

# Construção de um Conjunto de Dados para Análise Estática de Ransomwares

Marcelo M. Borges<sup>1</sup>, Arthur P. Labaki<sup>1</sup>, Renan G. Cattelan<sup>1</sup>, Rodrigo S. Miani<sup>1</sup>

<sup>1</sup>Faculdade de Computação – Universidade Federal de Uberlândia (UFU)  
Av. João Naves de Ávila, 2121 – 38.400-902 – Uberlândia – MG – Brazil

{marcelo.m.borges, arthur.labaki, renan, miani}@ufu.br

**Abstract.** *This paper presents the construction of a data set for static analysis of ransomwares. The resulting database consists of 338 PE files, separated into 21 different families and containing general information about the files and features of the ransomwares analyzed. The set produced, here made publicly available, can be used to identify and classify families of ransomware.*

**Resumo.** *Este trabalho apresenta a construção de um conjunto de dados para análise estática de ransomwares. A base de dados resultante consiste de 338 arquivos PE, separados em 21 diferentes famílias e contendo informações gerais dos arquivos e características dos ransomwares analisados. O conjunto produzido, aqui disponibilizado publicamente, pode ser utilizado para identificação e classificação de famílias de ransomware.*

## 1. Introdução

*Malware* ou software malicioso (em Inglês, *malicious software*) são programas desenvolvidos para infectar computadores, causando qualquer tipo de danos a usuários, equipamentos ou redes de comunicação. São, muitas vezes, usados como intermediários para a prática de golpes, a realização de ataques e a disseminação de *spam* [Sikorski and Honig 2012]. No Brasil, a criação de qualquer tipo de *malware* é considerado um crime digital. Os principais motivos que levam ao desenvolvimento e à propagação de códigos maliciosos são a obtenção de vantagens financeiras, a coleta de informações confidenciais, o desejo de autopromoção e o vandalismo.

Ataques de uma classe específica de software malicioso, conhecida como *ransomware*, têm se disseminado ao redor do mundo, prejudicando negócios, empresas e governos. No Brasil, a situação não é diferente e, de acordo com relatório da FortiGuard Labs [Fortinet 2020], o país ultrapassou a marca de 8,4 bilhões de ataques no ano de 2020. O mesmo relatório alerta ainda que tais ataques tendem a aumentar em 2021 e que os cibercriminosos estão empregando tecnologias cada vez mais avançadas, como inteligência artificial, em ataques direcionados e com alto grau de sofisticação e eficiência.

Esse cenário evidencia a importância da pesquisa em cibersegurança para a área de Sistemas de Informação e a necessidade de dados para melhor compreender e permitir combater esse tipo de ameaça. No entanto, a disponibilidade de dados para análise sobre *ransomware* ainda é limitada e, por isso, o objetivo deste trabalho foi aplicar técnicas de análise estática, seguindo uma metodologia bem definida, para criar um conjunto de dados robusto e confiável, com informações sobre famílias de *ransomware*, e disponibilizá-lo para a comunidade científica.

## 2. Fundamentação Teórica

Para identificar, classificar e eliminar um *malware*, é necessário analisar suas características, de modo a permitir o desenvolvimento de medidas preventivas. Com isso, alguns especialistas começaram a classificar os *malwares*, agrupando-os de acordo com as ações que realizam no hospedeiro. *Ransomware* é um tipo de *malware* que busca impedir o acesso do usuário ao sistema ou aos seus arquivos, normalmente por meio da cifragem de dados. Na grande maioria dos casos, o sistema só é liberado após o pagamento de um resgate.

O desenvolvimento de métodos de prevenção e remoção dos *malwares*, em geral, e dos *ransomwares*, em específico, depende da análise de suas características por pesquisadores e especialistas em cibersegurança. Essa análise tem como objetivo entender como o software malicioso funciona e pode ser realizada de forma dinâmica ou estática. Na análise dinâmica, o *malware* é monitorado durante sua execução, por meio de emuladores, *debuggers*, ambientes virtuais, diferentes ferramentas para monitoração de processos para capturar interações com o sistema e sua rede [Filho et al. 2010]. Já na análise estática, são extraídas informações gerais e características do *malware* sem executá-lo, através de análise de *strings*, *disassembling* e engenharia reversa, por exemplo [Filho et al. 2010]. Este trabalho tem como foco a análise estática de *ransomware*.

Diversos trabalhos propõem modelos de detecção de *ransomware* usando análise estática [Zhang et al. 2019] e [Egunjobi et al. 2019]. Contudo, poucos fornecem os conjuntos de dados usados para a validação dos modelos. Essa prática dificulta a reprodução de experimentos e comparação com outros trabalhos relacionados. Dois estudos se destacam neste aspecto. [Anderson and Roth 2018] e [Lashkari et al. 2018] propõem conjuntos de dados de análise estática para identificação e classificação de *malware*. Enquanto no trabalho de [Anderson and Roth 2018] não constam amostras de *ransomware*, em [Lashkari et al. 2018] só existem algumas amostras disponíveis para o sistema operacional Android.

## 3. Metodologia de Pesquisa

O processo de construção do conjunto de dados (Figura 1) iniciou-se pela Etapa 1, com a obtenção de diversas amostras de *malware* em formato executável a partir de uma ampla busca por múltiplos repositórios, especialmente o Virus Share<sup>1</sup>, que é um grande repositório de amostras de códigos maliciosos disponibilizado a pesquisadores. Ainda nessa etapa, foi empregada também a aplicação Virus Total<sup>2</sup>, que verifica a amostra utilizando diversos antivírus e *scanners*, assim confirmando se o *malware* capturado é, de fato, um *ransomware*.

Na Etapa 2, as amostras de *malware* com seus respectivos hashes MD5 foram renomeadas e foi então criada uma tabela de referência, contendo algumas informações dos *ransomwares* para maior controle. Nessa tabela, foi empregada novamente a aplicação Virus Total para se obter informações como *hash*, família de *ransomware* e a plataforma atuante. Também foi incluído, na tabela, o repositório em que a amostra foi obtida.

Para extração das características, na Etapa 3, foi utilizada a mesma metodologia aplicada no conjunto de dados EMBER [Anderson and Roth 2018], que consiste em usar

---

<sup>1</sup><https://virusshare.com>

<sup>2</sup><https://www.virustotal.com>

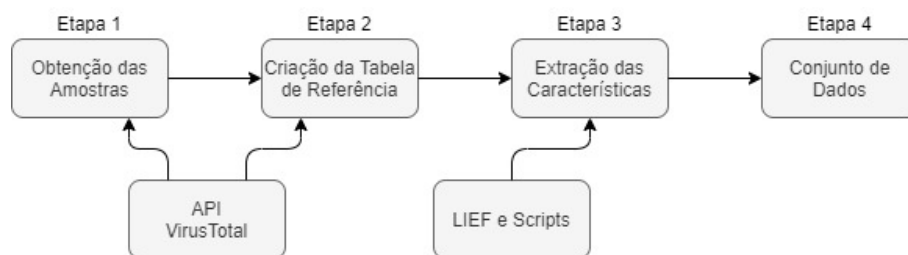


Figura 1. Etapas para a construção da base de dados

a biblioteca LIEF<sup>3</sup> (*Library to Instrument Executable Formats*) para analisar, modificar e abstrair arquivos executáveis do tipo *Portable Executable* (PE). Para isso, foram utilizados alguns *scripts* disponibilizados pelo próprio EMBER, que foram adaptados para que usassem as amostras ora geradas, que eram então convertidas de características brutas para um arquivo no formato JSON, chegando-se assim, na Etapa 4, ao conjunto de dados resultante, detalhado a seguir.

#### 4. Resultados

A base de dados resultante consiste de 338 arquivos PE. Esses arquivos são separados em 21 famílias diferentes, entre elas se destaca a *WannaCry*, com 57 amostras, que atraiu bastante visibilidade a esse tipo de ataque, com grande quantidade de ataques bem sucedidos no ano de 2017. Para mais informações sobre as famílias encontradas, veja o sumário no repositório disponibilizado.

A tabela de referência, cuja uma pequena fração é exemplificada na Tabela 1, foi definida com base em informações gerais sobre esses arquivos, com cinco colunas:

1. MD5: que consiste do *hash* MD5 do arquivo, servindo também como identificador de cada arquivo da base de dados;
2. *Family*: que representa a família a qual cada *ransomware* pertence;
3. *Platform*: que caracteriza do ambiente alvo de ação de cada *ransomware*, predominantemente Microsoft Windows;
4. *Sample by*: que equivale ao local (site) no qual o arquivo foi obtido; e
5. LIEF: simbolizando sucesso (V) ou fracasso (X) no processo de obtenção de características pelo LIEF.

MD5	Family	Platform	Sample by	LIEF
fbdbc39af1139aebba4da004475e8839	Bad Rabbit	Microsoft Windows (.EXE (Dropper), DLL (Trojan))	VirusShare	V
fec5a0d4dea87955c124f2eaa1f759f5	CryptoLocker	Microsoft Windows (.EXE)	VirusShare	V
796fdae3b1476ed20cdac74ca9d40973	Cryptowall	Microsoft Windows (.EXE)	VirusShare	V
7a9807d121aa0721671477101777cb34	GandCrab	Microsoft Windows (.EXE)	VirusShare	V
b9a84d52093a20975d44418e9eac631	Globeimposter	Microsoft Windows (.EXE)	VirusShare	V

Tabela 1. Fração da Tabela de Referência

As características extraídas são baseadas na própria estrutura de um arquivo PE. Tal formato é dividido em diversos tipos de cabeçalhos como *DOS Header*, *DOS Stub*, *COFF Header* e outros cabeçalhos opcionais. Seguindo lógica semelhante ao trabalho de [Anderson and Roth 2018], o conjunto de atributos extraídos pelo LIEF, diretamente do PE, é dividido em quatro grupos de características:

<sup>3</sup><https://lief.quarkslab.com>

1. Informações gerais do arquivo: engloba características comuns do arquivo como o tamanho e número de funções importadas e exportadas;
2. Informações de cabeçalho: reúne informações obtidas do *COFF Header* como a marcação de tempo (*timestamp*) e máquina alvo. De cabeçalhos opcionais, se obtém informações como características de DLL (bibliotecas), versões de sistemas e subsistemas, assinaturas e certificados;
3. Funções importadas: conjunto de funções importadas pelo PE separadas por biblioteca (DLL);
4. Informações de sessão: agrupa propriedades de cada sessão incluindo nome, tamanho, entropia, tamanho virtual e descrição por *string* das características.

Estão disponíveis, em um [repositório público](#), a tabela de referência, o sumário de famílias de *ransomwares* e os atributos extraídos pela ferramenta LIEF.

## 5. Conclusão

Este trabalho apresentou a construção de um conjunto de dados para análise estática de *ransomwares*. O conjunto produzido pode ser utilizado para diversas tarefas de classificação, como a identificação de amostras de *ransomware* dentre outros tipos de *malware* ou ainda classificar famílias de *ransomware*. A próxima etapa do trabalho envolve a criação de modelos de classificação para tais tarefas usando algoritmos de aprendizado de máquina.

## Referências

- [Anderson and Roth 2018] Anderson, H. S. and Roth, P. (2018). EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models. *ArXiv e-prints*.
- [Egunjobi et al. 2019] Egunjobi, S., Parkinson, S., and Crampton, A. (2019). Classifying ransomware using machine learning algorithms. In *International Conference on Intelligent Data Engineering and Automated Learning*, pages 45–52. Springer.
- [Filho et al. 2010] Filho, D. S. F., Grégio, A. R. A., Afonso, V. M., d. Santos, R. D. C., Jino, M., and de Geus, P. L. (2010). Análise comportamental de código malicioso através da monitoração de chamadas de sistema e tráfego de rede. In *Simpósio Brasileiro em Segurança da Informação e de Sistemas Computacionais*, pages 311–324.
- [Fortinet 2020] Fortinet (2020). A América Latina sofreu mais de 41 bilhões de tentativas de ataques cibernéticos em 2020. *Online: <https://www.fortinet.com/>*.
- [Lashkari et al. 2018] Lashkari, A. H., Kadir, A. F. A., Taheri, L., and Ghorbani, A. A. (2018). Toward developing a systematic approach to generate benchmark android malware datasets and classification. In *International Carnahan Conference on Security Technology*, pages 1–7.
- [Sikorski and Honig 2012] Sikorski, M. and Honig, A. (2012). *Practical malware analysis: the hands-on guide to dissecting malicious software*. No Starch Press.
- [Zhang et al. 2019] Zhang, H., Xiao, X., Mercaldo, F., Ni, S., Martinelli, F., and Sangaiah, A. K. (2019). Classification of ransomware families with machine learning based on n-gram of opcodes. *Future Generation Computer Systems*, 90:211–221.