM*. ISBN: 97870131270. [S.l.]: Instituto Brasileiro de Informação em Ciência e Tecnologia (IBICT), 2017. 164 p.
- 102 SKALSKI, P. *Make Sense*. 2019. Disponível em: <https://github.com/SkalskiP/make-sense/>.
- 103 KUMAR, N.; RAMDOSS, Y.; ORZACH, Y. *Network Analysis Using Wireshark 2 Cookbook: Practical recipes to analyze and secure your network using Wireshark 2*. [S.l.]: Packt Publishing Ltd. ISBN: 978-1-78646-167-4, 2018. 581 p.