



CENTRO UNIVERSITÁRIO DE BRASÍLIA
CIÊNCIA DA COMPUTAÇÃO

FELIPE BARCELOS DE CARVALHO (RA: 22350044)
MARCOS VINICIUS ROCHA (RA: 22352865)
JOÃO MARCELO GUIMARÃES DOURADO (RA: 22350653)
EDUARDO ARAÚJO UCHOA (RA: 22353207)
DAVI MAIA (RA: 22305561)

PROJECT AEGIS:
PAINEL DE INTELIGÊNCIA DE AMEAÇAS

BRASÍLIA
2025

Entregável 2: Scripts de ETL (Extração, Transformação e Carga)

Documentação dos Processos de ETL

1. Visão Geral do Pipeline O processo de ETL é executado automaticamente via GitHub Actions. Os scripts Python atuam como *crawlers* de APIs públicas, normalizam os dados e salvam o resultado em arquivos estáticos prontos para consumo pelo Frontend.

2. Script: Coletor CVE (coletor_cve.py)

- **Fonte:** NVD NIST API (National Vulnerability Database).
- **Transformação:** Filtra vulnerabilidades recentes; classifica automaticamente o tipo de falha (ex: SQL Injection, RCE) com base na descrição técnica e simplifica o score CVSS.
- **Carga:** Gera o arquivo site/cve_kpis.json.

```
database > ⌘ coletor_cve.py > ...
 1  # Importar bibliotecas necessárias
 2  import pandas as pd
 3  import requests
 4  from datetime import datetime, timedelta
 5  import time
 6  import json
 7  import os
 8
 9  # =====
10 # 0. CONFIGURAÇÃO DA API (NVD)
11 # =====
12 API_KEY = "49e988e1-7235-47d9-8f56-a2e7e4cae151"
13 API_URL = "https://services.nvd.nist.gov/rest/json/cves/2.0"
14 headers = {'apiKey': API_KEY}
15
16 # --- FUNÇÃO DE CLASSIFICAÇÃO ---
17 def classificar_vulnerabilidade(descricao):
18     """
19         Lê a descrição técnica e retorna uma categoria simplificada.
20     """
21     if not descricao:
22         return 'Outros'
23
24     desc = descricao.lower()
25
26     if 'remote code execution' in desc or 'execute arbitrary code' in desc:
27         return 'RCE (Execução Remota)'
28     elif 'denial of service' in desc or 'dos' in desc:
29         return 'DoS (Negação de Serviço)'
30     elif 'sql injection' in desc:
31         return 'SQL Injection'
32     elif 'cross-site scripting' in desc or 'xss' in desc:
33         return 'XSS'
34     elif 'privilege escalation' in desc or 'gain privileges' in desc:
35         return 'Escalada de Privilégio'
36     elif 'buffer overflow' in desc:
37         return 'Buffer Overflow'
38     else:
39         return 'Outros'
40
41 print("*"*80)
42 print("ATRIBUIÇÃO: Este produto usa dados da NVD API, mas não é endossado ou certificado pela NVD.")
43 print("*"*80 + "\n")
```

3. Script: Coletor de Vazamentos (coletor_vazamentos.py)

- **Fonte:** Have I Been Pwned API.
- **Transformação:** Ordena os vazamentos por data de divulgação; formata números grandes de contas afetadas; anonimiza dados sensíveis e padroniza datas.
- **Carga:** Gera o arquivo site/hibp_kpis.json.

```
database > ↵ coletor_vazamentos.py > ...
1  # Importar bibliotecas necessárias
2  import pandas as pd
3  import requests
4  import time
5  import json
6  import warnings
7
8  # Suprimir avisos de "UserWarning"
9  warnings.filterwarnings('ignore', category=UserWarning)
10
11 # =====
12 # 0. CONFIGURAÇÃO DA API (Have I Been Pwned - HIBP)
13 # =====
14 headers = {
15     'User-Agent': 'Project-Aegis-Coletor-Academico'
16 }
17 API_URL = "https://haveibeenpwned.com/api/v3/breaches"
18 SLEEP_TIME = 2
19
20 print("*80)
21 print("ATRIBUIÇÃO: Usando a Base de Dados 'Have I Been Pwned' (HIBP) v3.")
22 print("*80 + "\n")
23
24 # =====
25 # 1. DICIONÁRIO DE DADOS (NOSSO ALVO)
26 # =====
27 dicionario_dados_vazamento = {
28     "nome_vazamento": "String (Nome da empresa/site, ex: 'SocialNet')",
29     "data_vazamento": "Date (Data em que o vazamento ocorreu)",
30     "contas_afetadas": "Integer (Número total de contas vazadas)",
31     "setor": "String (Indústria da empresa, ex: 'Technology')"
32 }
33 print("--- 1. ESTRUTURA DO DICIONÁRIO DE DADOS (NOSSO ALVO) ---")
34 for chave, valor in dicionario_dados_vazamento.items():
35     print(f"{chave}: {valor}")
36 print("\n" + "*50 + "\n")
37
```

4. Script: Coletor de Geolocalização (coletor_paises.py)

- **Fonte:** AbuseIPDB (Blacklist de IPs maliciosos).
- **Transformação:** Consulta a lista de IPs reportados recentemente, realiza a geolocalização (IP -> País) e agrupa a contagem de ameaças por nação para alimentar o gráfico.
- **Carga:** Gera o arquivo site/paises_kpis.json.

```
database > ⚡ coletor_paises.py > ...
1  # Importar bibliotecas necessárias
2  import pandas as pd
3  import requests
4  import json
5  import time
6  from datetime import datetime
7
8  # =====
9  # 0. CONFIGURAÇÃO DA API (AbuseIPDB)
10 # =====
11 API_KEY = "4947aff3d5a68c6db3a3d38206e061b8df8033400f2c14c332acf48c2ae8bc0b0924798870703a3"
12
13 if "SUA_CHAVE" in API_KEY:
14     print("ERRO: Adicione sua chave de API.")
15     exit()
16
17 headers = {
18     'Key': API_KEY,
19     'Accept': 'application/json',
20 }
21 API_URL = "https://api.abuseipdb.com/api/v2/blacklist"
22 # Limitamos a 10.000 IPs (máximo gratuito)
23 PARAMS = {
24     'confidenceMinimum': 50,
25     'limit': 10000
26 }
27
28 # --- DICIONÁRIO DE MAPEAMENTO (Country Code -> Country Name) ---
29 COUNTRY_MAP = {
30     "US": "United States", "CN": "China", "RU": "Russia", "DE": "Germany",
31     "NL": "Netherlands", "FR": "France", "GB": "United Kingdom",
32     "IN": "India", "BR": "Brazil", "VN": "Vietnam", "UA": "Ukraine",
33     "KR": "South Korea", "JP": "Japan", "CA": "Canada", "AU": "Australia",
34     "PL": "Poland", "TR": "Turkey", "HK": "Hong Kong", "SG": "Singapore",
35     "TW": "Taiwan", "IR": "Iran", "RO": "Romania", "CZ": "Czech Republic",
36     "ID": "Indonesia", "ES": "Spain", "IT": "Italy", "TH": "Thailand"
37 }
38
39 print("*"*80)
40 print("ATRIBUIÇÃO: Usando a Base de Dados 'AbuseIPDB'.")
41 print("*"*80 + "\n")
42
```

5. Script: Coletor OTX (coletor_otx.py)

- **Fonte:** AlienVault Open Threat Exchange.
- **Transformação:** Agrega "Pulsos" de ameaças globais; remove duplicatas e extrai tags de categorização (ex: Phishing, Malware) para análise de tendências.
- **Carga:** Gera o arquivo site/otx_kpis.json.

```
database > ⇧ coletor_otx.py > ...
1  # Importar bibliotecas necessárias
2  import pandas as pd
3  import requests
4  import json
5  from datetime import datetime, timedelta
6  import time
7
8  # =====
9  # 0. CONFIGURAÇÃO DA API (AlienVault OTX)
10 # =====
11 API_KEY = "6dc4519f2fe20cb267b376d6f626b0ee64243693f27d30109f3e6b1802e78c2d"
12
13 if "SUA_CHAVE_OTX_AQUI" in API_KEY:
14     print("ERRO: Adicione sua chave de API da OTX.")
15     exit()
16
17 headers = {
18     'X-OTX-API-KEY': API_KEY
19 }
20
21 # Usando /subscribed pois mostrou-se mais estável que /activity
22 API_URL = "https://otx.alienvault.com/api/v1/pulses/subscribed"
23
24 # Configurações de Coleta
25 RESULTS_PER_PAGE = 50
26 MAX_PAGES = 200 # Limite de segurança para não rodar infinito (200 * 50 = 10.000 pulsos)
27 DIAS_PARA_BUSCAR = 365 # Queremos 1 ano de dados
28
29 print("*"*80)
30 print("ATRIBUIÇÃO: Usando a Base de Dados 'AlienVault OTX' (Open Threat Exchange).")
31 print("*"*80 + "\n")
32
33 # =====
34 # 1. COLETA REAL DE DADOS (ETL - 365 DIAS)
35 # =====
36 print("--- 2. COLETANDO DADOS REAIS DA API DA OTX (Janela de 1 Ano) ---")
37
38 # Calcular a data limite (1 ano atrás)
39 data_limite = datetime.now() - timedelta(days=DIAS_PARA_BUSCAR)
40 print(f"Buscando dados a partir de: {data_limite.strftime('%Y-%m-%d')}]")
41
42 dados_coletados = []
43 parar_coleta = False
```