

UNIVERSIDADE TECNOLÓGICA FEDERAL DO PARANÁ  
DEPARTAMENTO ACADÊMICO DE COMPUTAÇÃO  
CURSO DE BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

NELSON VIEIRA DA SILVA JÚNIOR

**ESTUDO E ANÁLISE DOS GRUPOS HACKERS QUE  
REALIZAM DESFIGURAÇÃO DE PÁGINAS WEB  
NO BRASIL**

MONOGRAFIA

CAMPO MOURÃO

2017

**NELSON VIEIRA DA SILVA JÚNIOR**

**ESTUDO E ANÁLISE DOS GRUPOS HACKERS QUE  
REALIZAM DESFIGURAÇÃO DE PÁGINAS WEB  
NO BRASIL**

Trabalho de Conclusão de Curso de graduação apresentado à disciplina de Trabalho de Conclusão de Curso 2, do Curso de Bacharelado em Ciência da Computação do Departamento Acadêmico de Computação da Universidade Tecnológica Federal do Paraná, como requisito parcial para obtenção do título de Bacharel em Ciência da Computação.

Orientador: Prof. Dr. Rodrigo Campiolo

**CAMPO MOURÃO**

**2017**



## ATA DE DEFESA DO TRABALHO DE CONCLUSÃO DE CURSO

Às **16:30** do dia **28 de novembro de 2017** foi realizada na sala **E102** da UTFPR-CM a sessão pública da defesa do Trabalho de Conclusão do Curso de Bacharelado em Ciência da Computação do(a) acadêmico(a) **Nelson Vieira da Silva Júnior** com o título **Estudo e análise de grupos hackers que realizam desfiguração de páginas Web no Brasil**. Estavam presentes, além do(a) acadêmico(a), os membros da banca examinadora composta por: **Prof. Dr. Rodrigo Campiolo** (orientador), **Prof. Dr. Luiz Arthur Feitosa dos Santos** e **Prof. Dr. Rogério Aparecido Gonçalves**. Inicialmente, o(a) acadêmico(a) fez a apresentação do seu trabalho, sendo, em seguida, arguido(a) pela banca examinadora. Após as arguições, sem a presença do(a) acadêmico(a), a banca examinadora o(a) considerou \_\_\_\_\_ na disciplina de Trabalho de Conclusão de Curso **2** e atribuiu, em consenso, a nota \_\_\_\_\_ (\_\_\_\_\_). Este resultado foi comunicado ao(à) acadêmico(a) e aos presentes na sessão pública. A banca examinadora também comunicou ao acadêmico(a) que este resultado fica condicionado à entrega da versão final dentro dos padrões e da documentação exigida pela UTFPR ao professor Responsável do TCC no prazo de **onze dias**. Em seguida foi encerrada a sessão e, para constar, foi lavrada a presente Ata que segue assinada pelos membros da banca examinadora, após lida e considerada conforme.

Observações: \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

Campo Mourão, **28 de novembro de 2017**

\_\_\_\_\_  
**Prof. Dr. Luiz Arthur Feitosa dos Santos**  
Membro 1

\_\_\_\_\_  
**Prof. Dr. Rogério Aparecido Gonçalves**  
Membro 2

\_\_\_\_\_  
**Prof. Dr. Rodrigo Campiolo**  
Orientador

**A ata de defesa assinada encontra-se na coordenação do curso.**

# Resumo

---

Vieira, Nelson. Estudo e análise dos grupos hackers que realizam desfiguração de páginas Web no Brasil. 2017. 57. f. Monografia (Curso de Bacharelado em Ciência da Computação), Universidade Tecnológica Federal do Paraná. Campo Mourão, 2017.

Desfiguração de páginas é uma forma comum de ataque em sítios Web, onde o conteúdo do sítio é totalmente ou parcialmente substituído por um atacante que inclui conteúdos como, textos, imagens e formas de assinatura do invasor ou manifestações ofensivas direcionada a política ou causas sociais. As consequências da desfiguração pode variar entre financeira e moral, portanto, investigar essas ações e registrar essas ocorrências são de grande valia para a segurança computacional, tornando-se medidas de prevenção ou alertas antecipados. Este trabalho tem como objetivo investigar grupos hackers direcionados a desfiguração de páginas no Brasil, com o intuito de extrair informações de inteligência para o monitoramento desses grupos e de padrões usados nas desfigurações que possibilitem o desenvolvimento de mecanismos de detecção automática e sistemas de alertas antecipados. Para isso foi coletado o HTML dos sítios que sofreram desfiguração e que foram registrados no Zone-H criando uma base de dados para que a mesma seja processada através de técnicas de mineração de texto. Através do processamento do HTML e da caracterização da base de dados, foi possível extrair os principais grupos/indivíduos que atuam no Brasil, caracterizando os domínios que esses tendem a atacar. Além disso foi obtido os horários e dias da semana que mais ocorrem desfiguração, foi extraído de redes sociais e confirmou-se que os ataques são divulgados pelas redes. Padrões individuais e gerais foram extraídos, obtendo também os principais termos utilizados nas desfigurações. Conclui-se que através das técnicas utilizadas foi possível evidenciar os grupos/indivíduos que mais atuam no Brasil e que a partir disso foi possível criar uma base de informações que poderá ser utilizada por mecanismos de detecção e de sistema antecipado de alerta.

**Palavras-chaves:** Zone-H. Web Crawler. Mineração de Texto

# Abstract

---

Vieira, Nelson. Study and analysis of hacker's groups that perform defacement of Web pages in Brazil. 2017. 57. f. Monograph (Undergraduate Program in Computer Science), Federal University of Technology – Paraná. Campo Mourão, PR, Brazil, 2017.

Deface or defacement is a common form of attack on Web sites, where the content of the site is totally or partially replaced by an attacker that includes content such as, texts, images and signature forms of the attacker or manifestations offensive policies or social causes. The consequences of page defacement may vary between financial and moral, so investigating such actions and recording such occurrences are of great value to computer security as this can become preventive measures or early warning. It aims to investigate hacker groups aimed at the disfigurement of pages in Brazil, in order to extract intelligence information for the monitoring of these groups and patterns used in the disfigurements that allow the development of mechanisms of automatic detection and early warning systems. For this, the HTML of the sites that have been disfigured were collected and recorded in Zone-H, creating a database to be processed through text mining techniques. Through the processing of the HTML and the characterization of the database, it was possible to extract the main groups/individuals that operate in Brazil, characterizing the domains that these tend to attack. In addition it was obtained the schedules and days of the week that most occur disfigurement, was extracted from social networks and it was confirmed that the attacks are released by the networks. Individual and general standards were extracted, also obtaining the main terms used in the disfigurements. It is concluded that through the techniques used it was possible to highlight the groups/individuals that most act in Brazil and that from this it was possible to create a base of information that could be used by detection mechanisms and early warning system.

**Keywords:** Zone-H. Web-crawler. Text Mining

# Lista de figuras

---

2.1	Desfigurações especiais no Zone-H . . . . .	10
2.2	Lista de desfiguração de páginas no Zone-H . . . . .	11
2.3	Exemplo de injeção SQL . . . . .	12
2.4	Exemplo de XSS armazenado. . . . .	13
2.5	Exemplo de lematização . . . . .	16
2.6	Análise de dados de cibersegurança obtidos de fontes de dados não estruturados	18
2.7	Proposta de arcabouço para fóruns hackers . . . . .	19
3.1	Fluxograma do método de pesquisa. . . . .	22
4.1	Gráfico caracterizando domínio por invasor/grupo . . . . .	34
4.2	Gráfico de desfiguração por dias da semana . . . . .	35
4.3	Gráfico de desfiguração por horário . . . . .	35
4.4	Padrão de desfiguração do grupo BRLZPoC . . . . .	38
4.5	Termos mais frequentes nas desfigurações . . . . .	38

# Sumário

---

<b>1</b>	<b>Introdução</b>	<b>6</b>
<b>2</b>	<b>Referencial Teórico</b>	<b>8</b>
2.1	Desfiguração de páginas . . . . .	9
2.1.1	Definição . . . . .	9
2.1.2	Desfiguração de páginas no Brasil . . . . .	9
2.1.3	Técnicas utilizadas para desfiguração . . . . .	11
2.1.4	Técnicas de prevenção . . . . .	14
2.2	Sítios de espelhamento de desfigurações de página . . . . .	14
2.3	Mineração de Texto . . . . .	15
2.3.1	Tokenização . . . . .	15
2.3.2	Lematização e radicalização . . . . .	16
2.3.3	Associação de palavras . . . . .	16
2.3.4	Reconhecimento de entidades nomeadas . . . . .	17
2.4	Trabalhos relacionados . . . . .	17
2.5	Considerações do capítulo . . . . .	20
<b>3</b>	<b>Método de Pesquisa</b>	<b>21</b>
3.1	Questões de pesquisa . . . . .	21
3.2	Materiais e métodos . . . . .	22
3.2.1	Zone-H . . . . .	23
3.2.2	Coletor . . . . .	23
3.2.3	Dados coletados . . . . .	24

3.2.4	Estatísticas . . . . .	24
3.2.5	Espelho do sítio desfigurado . . . . .	26
3.2.6	Mineração de texto . . . . .	27
3.2.7	Análise especialista . . . . .	29
3.2.8	Padrões, características e perfis . . . . .	30
3.2.9	Bases de inteligência / Perfis de monitoramento . . . . .	30
3.3	Considerações do capítulo . . . . .	31
<b>4</b>	<b>Resultados e Discussões</b>	<b>32</b>
4.1	Caracterização dos resultados . . . . .	32
4.2	Investigação de redes sociais . . . . .	36
4.3	Análise de características e padrões de ataques . . . . .	36
4.4	Avaliação das questões de pesquisa . . . . .	40
4.5	Avaliação dos procedimentos . . . . .	41
4.6	Considerações do capítulo . . . . .	43
<b>5</b>	<b>Conclusões e Trabalhos Futuros</b>	<b>44</b>
	<b>Referências</b>	<b>46</b>
	<b>Apêndices</b>	<b>49</b>



---

## Introdução

---

A desfiguração de páginas (*deface* ou *defacement* como é conhecido popularmente) é uma forma comum de ataque em sítios Web. Nesses ataques, o conteúdo do sítio é totalmente ou parcialmente trocado pelo atacante, que inclui desde conteúdos embaraçosos, como imagens perturbadoras e formas de assinatura do invasor, até manifestações ofensivas, por exemplo, direcionadas à política (Davanzo *et al.*, 2008).

O ato de desfiguração pode ser realizado através de técnicas que exploram vulnerabilidades presentes em um determinado sítio Web, dando acesso a informações ou dados sensíveis (como: CPF, número de cartão de crédito e senhas) que não deviam aparecer para usuários não autorizados, possibilitando inserção de códigos, que além da desfiguração da página pode causar o roubo de contas e exclusão de dados. A desfiguração é um ato que não pode ser ignorado pois, além dos problemas citados, também causa problemas à reputação das organizações atacadas, já que foi exposto que as mesmas estão vulneráveis a ataques.

O Centro de Estudos, Resposta e Tratamento de Incidentes de Segurança no Brasil (Cert.br) é uma organização que tem o objetivo de registrar e tratar da segurança em computadores que envolvam redes conectadas à Internet brasileira. O Cert.br registrou mais de 55441 (Cert.br, 2017) casos em 2016 classificados como “Web”, que seria, um caso particular de ataque visando especificamente o comprometimento de servidores Web ou desfigurações de páginas na Internet.

Além do Cert.br, há serviços que possibilitam que atacantes registrem suas desfigurações e estes podem ser usados para monitoração dos ataques, como é o caso do **Zone-H**, **Golgeler**<sup>1</sup> e **Hack-Mirror**<sup>2</sup>. Algumas das informações registradas são: *autor*, *alvo*, *espelho da desfiguração*,

---

<sup>1</sup> <http://golgeler.net/>

<sup>2</sup> <http://hack-mirror.com>

*data e hora do registro e sistema operacional.*

O domínio principal do **Zone-H** (zone-h.com) contém registros de sítios de diversos países e o domínio “zone-h.com.br” registra ataques apenas em domínios brasileiros. No “.com” o sítio está alcançando 13 milhões de registros (exatamente 12.957.582 em 02/10/2017) já o domínio “.com.br” tem registrado nessa mesma data, 507 mil desfigurações, sendo aproximadamente 30 mil só em 2017 (Zone-h, 2017). Portanto, o **Zone-H** oferece um conjunto de dados considerável sobre desfiguração e, por isso, foi selecionado como fonte para a construção da base de dados usada para a investigação de desfigurações no Brasil.

Apesar de existir um registro do **Zone-H** específico para sítios Web brasileiros, não há estudos que explorem além das simples estatísticas sobre a desfiguração de páginas no Brasil. Logo, há questões em aberto como: Quem são os grupos ou indivíduos que realizam esse ataques? Quais são suas motivações? Quais os padrões de ataque? Onde eles se organizam ou divulgam suas ações, por exemplo, possuem redes sociais? Essas redes são utilizadas para combinar ataques ou interagir com outros atacantes e grupos? É possível extrair essas redes sociais dos ataques de desfiguração realizadas pelos mesmos?

Tem-se como objetivo geral nesse trabalho o estudo e análise dos grupos *hackers* direcionados à desfiguração de páginas no Brasil, com o intuito de extrair informações para o monitoramento desses grupos e de padrões usados nas desfigurações que possibilitem o desenvolvimento de mecanismos de detecção automática de desfiguração de páginas, de sistemas de alertas antecipados ou a resposta mais rápida a esse tipo de incidente de segurança.

Como objetivos específicos têm-se:

- Automatizar a coleta de publicações em serviços de registro de desfigurações.
- Identificar e investigar os principais autores e padrões de desfigurações de página Web no Brasil.
- Avaliar técnicas de Mineração de Texto que podem ser usadas para detecção de desfigurações de páginas Web em português.
- Gerar bases de inteligência para monitoramento e verificação de desfigurações de páginas Web.

O trabalho está organizado da seguinte forma: O Capítulo 2 apresenta o referencial teórico com conceitos gerais que compõe a pesquisa. O Capítulo 3 expõe as questões de pesquisa levantadas por esse trabalho e a metodologia utilizada para respondê-las. O Capítulo 4 exhibe os resultados obtidos através do método e a discussão sobre cada um deles. Por fim, o Capítulo 5 aponta conclusões sobre este estudo e os trabalhos futuros a serem realizados.

---

## Referencial Teórico

---

Hoje com a evolução das aplicações na Internet, muitas das atividades diárias são realizadas através de sítios Web, como compras, troca de mensagens, busca por informações, etc. Com isso, as empresas têm aderido a aplicações Web e têm melhorado seus investimentos e seu relacionamento com os clientes, pois torna-se um meio mais fácil e acessível de comunicação ou execução de atividades e serviços.

Muitas vezes é necessário que essas aplicações estejam conectadas a Sistemas Gerenciadores de Banco de Dados (SGBD), que armazenam informações e dados do cliente, informações financeiras e afins. Devido a isso, esses sítios têm se tornado alvos de usuários maliciosos que buscam acesso a esses dados para então ganhar dinheiro com eles (por exemplo, conseguindo acesso a números de cartões e senhas), ou apenas expressar sua revolta ou protesto contra determinada organização.

O ato de desfigurar páginas há tempos é uma realidade na Internet, assim como outros ataques que exploram vulnerabilidades presentes em sítios Web, como, *phishing*, *worms*, *negação de serviço* e similares (Bartoli *et al.*, 2009).

Este capítulo aborda sobre desfiguração de páginas, sua definição, o que vêm acontecendo no Brasil em relação as desfigurações de páginas, um visão geral das técnicas utilizadas para tal e meios para se prevenir desse problema, após isso são apresentados sítios Web de espelhamento desses ataques, além de dar noções essenciais sobre mineração de texto e, por fim, os trabalhos relacionados com esta pesquisa.

## 2.1. Desfiguração de páginas

Esta seção aborda sobre definições associadas à desfiguração de páginas, às técnicas usadas por atacantes e os meios de prevenção contra a desfiguração.

### 2.1.1. Definição

Desfiguração de páginas é uma forma comum de ataque em sites Web. Neste tipo de ataque, o conteúdo do site é totalmente ou parcialmente substituído por um atacante que inclui desde conteúdos embaraçosos, como imagens perturbadoras e formas de assinatura do invasor, até manifestações ofensivas, como mensagens agressivas direcionadas ao governo ou a movimentos sociais (Davanzo *et al.*, 2008).

A desfiguração é classificada em duas categorias primárias: *desfigurações substitutivas* (*substitutive defacements*) e *desfigurações aditivas* (*additive defacements*) (Bartoli *et al.*, 2009).

Bartoli *et al.* (2009) descreve que a *desfiguração substitutiva* é caracterizada pela substituição do conteúdo presente no site, ou seja, existe o conteúdo padrão do determinado site e o atacante substitui o mesmo por imagens de revolta, assinatura própria ou do grupo que pertence, a grosso modo, trata-se de uma “pichação virtual”.

Bartoli *et al.* (2009) também descreve que a *desfiguração aditiva* é a adição de uma página ou conteúdo qualquer dentro do site Web atacado, essa página ou conteúdo pode redirecionar os clientes/usuários para locais de controle do atacante e, nesses ambientes controlados, o atacante pode fazer novos ataques e capturar dados sensíveis do usuário sem que este perceba, já que para o usuário a página em que ele está navegando faz parte do site “original”.

Além das categorias primárias, outro tipo importante é a *desfiguração em massa* (do inglês, *mass defacement*), que consiste em desfigurar um conjunto de sites Web explorando uma vulnerabilidade comum, seja essa ação pelo atacante ou por ferramentas automáticas.

### 2.1.2. Desfiguração de páginas no Brasil

A desfiguração de páginas no Brasil é uma realidade constante. Um dos alvos comuns é a desfiguração de sites governamentais. Por exemplo, em 2010, o governo brasileiro sofreu uma grande onda de ataques associados à negação de serviços e desfiguração de páginas (Salatiel, 2011). Sites Web como o da Presidência da República, Portal Brasil, Receita Federal, Petrobras, ministérios do Esporte e da Cultura e o Instituto Brasileiro de Geografia e Estatística (IBGE) foram “vítimas” desses ataques. De acordo com o Serviço Federal de

Processamento de Dados (Serpro), tem-se uma estimativa de que 20 portais do governo federal e 200 sites municipais foram afetados.

Ainda hoje é possível observar esse grande número de desfigurações voltadas a aplicações governamentais, como prefeituras e Universidades, através do Zone-H. O mesmo permite filtrar a desfiguração como “Especial” ou pelo domínio “.gov.br” e logo retorna milhares desses sítios afetados. A Figura 2.1 mostra o resultado de uma pesquisa no Zone-H de desfigurações “Especiais”, ou seja, desfiguração em sítios considerados relevantes.

Total invasões: **29,853** das quais **10,143** única(s) no ip **19,710** ataques em massa

Legenda:

H - Página inicial

M - Ataque em massa (clique para ver todos os ataques neste IP)

R - Re-desfigurado (clique para ver todos os ataques deste site)

★ - Invasão especial (especial quer dizer que é um site importante)

Date	Invasor	H	M	R	★ Domínio	OS	Ver
2017/12/05	G.H.G				★ sisge.pm.pb.gov.br/revista/pub...	Linux	cópia
2017/12/05	BrazilObscure	H		R	★ www.camarasam.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ sisdengue.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ sisadm.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ fiscal.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ www.cpl.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ planodesaneamento.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ www.jfantifumo.pjf.mg.gov.br	Linux	cópia
2017/12/05	Umbrella Gang	H	M	R	★ www.saoraimundododocabezerra.m...	Linux	cópia
2017/12/05	Umbrella Gang	H	M		★ www.mataroma.ma.gov.br	Linux	cópia
2017/12/05	Umbrella Gang	H	M	R	★ www.governadornunesfreire.ma.g...	Linux	cópia
2017/12/05	PYS404		M		★ www.dcl.osasco.sp.gov.br/z.htm	Win 2008	cópia
2017/12/05	PYS404		M		★ bic.osasco.sp.gov.br/z.htm	Win 2008	cópia
2017/12/05	PYS404				★ compras.portovelho.ro.gov.br/z...	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ canalescola.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ sau.pjf.mg.gov.br	Linux	cópia
2017/12/05	DARKRON	H	M	R	★ www.demlurb.pjf.mg.gov.br	Linux	cópia

**Figura 2.1.** Desfigurações especiais no Zone-H (07/12/2017). Fonte: <http://br.zone-h.org/archive/special=1?zh=1>

Não apenas sítios governamentais sofrem com desfiguração, mas também todo tipo de sítio como: imobiliárias, lojas online, *blogs*, bandas, entre outros. A Figura 2.2 apresenta uma lista de desfiguração em sítios como esses citados. Apesar do grande número de registros presentes no Zone-H e em outras aplicações voltadas para esse fim, nem toda desfiguração é notificada ou registrada.

Time	Invasor	H	M	R	★ Domínio	OS	Ver
15:32:36	darkbio				www.bandamarcopolo.com.br/inde...	Linux	cópia
15:02:13	Ashiyane Digital Security Team				www.imagemarte.com.br/arquivos...	Linux	cópia
15:00:54	Erza Jullian				www.construtorajunior.com.br/x...	Linux	cópia
14:53:31	H4R3M				www.inglesparaviagem.net.br/dk...	Linux	cópia
14:37:42	Shade				tecniserv.com.br/sh.html	Win 2012	cópia
14:34:59	./UnIX				www.cabalsuffle.hospedagemlivr...	Linux	cópia
14:17:50	Aris Dot ID				www.amarestufas.pt/x.txt	Linux	cópia
13:41:21	Sukabumi Hackers				radiolegiaourbana.com.br/lisen...	Linux	cópia
13:02:53	Dx_Cyber				evillarodrigues.ntr.br/mimpiin...	Linux	cópia
12:28:30	GokTurk				criareinfo.com.br/loja2/image/...	Unknown	cópia
10:11:39	ProtoWave Reloaded	H			★ www.codasp.agricultura.sp.gov.br	Linux	cópia
10:09:36	ProtoWave Reloaded	H			★ www.codasp.sp.gov.br	Linux	cópia
03:14:34	4nonymous w4s				www.portaldoeconomista.org.br/...	Linux	cópia
01:45:40	DEF4CER				ellyjeans.com.br/elly.php	Win 2008	cópia
00:40:57	DEF4CER			R	contafisco.com.br/mais.php?&pa...	Linux	cópia
00:29:15	TheSun				mercadolivre.idealetroshop.c...	Linux	cópia
00:27:14	TheSun				loja.idealetroshop.com.br/fe...	Linux	cópia
22:31:32	Saz0n			R	★ www.potirendaba.sp.gov.br/admi...	Win 2008	cópia
22:03:24	@Sprek3rsSec	H		R	radiomocadengosa.com.br	Linux	cópia
21:48:12	@Sprek3rsSec	H			radiogracaevida.com.br	Linux	cópia
21:41:44	@Sprek3rsSec	H		R	www.maisalianca.com.br	Linux	cópia
21:31:08	@Sprek3rsSec	H			ieoa.produtoraalphanet.com.br	Linux	cópia
21:23:21	@Sprek3rsSec	H			pleroma.produtoraalphanet.com.br	Linux	cópia
21:13:04	Claden Hackeando	H			www.crb2.org.br	Win 2003	cópia
21:02:55	@Sprek3rsSec	H		R	www.idtp.produtoraalphanet.com.br	Linux	cópia

Figura 2.2. Lista de desfiguração de páginas no Zone-H

No Brasil, a legislação poderia ser usada para punir os atacantes que realizam desfiguração de páginas. Segundo o Prof. Emerson Wendt, Delegado de Polícia Civil do Estado do Rio Grande do Sul, em uma matéria em seu *blog* (Wendt, 2011) concluiu que o atacante dependendo de sua ação pode pegar de seis meses até cinco anos de prisão, uma vez que o atacante indisponibiliza um sítio de utilidade pública ou danifica um patrimônio da União, Estado, Município, empresa concessionária de serviços públicos ou sociedade de economia mista, além de divulgar, sem justa causa, informações sigilosas ou reservadas contidas ou não nos sistemas de informações ou banco de dados da Administração Pública.

### 2.1.3. Técnicas utilizadas para desfiguração

Para a realização da desfiguração de página, é necessário explorar possíveis vulnerabilidades e brechas presentes em sítios Web, para isso são aplicadas técnicas como Injeção de Linguagem de Consulta Estruturada (Injeção SQL, do inglês *Structured Query Language*) e *Cross-site scripting* (XSS), que são duas das técnicas mais usadas. Além disso, uma maneira comum utilizada para encontrar sítios vulneráveis é através de “*google dorks*” que são consultas especializadas utilizando o motor de busca Google, que retornam esses sítios vulneráveis.

Segundo o Projeto Aberto de Segurança em Aplicações (Owasp, do inglês *Open Web Application Security Project*), em 2010, Injeção SQL e XSS eram as duas primeiras técnicas, respectivamente, mais utilizadas para desfiguração de páginas (Owasp, 2010). Em 2016, as duas técnicas ainda permanecem nos primeiros lugares (Hacking, 2016).

A Injeção SQL consiste na inserção de uma consulta SQL através de um dado de entrada do usuário na aplicação. Uma exploração de Injeção SQL bem-sucedida conseguirá ter acesso a dados sensíveis do banco de dados e a partir disso é possível ler, modificar e excluir conteúdos presentes na base, além de executar operações de administrador, como excluir todo banco de dados. Injeção de SQL portanto, é a inserção de comandos SQL em uma entrada de dados presente na aplicação que afetará a execução predefinida de comandos SQL (Owasp, 2010).

A Figura 2.3 apresenta um simples exemplo de injeção SQL e como esse ataque seria executado no SGBD. É possível observar que a primeira consulta é feita normalmente, não fugindo dos padrões e do já esperado pelo SGBD, por ser um usuário e senha comum. Já a segunda consulta, os dados inseridos no campo usuário realizará uma comparação lógica OU no qual a sentença sempre retornará verdadeiro, a seguir, a consulta se encerra com um ponto e vírgula e a adição de um comentário desconsiderando a senha.

Assim, um usuário sem identificação e senha passa a ter acesso à aplicação. A partir desse momento, o invasor pode executar um código malicioso e fazer a inserção de imagens e conteúdo no sítio, assim como o XSS que será apresentado a seguir.

```

Usuário: Nelson
Senha: 123
SELECT * FROM Users WHERE user_id = 'Nelson' AND password = '123'
Usuário: ' OR 1 = 1; /*
Senha: */--
SELECT * FROM Users WHERE user_id = " OR 1 = 1; /*'AND password='*/--'

```

**Figura 2.3.** Exemplo de injeção SQL

O XSS é um ataque que realiza a inserção de códigos onde scripts maliciosos são injetados em sítios confiáveis. Os ataques ocorrem quando um invasor usa um navegador Web para enviar um código JavaScript mal-intencionado para um usuário final (Owasp, 2016).

A vulnerabilidade XSS é executada no navegador da vítima sem o consentimento dela. Com isso, é possível enviar requisições para o servidor usando credenciais de permissão da vítima atacada. Também é possível fazer o sequestro de sessão, que permite acessar o sistema com autenticação da vítima. Portanto, as finalidades do XSS são roubar identificadores de sessão do navegador Web, enganar o usuário fazendo com que o mesmo acesse um conteúdo que pareça o sítio real, mas na verdade é outro com total controle do atacante, proporcionar ataques de negação de serviço e desfiguração de páginas.

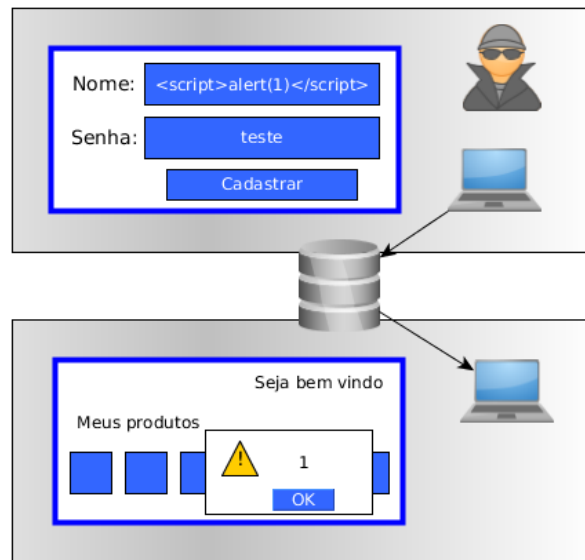
Existem diferentes tipos de ataque XSS, os principais são os tipos chamados

de "Persistente" e "Não Persistente" (Johari; Sharma, 2012). XSS Refletido e Armazenado são tipos de XSS Não Persistentes e Persistentes respectivamente.

O XSS Refletido é aquele em que o servidor reflete o que é enviado para ele, sem filtrar aquele determinado parâmetro, por exemplo, ao preencher o campo usuário, o servidor repete o que foi escrito no código-fonte da página sem tratar o que foi inserido, causando a vulnerabilidade.

Assim, para identificar o XSS Refletido, basta inserir um *script* e verificar se o sistema reproduzirá esse *script*.

O XSS Armazenado é semelhante ao Refletido, porém a aplicação não reflete a entrada diretamente no sítio mas a armazena em seu banco de dados. A Figura 2.4 nos mostra um exemplo de cadastro na qual, ao invés de digitar um nome comum, foi inserido um *script* malicioso no campo nome e como a aplicação não filtrou os caracteres nocivos ao persistir e imprimir o dado inserido, a vulnerabilidade foi explorada.



**Figura 2.4.** Exemplo de XSS armazenado.

Dado o sucesso na exploração da vulnerabilidade, se o atacante ao invés de executar um *script* que gere um alerta, ele pode inserir imagens, realizando assim a desfiguração da página.

Código para inserção de imagem: `<IMG SRC="http://mywebsite.com/defacementpic.jpg">`

Apesar de serem citados apenas a Injeção SQL e o XSS, existem outras técnicas que podem ser exploradas como, Autenticação Violada (Broken Authentication) e Gerenciamento de Sessão (Session Management), CSRF (do inglês Cross-site Request Forgery - Falsificação de Solicitação entre Sites), entre outras<sup>1</sup>.

<sup>1</sup> [https://www.owasp.org/index.php/Top\\_10\\_2013-Top\\_10](https://www.owasp.org/index.php/Top_10_2013-Top_10)



### 2.1.4. Técnicas de prevenção

Existem diversas maneiras de evitar esses ataques apresentados, as principais delas são as boas práticas de programação, já que algumas vulnerabilidades são exploradas através de dados que não são verificados antes de realizar uma ação, como ocorre na injeção de SQL e no XSS.

Outra maneira é a realização de **Teste de Intrusão** (do inglês **Penetration Test** ou **pentest**), que é um método que avaliará a segurança do sítio, tentando de uma forma controlada invadi-lo a fim de detectar vulnerabilidades (Bertoglio; Zorzo, 2015).

O **Teste de Intrusão**, normalmente, é composto pelas seguintes atividades: coleta de informação sobre o sistema alvo; escaneamento e descoberta dos serviços; identificação de sistemas e aplicações; descoberta de vulnerabilidades e exploração de vulnerabilidades (Henry, 2012). Com isso, o teste é robusto o suficiente para fornecer um nível de segurança satisfatório com um alto detalhamento referente às fraquezas do alvo.

Além disso, é recomendado a atualização de *software* e de arcabouços de desenvolvimento Web, uma vez que as tecnologias estão em constante mudanças. Logo, manter o sistema atualizado é essencial para prover um nível de segurança básico contra desfigurações.

## 2.2. Sítios de espelhamento de desfigurações de página

Os sítios de espelhamento são aplicações que registram incidentes de desfigurações de páginas. Esses incidentes são coletados *online* de fontes públicas ou anonimamente. Normalmente os próprios atacantes fazem o registro da página desfigurada. Essas aplicações não têm responsabilidades sobre os ataques que são registrados e nem possuem ligações com os atacantes.

Normalmente esses sítios são estruturados separando os tipos de desfiguração, como “especial” que são os ataques a sítios julgados importantes, e os “normais” que são os demais sítios. Para realizar essa classificação, as desfigurações precisam ser avaliadas pelos administradores, já que pode existir registros falsos de desfigurações. Além dessa separação, existe o “Onhold” que a todo tempo está recebendo registros e já faz a exposição dos mesmos, porém podem existir registros falsos, pois ainda não foram validados.

Dos sítios de espelhamento encontrados, outra informação comum é a divulgação de estatísticas sobre os os atacantes que realizam mais desfigurações, os grupos, a quantidade de registros já feitos. Dentro os sítios de espelhamentos estão: **Golgeler**, **Hack-Mirror** e o **Zone-H**.

O Zone-H foi escolhido como fonte de dados desta pesquisa já que o mesmo é um dos maiores de intrusão Web e é publicado em diversas línguas, uma delas é o português que apresenta desfigurações feitas apenas com o domínio .br, que é de interesse dessa pesquisa. Além de realizar o espelhamento das desfigurações, o Zone-H é um portal de segurança da Internet, contendo informações de segurança da informação (Preatoni, 2017). O Zone-H também segue a estrutura organizacional já descrito, como os demais sítios.

## 2.3. Mineração de Texto

A Mineração de Texto é definida como o processo de encontrar padrões eficientes, modelos, direções, tendências e regras em um conjunto de dados textuais (Nahm; Mooney, 2002).

Assim, combinando técnicas de Mineração de Dados, Aprendizagem de Máquina, Processamento de Linguagem Natural, Recuperação de Informações e Gerenciamento de Conhecimento, é possível a extração dessas informações relevantes, utilizando técnicas que são capazes de identificar e explorar padrões de interesse presentes no conjunto de dados textuais (Feldman; Sanger, 2007).

Esse conjunto de dados textuais podem ser chamados de não-estruturados, semi-estruturados e estruturados. Os textos que não obedecem a um padrão de formatação são os não-estruturados, os que seguem algum padrão, como textos científicos e livros são semi-estruturados, por fim, os textos representados em uma linguagem de marcação são os estruturados (Conrado, 2009).

Nas seções seguintes, são apresentados os conceitos sobre algumas técnicas de mineração de texto para explanação desse processo, onde essas técnicas, exceto *lematização* e *radicalização*, foram utilizadas pelo trabalho.

### 2.3.1. Tokenização

A *Tokenização* tem como objetivo extrair unidades mínimas no texto. Cada unidade é chamada de *token* e corresponde a uma palavra do texto, que pode estar relacionada também a símbolos e caracteres de pontuação (Manning *et al.*, 2008).

Por exemplo, a frase “Amanhã chove em Campo Mourão!”, poderá ser dividida em seis *tokens*: [Amanhã] [chove] [em] [Campo] [Mourão] [!].

Quando gerado os *tokens*, o “espaço” sempre é descartado. No caso de páginas Web, se o texto não estiver pré-processado, é necessário desconsiderar as TAGs HTML além do espaço.

### 2.3.2. Lematização e radicalização

A *radicalização* (“Stemmização” ou *Stemming* como é conhecido) tem como objetivo reduzir as palavras às suas formas inflexionáveis e às vezes reduzir às suas derivações (Manning *et al.*, 2008). A *radicalização* reduz cada palavra do texto ao seu provável radical, ou seja, palavra raiz (*stem*), em que cada palavra é analisada isoladamente (Conrado, 2009).

A palavra raiz não é necessariamente idêntica à raiz morfológica da palavra, mas é suficiente para relacionar e mapear palavras a ela. Exemplo: *Análise dos grupos hackers no Brasil*. Considerando a remoção de “stopwords” (palavras consideradas irrelevantes para um conjunto de resultado) o resultado da radicalização é a seguinte: **anális grup hack brasil**.

Ao realizar *radicalização* é necessário ter cuidado com os efeitos *overstemming* e *understemming*. *Overstemming* ocorre quando o resultado extraído não é um sufixo, mas sim parte do radical. Por exemplo, a palavra «gramática», após o processamento é reduzida para «grama», o que não representa o seu radical, que é «gramat». *Understemming* ocorre quando o sufixo não é removido totalmente. Por exemplo, a palavra «referência», após o processamento é reduzida para «referênc», ao invés de «refer», que é o radical correto (Morais; Ambrósio, 2007).

A técnica de *lematização*, ou *Redução à Forma Canônica* como conhecida, transforma verbos para sua forma no infinitivo, e substantivos e adjetivos para o masculino singular (Conrado, 2009). A Figura 2.5 apresenta a redução de palavras para o seu lema:

Lema	SingularFem.	PluralFem.	PluralMasc.
brasileiro	brasileira	brasileiras	brasileiros
pesquisa	pesquisa	pesquisas	pesquisas
perfil	perfil	perfis	perfis
estudante	estudante	estudantes	estudantes

**Figura 2.5.** Exemplo de lematização. Fonte: Moraes e Ambrósio (2007)

Com essas técnicas é possível padronizar palavras e classificá-las como uma informação relevante ou não.

### 2.3.3. Associação de palavras

O objetivo da associação de palavras é a ligação automática de documentos texto a uma determinada classe, pertencente a um conjunto predefinido de classes. Conceitos importantes dentro da associação são: *análise paradigmática* e *análise sintagmática* e *colocações*.

A *análise paradigmática* e *sintagmática* seriam, as relações de seleção e as relações de combinação entre os elementos linguísticos, respectivamente.

A *análise paradigmática* busca uma série de elementos linguísticos que possam expressar o mesmo sentido, ou seja, na frase: “Foi teu avô”. No lugar de “teu”, é possível figurar, se o sentido do enunciado fosse outro, os termos seu, meu, nosso, o, etc.

Por outro lado, na *sintagmática* não se combina qualquer elemento aleatoriamente. Ela é vista como uma unidade formada por uma ou várias palavras que, juntas, desempenham uma função na frase.

As *colocações* têm o objetivo de agrupar palavras onde o significado é a soma dos significados das partes, além de algum componente semântico adicional. Como exemplo, pode-se citar: cabelo branco, pele branca e vinho branco, tal que o branco do cabelo é cinza, o branco da pele é rosado e o branco do vinho é amarelo (Santos, 2002).

### 2.3.4. Reconhecimento de entidades nomeadas

O Reconhecimento de Entidades Nomeadas (REN) consiste na tarefa de identificar as entidades nomeadas (EN), na sua maioria nomes próprios, a partir de textos de forma livre e classificá-las dentro de um conjunto de tipos de categorias predefinidas, tais como *pessoa*, *organização* e *local*, as quais remetem a um referente específico (Mota *et al.*, 2007).

Segundo Sureka *et al.* (2009), o REN e a posterior classificação de tais entidades é uma técnica amplamente utilizada no PLN e consiste na identificação de nomes de entidades-chave presentes na forma livre de dados textuais. A entrada para o sistema de extração de entidade nomeada é o texto de forma livre, e a saída é um conjunto das chamadas anotações, ou seja, grupo de caracteres extraídos de trechos do texto de entrada. A saída do sistema de extração de entidades nomeadas é, basicamente, uma representação estruturada a partir da entrada de um texto não estruturado.

As três principais abordagens para extração de entidades nomeadas são: *sistemas baseados em regras*, *sistemas baseados em aprendizado de máquina* e *abordagens híbridas*. *Sistemas baseados em regras* ou *sistemas baseados no conhecimento* consistem em definir heurísticas na forma de expressões regulares ou de padrões linguísticos. *Sistemas baseados em aprendizado de máquina* utilizam algoritmos e técnicas que permitam ao computador aprender a reconhecer entidades com base em textos. E as abordagens *híbridas* combinam elementos das duas abordagens anteriores (Amaral; Vieira, 2013).

## 2.4. Trabalhos relacionados

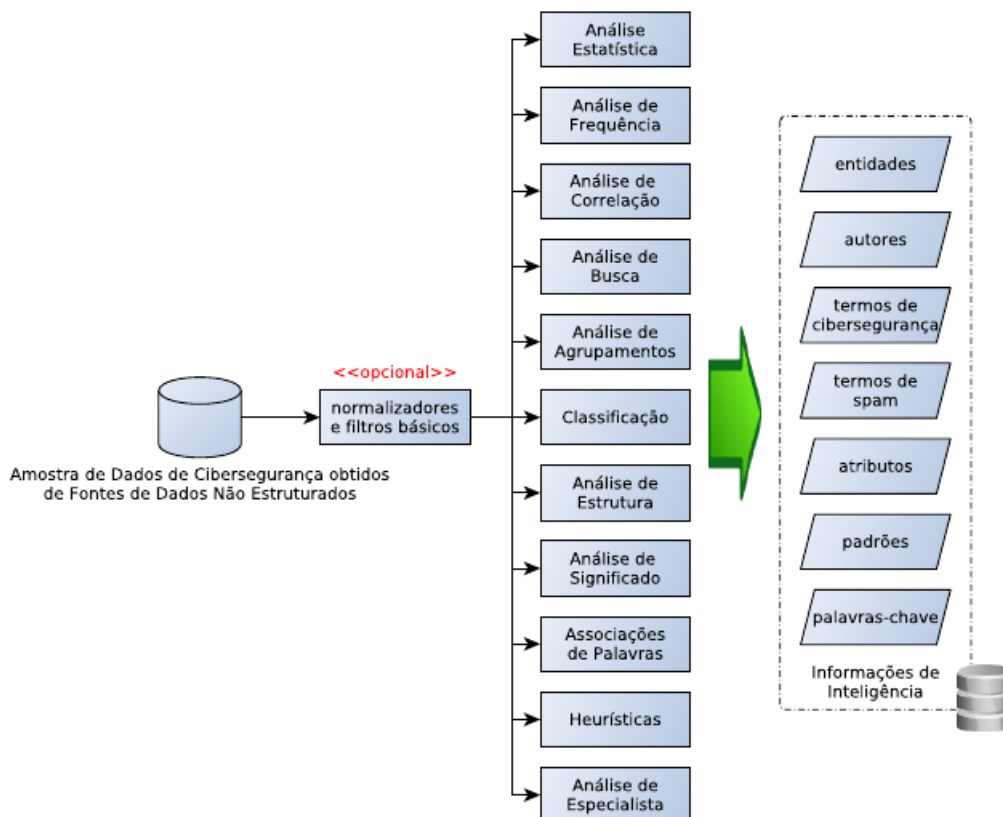
Essa seção aborda os trabalhos relacionados à pesquisa, formas de como analisar dados não estruturados, a importância e técnicas para detecção automática de desfiguração de

páginas e a investigação dos indivíduos que realizam o ato de desfigurar através de redes sociais.

Em Campiolo (2016) é proposto um arcabouço que apresenta várias técnicas para a análise de dados não estruturados associados à cibersegurança. O resultado da análise são conjuntos de elementos usados para a geração de bases de inteligência que servirão para a identificação e extração de alertas de cibersegurança da fonte analisada.

A Figura 2.6 apresenta os diferentes tipos de análise que podem ser realizadas e o que pode ser extraído que é relevante para a geração de alertas. Essas informações são usadas para criação da base de inteligência e algoritmos de extração de alertas. As saídas dos processos de análise propostas no arcabouço são categorizadas em sete grupos de interesse: *entidades*, *autores*, *termos de cibersegurança*, *termos de spam*, *atributos*, *padrões* e *palavras-chave*.

O presente estudo utilizou algumas técnicas apresentadas por Campiolo (2016) como: *Análise Estatística*, *Análise de Frequência*, *Análise de Correlação*, *Associações de Palavras* e *Heurísticas*; e produziu informações de inteligências como: *entidades*, *autores*, *termos de spam*, *padrões* e *palavras-chave*.



**Figura 2.6.** Análise de dados de cibersegurança obtidos de fontes de dados não estruturados.  
Fonte: (Campiolo, 2016)

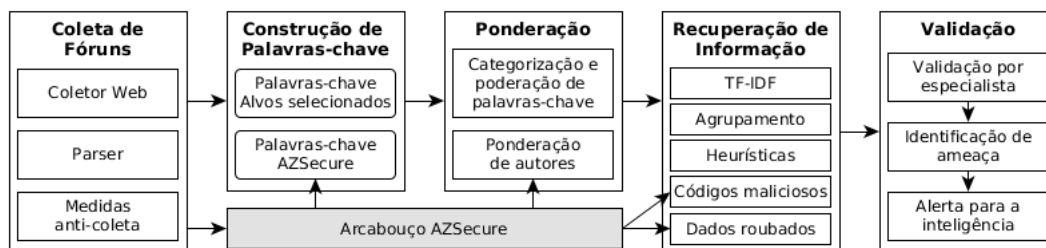
Davanzo *et al.* (2008) realizou um estudo sobre técnicas para detecção automática de sítios

que sofreram desfiguração de páginas e constatou que 40% dos sítios manteve a desfiguração por até 1 semana e 37% foram corrigidos apenas depois de duas semanas. Em seu trabalho foi proposto o uso de técnicas de detecção de anomalias para detectar automaticamente essas páginas e alertar aos proprietários o ocorrido. A característica crucial da proposta é fazer com que o mecanismo de detecção não tenha dependência do sítio a ser monitorado e seus envolvidos, e assim, criar um perfil do sítio monitorado, para que o sistema emita um alerta quando aparecer algo incomum. Como resultados, o estudo evidenciou a quantidade de registros feitos por desfiguração de páginas, o tempo em que os sítios normalmente demoram para identificar a desfiguração, como apresentado anteriormente. Entretanto, não abordou o lado do indivíduo que realiza a desfiguração, o que é compreensivo já que seu foco foi a criação de um mecanismo de detecção automática.

Muitas vezes o foco de investigação sobre crimes cibernéticos se direciona para vulnerabilidades do sistema e para mecanismos de detecção antecipada, porém muito pouco trabalho tem sido feito para ir além das questões tecnológicas e investigar o indivíduo por trás desses crimes ou ataques.

Benjamin *et al.* (2015) apresentam um foco semelhante a nossa pesquisa em seu estudo, que é buscar um maior conhecimento sobre os grupos *hackers*. Seu objetivo foi desenvolver uma metodologia automatizada para identificar evidências tangíveis e verificar ameaças potenciais em fóruns *hackers*, canais IRC (do inglês, *Internet Relay Chat*, utilizado basicamente como bate-papo (*chat*) e para troca de arquivos, que permitem a conversa em grupo ou privada) e *carding shops* (sítios onde se comercializa número de cartões de crédito). Com isso, sua abordagem permitiu refinar ameaças em potenciais dos conteúdos *hackers* recolhidos.

Para isso, foi necessário estudar cada plataforma, pois cada uma tem suas próprias complexidades e requer diferentes estratégias para a investigação. A Figura 2.7 mostra a abordagem usada nos fóruns *hackers*.



**Figura 2.7.** Proposta de arcabouço para fóruns hackers. Adaptado de (Benjamin *et al.*, 2015)

O primeiro processo é a coleta de informações e mensagens do fórum, que foi feito por meio do uso de *Web Crawler* e *Parsers*. O segundo passo foi a criação de palavras-chave, que ajudam a identificar potenciais ameaças. Para isso, foi utilizado o AZSecure (arcabouço criado anteriormente pelos autores (Li; Chen, 2014)) para realizar a seleção de palavras. O

terceiro passo foi categorizar e dar peso as palavras-chave e para os autores. O quarto passo se concentrou especificamente na identificação de artigos contendo ameaças potenciais, e em classificar postagens com base na relevância e urgência.

A coleta e análise do dados do IRC começou com a identificação dos canais *hackers* no IRC através de uma coleção de palavras-chave. Dado a identificação, aplicativos automatizados (*bots*) foram implantados nos chats para coletarem dados em tempo real. As palavras-chave também foram utilizadas na identificação dos temas mais populares entre os participantes. Esse processo foi útil, pois forneceu um rápido resumo das conversas que estão presentes nas comunidades *hackers* do IRC.

Como todas as mensagens eram transmitidas publicamente, foi possível usar endereçamento direto para calcular a rede social entre os participantes para cada comunidade. Com isso foi possível identificar os **atores-chave** que puderam fornecer mais evidências de potenciais ameaças.

Para a coleta e análise dos metadados dos *carding shops* foi proposto outro arcabouço que aproveita os dados coletados dos fóruns unindo-os com os metadados dos *carding shops*. Após isso foram geradas novas palavras-chave e foi feita a extração de características delas, finalizando a coleta com a classificação dos textos. Na fase de análise dos metadados, foram obtidas a classificação do vendedor e o seu perfil.

Alguns *carding shops* mostram informações detalhadas de localização, tais como código postal, porém muitos deles listam apenas país ou estado. Com isso, foi possível observar que fóruns de *hackers*, canais IRC e lojas de cartão de créditos contém uma variedade de conteúdos relevantes para a descoberta de ameaças cibernéticas atuais e emergentes, além de informar os *hackers* atuantes.

## 2.5. Considerações do capítulo

Este capítulo apresentou o conceito de desfiguração de página, a realidade desse ato no Brasil, algumas técnicas utilizadas para desfiguração de páginas e técnicas de prevenção. Além disso, foram apresentadas as técnicas de mineração de texto aplicadas no estudo, como tokenização, lematização e radicalização, associação de palavras e reconhecimento de entidades nomeadas. Foram expostos os trabalhos relacionados a este, os quais apresentaram diferentes formas (arcabouços) de coleta de informações inteligentes e relevantes que contribuem para sistemas de alerta antecipado, além de variados meios de coleta da base de dados, como Twitter, fóruns *hackers*, canais IRC e *carding shops*. Essa pesquisa assemelha-se aos os trabalhos apresentados, porém com uma nova fonte de coleta, o **Zone-H**. Com isso o tratamento do conteúdo se difere, já que este não é um texto plano, mas sim HTML.

---

## Método de Pesquisa

---

Esse capítulo apresenta o método de pesquisa aplicado no estudo, iniciando com a apresentação das questões de pesquisa que nortearam o desenvolvimento deste trabalho e uma breve discussão de como elas foram investigadas. Além disso, apresenta-se um fluxograma que abrange o passo a passo tomado para o caracterização e discussão dos resultados, além dos materiais e códigos desenvolvidos durante a pesquisa.

### 3.1. Questões de pesquisa

Neste projeto objetivou-se investigar grupos/indivíduos que atuam em desfiguração de páginas no Brasil e, por meio do registro de seus ataques, extrair informações que possam colaborar no monitoramento dos mesmos. Para tal, foram estabelecidas as seguintes questões de pesquisa:

- Q1.** Quais os grupos/indivíduos *hackers* mais ativos em desfiguração de páginas no Brasil?
- Q2.** É possível relacionar os grupos/indivíduos *hackers* com perfis ou páginas em redes sociais?
- Q3.** É possível estabelecer um padrão de desfiguração (características na página) segundo um grupo/indivíduo?
- Q4.** É possível identificar um padrão nas páginas para caracterizar uma desfiguração?

Ao identificar os grupos/indivíduos mais ativos (Q1), considerando que os mais ativos são aqueles que possuem o maior número de desfigurações registradas, direciona-se a investigação e caracterização para esses, pois necessitam de um maior cuidado, gerando características particulares que colaboraram no monitoramento dos mesmos.

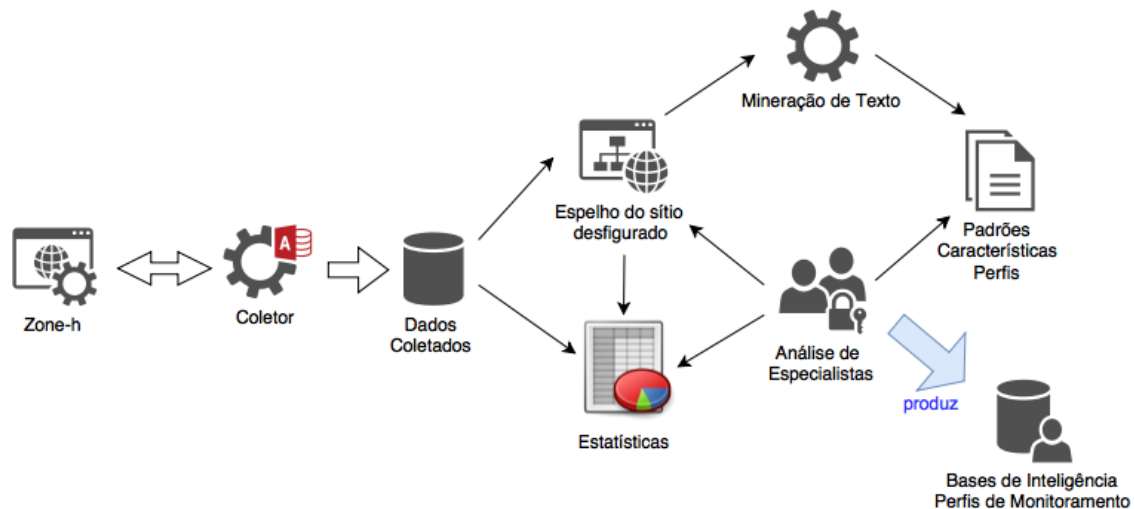


Com o objetivo de gerar informações de inteligência, as redes sociais (Q2) são essenciais, já que são uma grande fonte de investigação e monitoramento. Estudos como o de Campiolo (2016), evidenciam isso e mostram que é possível obter notificações antecipadas de ataques através delas.

Identificar padrões na desfiguração, tanto individuais (Q3) quanto gerais (Q4), é interessante e relevante, já que esses padrões podem ser usados para auxiliar a identificação dos ataques nas páginas Web brasileiras. A Seção 3.2 apresenta os materiais e métodos para a investigação das questões de pesquisa.

## 3.2. Materiais e métodos

Para responder as questões de pesquisa propõe-se o método de pesquisa apresentado na Figura 3.1. Além da explanação do método, também são apresentados os materiais utilizados no estudo.



**Figura 3.1.** Fluxograma do método de pesquisa.

O método consiste na coleta de dados através do sítio Web Zone-H, criando assim uma base de dados. Dentre esses dados estão os espelhos de sítios desfigurados que, contém padrões, características e perfis a serem monitorados, que podem ser extraídos com técnicas de mineração de texto. Dado a extração, é incluído uma **análise especialista** para validar resultados obtidos como, redes sociais, **associação de palavras** dentro de um contexto ou **similaridade** entre desfigurações. Além disso, é possível através da **análise estatística** extrair os grupos que mais atuam em desfiguração e caracterizar as tendências de domínios-alvos por cada um, tudo isso é armazenado nas bases de inteligência que, futuramente, poderão ser utilizadas por mecanismos de detecção de desfigurações. As próximas subseções detalham cada um dos itens citados no fluxograma.

### 3.2.1. Zone-H

O Zone-H, sítio Web escolhido para criação da base de dados por seu vasto número de registro de desfigurações, possui ligações (*links*) para recursos que possibilitam realizar a coleta desses registros, como, “Arquivo” e “Arquivo Especial” em que os registros já foram analisados pelos administradores e confirmados como uma desfiguração e o “Onhold” na qual as desfigurações são registradas e já são disponibilizadas, proporcionando rapidez na divulgação dos ataques, apesar da possibilidade do ataque não ser de fato uma desfiguração. A fonte do nosso coletor vem do “Onhold”.

Como já dito na Seção 2.2, o Zone-H é publicado em diversas línguas. Foi utilizado o “*zone-h.com.br*”, pois o mesmo nos retorna apenas desfigurações em domínios .br, sem a necessidade de filtros para recuperar esses domínios.

### 3.2.2. Coletor

O Coletor consiste em um *script* em Python (versão 2.7.12) que utiliza a API Selenium<sup>1</sup>, na versão 2.53.6, para realização de Web Crawling no Zone-H. O Selenium simula a execução de um navegador Firefox, na versão 39.0.3, assim, o coletor visita o site Zone-H periodicamente (30 em 30 minutos) procurando por novos registros de desfiguração, criando assim uma base dados para extração de informações. O Código 3.1 (o código-fonte do coletor escrito em Python se encontra no Apêndice A) descreve em pseudo-algoritmo o funcionamento do coletor.

```
def coletor:
    abre o firefox e acessa o zone-h

    for cada registro na pagina:
        if registro nao existe:
            salva no banco de dados
        else:
            interrompe execucao do coletor
        if registro possui popup:
            ignora o registro e vai para o proximo
        if captcha solicitado:
            gera notificacao de captcha solicitado
```

**Código 3.1.** Coletor automático

Um problema encontrado no coletor foi a questão da inserção de *popups* pelo invasor, pois o coletor busca por campos determinados e não consegue achá-los, com isso aquele registro é

<sup>1</sup> [http://www.seleniumhq.org/docs/03\\_webdriver.jsp](http://www.seleniumhq.org/docs/03_webdriver.jsp)

perdido e o coletor gera uma exceção e continua sua execução para o próximo registro. Outro problema é referente ao *captcha*, que quando solicitado, a execução do coletor é interrompida.

### 3.2.3. Dados coletados

Os dados coletados geraram a base de dados que consiste em registros do dia 12/01/2017 até o dia 12/11/2017 totalizando em uma base com 7184 registros. Esses registros estão armazenados no sistema gerenciador de banco de dados (SGBD) MySQL. Alguns dados possuem a data anterior a especificada pois estes não vieram do “Onhold” mas sim do “Arquivo Especial”, porém esses dados são minoria (55 registros, sendo 43 de 2015 e 12 de 2016). Vale considerar a perda de alguns dados devido a *popups* no registro, como descrito no Coletor. Essa perda não foi mensurada.

Os dados coletados foram:

- **Data do registro:** indica o dia, mês, ano e horário em que foi registrada a desfiguração no Zone-H. Essa data não representa o momento em que ocorreu a desfiguração de fato.
- **Invasor:** identifica o grupo ou invasor que realizou a desfiguração.
- **Domínio:** exibe a URL do sítio Web desfigurado.
- **Endereço IP:** representa o endereço do sítio Web desfigurado.
- **Sistema:** Sistema Operacional presente no servidor do sítio desfigurado.
- **Web Service:** serviço Web presente no sítio desfigurado.
- **HTML da página:** essa informação é recuperada através da visita no domínio (disponibilizado no registro) do sítio Web desfigurado.
- **HTML registrado no Zone-H:** todo registro possui um espelho do sítio desfigurado e essa informação também é recuperada.

Os dados coletados são no formato HTML, e eles passam por um pré-processamento para a remoção das marcações HTML, códigos Javascript e CSS através da biblioteca BeautifulSoup<sup>2</sup>.

### 3.2.4. Estatísticas

A análise estatística consiste na extração de informações que descreveram a base em números que viabilizaram observar, analisar e associar os dados presentes na base de dados (o Zone-H possui um espaço para estatísticas, porém, nada é apresentado em relação aos invasores e grupos, apenas os números diários/mensais/anuais de desfigurações registradas).

---

<sup>2</sup> <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

Nessa análise foi feito o levantamento dos grupos/indivíduos mais ativos no Brasil, com base nos dados coletados. A partir desses invasores foi realizada a associação de cada um com os domínios de cada desfiguração realizada pelos mesmos, conseguindo identificar tendências de ataques para períodos analisados. Por meio das datas, foi possível observar os dias da semana e as horas em que mais ocorrem desfigurações.

Devido à organização e disponibilidade de dados realizou-se a extração dessas estatísticas através de consultas SQL. A descrição estatística da base de dados se encontra no Apêndice B.

Para obter os grupos e indivíduos mais ativos foi realizada uma consulta ao banco de dados que é apresentado pelo Código 3.2.

```
SELECT invasor , COUNT(*)
FROM dataDefacers
GROUP BY invasor
ORDER BY COUNT(*) DESC LIMIT 20
```

**Código 3.2.** Consulta SQL hackers mais ativos

Com isso, foi possível resgatar os principais grupos/indivíduos presentes em nossos registros. A amostra fixada para a análise foi de 20 grupos em um total de 866. Essa amostra corresponde a 2,31% dos grupos e suas respectivas desfigurações à 38,1% do total da base. A ordenação com os grupos/indivíduos que realizaram no mínimo 10 desfigurações está disponível no Apêndice C. Este resultado foi armazenado em um arquivo CSV (do inglês *Comma-separated values*) com os valores: invasor e quantidade de desfigurações.

Após resgatar grupos/indivíduos presentes na base de dados, foram caracterizados os domínios em que os mesmos costumam realizar seu ataques. O Código 3.3 apresenta um pseudo-algoritmo e demonstra como esse processo foi realizado.

```
def caracterizaDominio:
  for cada invasor:
    SELECT dominio
    FROM dataDefacers
    WHERE invasor in VinteMaisAtivos

    for cada tipo de dominio:
      if o tipo == ".com.br" or ".gov" or "demais dominios":
        incrementa variavel do dominio respectivo

    escreve no arquivo csv o invasor e a frequencia de cada dominio
```

**Código 3.3.** Caracterizando domínios por invasor

A partir do resultado obtido com o Código 3.3, foi observado e registrado a preferência de cada atacante, onde os domínios caracterizados foram: “.com.br”, “.gov”, “.blog”, “.edu.br”.

“br (Universidades)” e “outros domínios”. O resultado desta consulta também foi armazenado em um arquivo CSV com o conteúdo: invasor, domínios e quantidade de desfigurações nos respectivos domínios.

Dada a identificação dos grupos/indivíduos mais ativos e seus domínios de preferência, restou apenas a caracterização dos horários e dias da semanas em que mais acontecem desfigurações. Vale ressaltar que essa informação é retirada sobre a data e hora capturada do Zone-H, que não é exatamente o momento em que ocorreu a desfiguração, mas sim a data e hora em que o ataque foi notificado e registrado no Zone-H. Como o atacante busca divulgar seu ataque antes que o mesmo seja removido, têm-se essa hora como próxima a realização do ataque. Os Códigos 3.4 e 3.5 apresentam pseudo-algoritmos que mostram como essa caracterização descrita foi obtida.

```
def diasDaSemana:
    #Seleciona a data de todos os registros da base de dados
    SELECT dataRegistro FROM dataDefacers

    for cada data:
        #DAYOFWEEK retorna um numero inteiro que corresponde ao dia da semana
        dia = SELECT DAYOFWEEK(data)
        if dia == "Domingo, ..., Sabado":
            incrementa a variavel corresponde ao dia da semana

    escreve no arquivo csv o dia da semana e a frequencia de dia
```

**Código 3.4.** Quantidade de desfiguração por dia da semana

```
def horarios:
    //hora = 01:00, 02:00, ..., 23:00
    for cada hora:
        SELECT COUNT(*) FROM dataDefacers WHERE HOUR(dataRegistro) = hora
        salva o COUNT retornado pela consulta e associe com a hora

    escreve no arquivo csv a hora e a frequencia de cada hora
```

**Código 3.5.** Quantidade desfiguração por hora

### 3.2.5. Espelho do sítio desfigurado

O espelho do sítio desfigurado consiste em uma cópia do sítio desfigurado, isto é, o código fonte do HTML. Esse código é usado para a mineração de texto visando identificar padrões nos ataques, extrair novas estatísticas, identificar os termos e entidades mais comuns e fazer a extração de redes sociais do atacante.

O espelho do sítio foi resgatado através do Web Crawler que utiliza da linguagem de consulta XPath, que possibilita a localização da URL do espelho no Zone-H e assim é realizado

o *download* do HTML. Todo registro no Zone-H possui um espelho da desfiguração. Além do espelho da desfiguração, foi coletado o espelho do sítio original.

O processamento do espelho do sítio desfigurado é usado para responder as questões de pesquisa, exceto a primeira. Portanto, identificar redes sociais, caracterizar desfigurações de grupos ou indivíduos e identificar termos que possam caracterizar uma desfiguração, só é possível através do espelho do sítio desfigurado.

Apesar de ser coletado, o espelho original do sítio não está sendo utilizado, pois o mesmo não é necessário para alcance do objetivo do trabalho. O intuito com ele é fazer uma verificação através de técnicas de similaridade, para saber se o sítio original ainda possui a desfiguração em seu conteúdo, porém, isso foi proposto como um trabalho futuro.

### 3.2.6. Mineração de texto

A mineração de texto consiste na aplicação de técnicas como **tokenização**, **associação de palavras**, **análise de frequência**, identificação de entidades e termos e extração de redes sociais para construção de uma base de informações de inteligência. Todo esse processo é feito através da biblioteca NLTK (a extração de redes sociais possui um algoritmo extra criado pelo autor da pesquisa, como já citado na seção 3.1).

A técnica de **tokenização** foi realizada através da função `"word_tokenize"`, com isso é possível extrair unidades mínimas do texto. Uma vez que os dados estão separados em *tokens*, é possível fazer a **análise de frequência**, **identificação de entidades** e **associação de palavras**, ou seja, a **tokenização** é o primeiro passo e que é utilizado nas demais técnicas.

A **associação de palavras** foi realizada por meio da função `"ngrams"` que possibilita buscar uma palavra-alvo e obter outras palavras que a acompanham, assim é possível fazer identificação de redes sociais e caracterizar palavras que frequentemente aparecem conjuntas, entendendo seu contexto, se necessário.

A análise de frequência foi realizada através do `"FreqDist"` e `"Counter"` que possibilitam obtenção dos termos mais frequentes presentes no HTML. Nesta etapa é importante a remoção de *stopwords* (termo explanado na Sub-subseção 2.3.2). As entidades são recuperadas através da utilização de classificação gramatical por meio do corpus *Mac\_Morpho*.

Corpus *Mac\_Morpho* é um conjunto de documentos, formado por artigos publicados no jornal Folha de São Paulo, em 1994, contendo mais de 1 milhão de palavras, anotadas pelo etiquetador de palavras (BICK 2000)<sup>3</sup>. Ele é utilizado no NLTK como treino para

---

<sup>3</sup> <https://sites.google.com/site/linguacorporus/acdc/mac-morpho>

classificação gramatical dos tokens, outro corpus disponível, porém não utilizado, é o "Floresta Sintática"<sup>4</sup>.

A extração de redes sociais foi realizada através da técnica de **associação de palavras**, onde é possível definir a palavra-alvo como: "Facebook", por exemplo. Desta maneira, caso seja divulgado o Facebook do atacante no HTML é possível resgatar o seu perfil.

Em específico, o primeiro passo para a extração da redes é a normalização do HTML e a criação de tokens. Depois disso é definido palavras-alvo, que são as redes sociais (Facebook, Twitter, Skype, IRC) e meios de comunicação (email), por fim, é utilizado `FreqDist` e `ngrams` para devolver palavras associadas às palavras-alvo. Com isso é possível observar identificadores e *links* associados as redes sociais. O Código 3.6 representa essa especificação.

```
def redeSociais(tokens):
    #tokens sao as palavras do html da desfiguracao
    #Palavras que serao procuradas
    target_words = skype, facebook, twitter, irc, gmail/hotmail/yahoo
    fd = FreqDist(ng
        for ng in ngrams(tokens, 5)
            if target_words in ng)
    for hit in fd:
        print(' '.join(hit))
```

**Código 3.6.** Extração de redes sociais NLTK

É passado para o `ngrams` a lista de *tokens* do HTML e um valor inteiro que é a quantidade de palavras que estará associado com a palavra-alvo. É necessária a análise do especialista, já que pode ser resgatado redes sociais presentes na página que não está associado a um grupo ou indivíduo mas sim ao proprietário do sítio por exemplo, portanto, cabe ao especialista remover essas informações irrelevantes.

O algoritmo desenvolvido neste trabalho que colabora na extração de redes sociais é apresentado pelo Código 3.7 em pseudo-código.

```
def redeSociais(html):
    listaRedesSociais
    #html esta normalizado e em caixa baixa
    for cada palavra no html:
        if palavra in {"facebook", "fb", "skype", "gmail",
            "yahoo", "hotmail", "irc", "twitter"}
            if palavra not in listaRedesSociais
                listaRedesSociais = palavra + palavra[1] + palavra [2]
    return listaRedesSociais
```

**Código 3.7.** Extração de redes sociais

<sup>4</sup> <http://www.linguateca.pt/Floresta/>

Os padrões de desfigurações, tanto para um grupo/indivíduo como para padrões gerais, também foram utilizadas técnicas de análise de frequência, associação de palavras, além de testes de similaridade.

Para padrões individuais foi resgatado o espelho da desfiguração dos principais grupos/invasores e assim foram realizadas a análise e a extração das principais palavras utilizadas pelo grupo, o Código 3.8 apresenta o pseudocódigo da solução desenvolvida.

```
#Principais grupos/individuos
invasores = cursor.execute("SELECT COUNT(*), invasor
                            FROM dataDefacers
                            GROUP BY invasor
                            ORDER BY COUNT(*) DESC LIMIT 20")
para cada invasor em invasores:
    htmls = todas as desfiguracoes desse invasor
    remocao de 'stopwords' dos htmls
    criacao de uma lista com todas as palavras
    funcao Counter recebe as palavras e retorna as palavras mais frequentes
```

**Código 3.8.** Extração das palavras mais frequentes por grupo/indivíduo

Com as palavras mais frequentes, foi identificado alguns padrões, como vocabulário usado nas desfigurações. Através da análise especialista, que é explicada na subseção 3.2.7, selecionou-se palavras para entender o contexto e caso necessário usar **associação de palavras com bigramas e trigramas** como já explicado anteriormente.

Para análise de similaridade foi utilizada a biblioteca em Python "**Scikit-learn**" e em específico a função **tf-idf**(do inglês *term frequency-inverse document frequency*), assim obteve-se padrões de desfigurações utilizado pelos grupos, como por exemplo, a assinatura padrão do grupo. Foi feita uma análise temporal para verificar se esse padrão sofria mudanças com o tempo.

Para padrões gerais, foi aplicado a análise de frequência de palavras novamente, extraindo as palavras mais frequentes em todas as desfigurações. Uma vez que são obtidos esses termos, é feito uma análise de frequência confirmando o termo como um padrão segundo o percentual de ocorrências do mesmo nas desfigurações. Além dos termos, o mesmo processo foi repetido com *nicknames* dos invasores e com os nomes de grupo coletado.

### 3.2.7. Análise especialista

A análise especialista consiste na análise manual dos resultados por um especialista em segurança com o objetivo de identificar e filtrar padrões, características e perfis (Seção 3.2.8) e criar as bases de inteligência (Seção 3.2.9).



A análise foi feita por um único especialista, cujos critérios utilizados foram aplicados com a ajuda da biblioteca NLTK para validação dos termos mais frequentes. Uma vez que têm-se os termos, eles são usados como alvo para verificar se o que acompanha esses termos estão dentro do contexto de uma desfiguração, tornando-o um padrão ou característica de uma desfiguração. Os termos mais explícitos, que não necessitam de validação, foram excluídos diretamente, por exemplo, *stopwords*.

No caso da validação de redes sociais, avaliou-se se o perfil estava associado ao grupo/invasor, se não possuísse ligação com o grupo/invasor, o perfil era descartado. Além disso, se os links divulgados fossem inválidos, os mesmos também eram descartados.

### 3.2.8. Padrões, características e perfis

Padrões, características e perfis são resultados da aplicação de técnicas de mineração de texto já descritas e análise de especialista. Destacam-se os padrões encontrados que caracterizam ataques de grupos, termos que podem caracterizar uma desfiguração e perfis em redes sociais sujeitos a monitoramento.

Os padrões e características, por exemplo, podem ser palavras que na maioria das vezes estão presentes nas desfigurações, como: «hacked» e «by». Esses termos podem ser validados através da associação de palavras. Essa técnica evidencia que as palavras que acompanham aquele termo se referem a desfiguração, garantindo que o mesmo pode se tornar um padrão ou característica. Além disso, determinados grupos sempre que realizam um ataque passam uma mesma mensagem, ou injetam imagens que caracterizam aquele grupo.

Além desses padrões e características, perfis de redes sociais podem ser extraídos dos HTMLs, pois alguns grupos costumam deixar registrado suas redes em seus ataques.

### 3.2.9. Bases de inteligência / Perfis de monitoramento

As bases de inteligência consistem em bases que armazenam informações relevantes que podem ser usadas a favor de mecanismos de segurança. Essas informações são termos relacionados à desfiguração de páginas, nomes de grupos que normalmente estão nas desfigurações, identificadores de invasores e perfis de monitoramento, que são URLs para redes sociais. As redes sociais são importantes pois possibilitam realizar o monitoramento daquele indivíduo ou grupo, pois os mesmos podem orquestrar ataques através das redes sociais, além de que nem sempre uma desfiguração pode ser registrada no **Zone-H** por exemplo, apesar de não ser comum, um atacante pode optar por não registrar sua desfiguração.

### **3.3. Considerações do capítulo**

Neste capítulo foram apresentadas as questões de pesquisa e os métodos para investigação dessas questões. Um detalhamento de cada passo a ser tomado em cada processo do método foi apresentado e alguns exemplificados com códigos e consultas SQL. No Capítulo 4 apresentamos os resultados obtidos com a aplicação dos métodos que definimos neste capítulo, além de outros materiais relacionados.

## Resultados e Discussões

Esse capítulo descreve os resultados obtidos por meio do processamento da base de dados coletada no Zone-H, com aplicação do método e técnicas descritas no capítulo 3, além disso, apresenta uma análise dos resultados, e uma discussão dos pontos positivos e negativos dos métodos aplicados. Por fim, apresenta as respostas obtidas para as questões de pesquisa.

### 4.1. Caracterização dos resultados

A base de dados, coletada no Zone-H no período de 12 de janeiro de 2017 a 12 de novembro de 2017, possui as seguintes características:

- 7184 registros.
- 866 invasores distintos.
- 13 sistemas operacionais distintos.
- 18 servidores Web distintos.

A tabela 4.1 apresenta números de registros obtidos em cada mês.

**Tabela 4.1.** Quantidade de desfigurações por mês

Mês/Dias	Jan/19	Fev/28	Mar/31	Abr/30	Mai/31	Jun/30	Jul/31	Ago/31	Set/30	Out/31	Nov/12
QtdRegistro	38	1688	608	702	659	501	308	604	903	619	499

Observando a Tabela 4.1, verifica-se que janeiro está apenas com 19 dias e novembro com 12, devido ao período de coleta. Com esses dados, observa-se uma média diária de  $\simeq 23$  e uma média mensal de  $\simeq 648$  desfigurações. Lembrando que há registros descartados pois não foi possível monitorar devido as limitações do coletor com CAPTCHA e *popups*.

A Tabela 4.2 e as Figuras 4.1, 4.2 e 4.3 apresentam os grupos/indivíduos mais ativos no Brasil (aqueles que realizaram o maior número de desfigurações), caracterizam ainda os domínios que eles costumam atacar e analisam os dias da semana e os horários que mais ocorrem desfigurações.

**Tabela 4.2.** Principais grupos/indivíduos da base de dados criada

Invasor/Grupo	Quantidade de deface
VandaTheGod	461
MuhmadEmad	283
SA3D HaCk3D	214
Sillent_Attack	154
DARKRON	151
Tsunami Faction	145
Yunkers Crew	133
Zedan-Mrx	129
Tux Society	126
BRLZPoC	124
suliman_hacker	93
Ashiyane Digital Security Team	91
BALA SNIPER	88
XwoLfTn	83
Alarg53	82
GeNErAL	80
./CryptonKing	77
Umbrella Gang	75
Mo3Gza HaCkEr	73
ProtoWave Reloaded	73

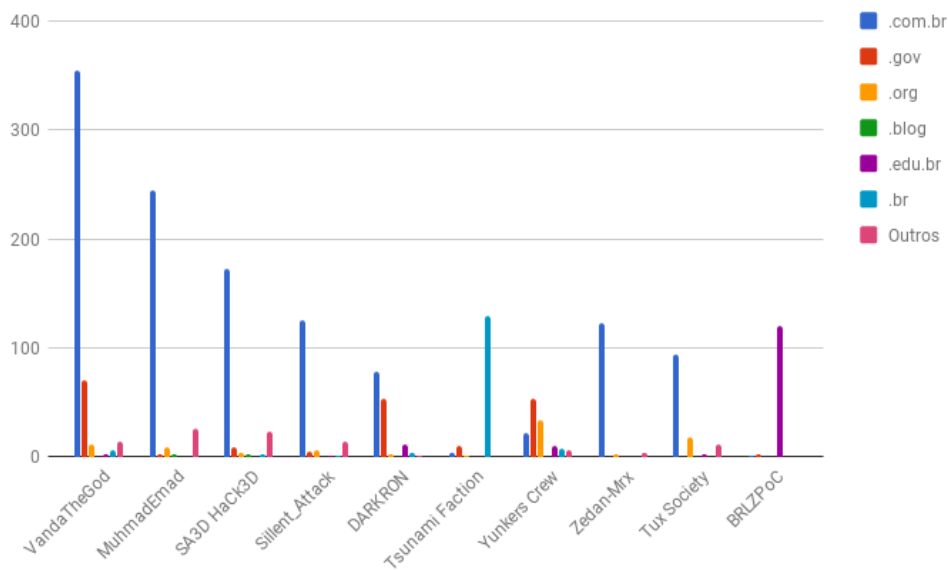
Observa-se na Tabela 4.2 que foram listados 20 grupos/indivíduos que mais realizaram desfigurações no período analisado. O grupo **VandaTheGod** além de aparecer no topo de números de desfigurações, se manteve constante na realização de seus ataques, ou seja, não foi apenas um *mass deface* em alguns dias, mas sim de mês em mês, desde abril até novembro. Diferente do grupo “SA3D HaCk3D” que realizou todas 214 desfigurações no mês de fevereiro e não apareceu mais. Portanto vale ressaltar o significado de ativo considerado pela pesquisa, que não quer dizer necessariamente um grupo que vem mantendo a realização de desfiguração.

Para ilustrar os grupos/indivíduos que vêm realizando pelo menos um ataque por mês e contendo desfigurações no mínimo em setembro, apresenta-se a Tabela 4.3.

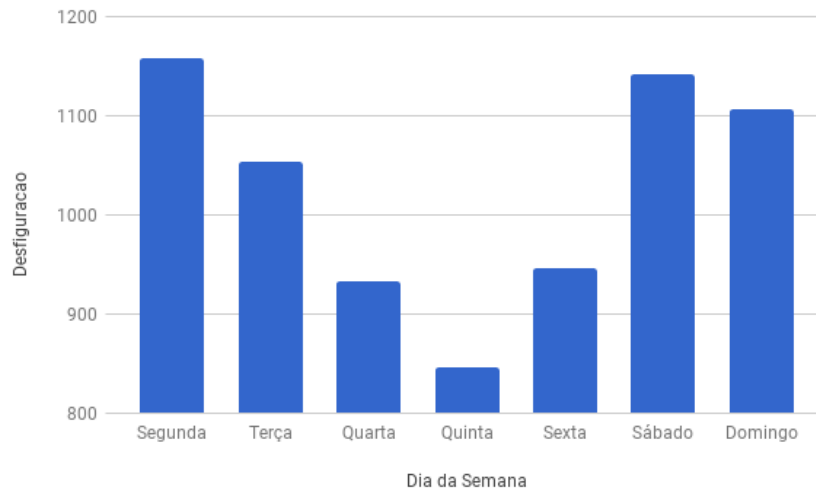
**Tabela 4.3.** Desfigurações mensais realizadas por grupos/indivíduos.

Invasor/Mês	Fev	Mar	Abr	Mai	Jun	Jul	Ago	Set	Out	Nov
VandaTheGod		x	x	x	x		x	x	x	x
MuhmadEmad	x			x		x		x	x	
Sillent_Attack				x	x	x	x	x		x
DARKRON	x	x	x	x			x	x	x	x
Zedan-Mrx	x	x	x	x	x		x	x		
Tux Society	x	x	x	x	x	x	x	x		
ProtoWave Reloaded		x		x	x	x	x	x	x	x

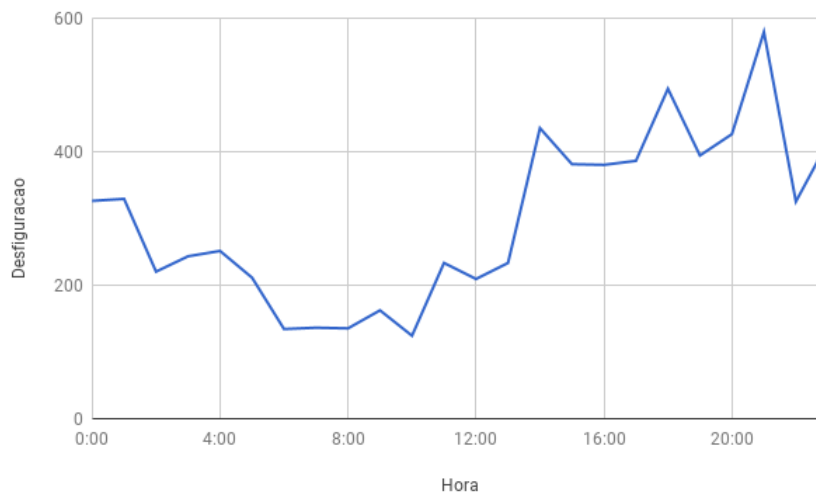
Caracterizado os grupos ativos foram extraídos os tipos de domínios em que os grupos costumavam realizar seus ataques, para tanto foram checados os seguintes domínios: «.com.br», «.gov», «.org», «.blog», «.edu.br», «.br (Universidades)» e «outros domínios». A Figura 4.1 apresenta resultado de alguns grupos citados como principais, e foi possível observar a tendência de domínios alvo nesse período, por exemplo, o grupo **Yunkers Crew** visivelmente tem o seu foco em domínios governamentais e o grupo **BRLZPoC** em domínios educacionais.

**Figura 4.1.** Gráfico caracterizando domínio por invasor/grupo

Os gráficos apresentados nas Figuras 4.2 e 4.3 apresentam os dias da semana e os horários mais comuns em que as desfigurações ocorrem.



**Figura 4.2.** Gráfico de desfiguração por dias da semana



**Figura 4.3.** Gráfico de desfiguração por horário

Observando os gráficos nas Figuras 4.2 e 4.3, verificou-se que apesar da perda de alguns registros devido a limitações do coletor, obteve-se um número considerável de desfigurações por dia e por mês, mostrando que esse ato é uma realidade na Internet brasileira e que existem muitos sítios vulneráveis a esse tipo de ataque. Notou-se os principais grupos/indivíduos e sua tendência de domínios alvo, onde a maioria acaba por visar domínios padrão “.com.br”. E apesar de alguns picos tanto nos dias da semanas quanto nos horários, essas duas características estão basicamente na mesma média.

## 4.2. Investigação de redes sociais

Aplicando a associação de palavras descrita na Subseção 3.2.6, identificou-se as redes sociais de alguns grupos a partir do processamento dos espelhos das páginas desfiguradas.

A seguir, a Tabela 4.4 apresentada alguns dos principais grupos e as respectivas redes sociais identificadas:

**Tabela 4.4.** Redes sociais de grupos/indivíduos

Grupo/Indivíduo	Rede Social	E-mail
Ashiyane Digital Security Team	twitter.com/shakerihassa	cyberhacker560@gmail.com
BALA SNIPER	facebook.com/balasniper007	balasniper17@gmail.com
DARKRON	twitter.com/@d4rkr0n	-
MuhmadEmad	-	kurdlinux007@gmail.com
ProtoWave Reloaded	facebook.com/pwave01	-
Sillent_Attack	facebook.com/sillent.attack skype: live:gliphacking	-
Tsunami Faction	fb.com/tsunamifaction	-
Tux Society	facebook.com/tuxxsociety	-
Umbrella Gang	twitter.com/hanneswho hanneswho@nigge.rsirc.priv8.co	-
VandaTheGod	twitter.com/vandathegod facebook.com/BrazilianCyberArmy irc.privbr.com	-
Yunkers Crew	www.facebook.com/yunkers01	-
Zedan-Mrx	skype: live:zedan-mrx skype: scan-suisse	-

Todas as redes sociais apresentadas foram validadas e estão ativas, apenas os e-mails não foram validados, devido a questões de sigilo da investigação não enviamos mensagens aos invasores. Apesar de apresentar as redes sociais do principais invasores, o processamento em toda base obteve rede sociais de 72 grupos/indivíduos.

## 4.3. Análise de características e padrões de ataques

Para a análise de características e padrões de ataques individuais e gerais, foi aplicado, como já explicado na Subseção 3.2.6, uma análise de frequência para extrair os termos mais

utilizados por um grupo/indivíduo e a partir disso, explorar entidades, associação de palavras e similaridades em busca desses padrões. A aplicação de janelas de tempo se fez necessária, para identificar se esse padrão varia conforme o tempo.

Aplicada essa análise, foi possível observar, por exemplo, que o grupo SA3D HaCk3D na maioria dos ataques (68,7% - 147 de 214 desfigurações), o primeiro texto do conteúdo é: “Hacked By SA3D HaCk3D”, isso foi observado através do HTML normalizado que não possui marcação HTML. Após a identificação do padrão, foi analisado como esse era apresentado no código HTML, e então foi possível concluir que esse conteúdo padrão se encontra no título da página Web, portanto, uma característica desse grupo é modificar o título da página e inserir o texto: “Hacked By SA3D HaCk3D”.

Grupos como o `"/Cryptonking"` frequentemente o padrão é específico por membro. Por meio de teste de similaridade e de uma análise no HTML normalizado foi extraído o padrão de cada um. Por exemplo, o invasor `"Sh40Cr1m1n0s0"` em todos seus ataques deixou a mensagem do quadro:

*“ brasil ta foda parceiro seu site acaba de ser penetrado #1337 por ./cryptonking  
 ℰℰ sh40cr1m1n0s0 se liga rs kkkkk to pesadao em vcs ”*

O `"Suicide Ghost"` nos 21 e 24 de abril realizou 4 desfigurações com a mensagem:

*“ vishe prayboy .. :: suicide ghost :: .. ./cryptonking vunse voi rakeadu @\_@ ”*

Depois, nos dias 26 e 27, realizou uma mudança simples, porém mantendo o padrão de colocar o nome e o grupo:

*“ .. ;; suicide ghost ;; .. { ./cryptonking } ”*

O grupo `"vbsdz17"`, apesar de não aparecer entre os mais ativos, entre o fim de fevereiro (26/02) e começo de março (09/03), realizou 18 desfigurações, todas com o mesmo padrão:

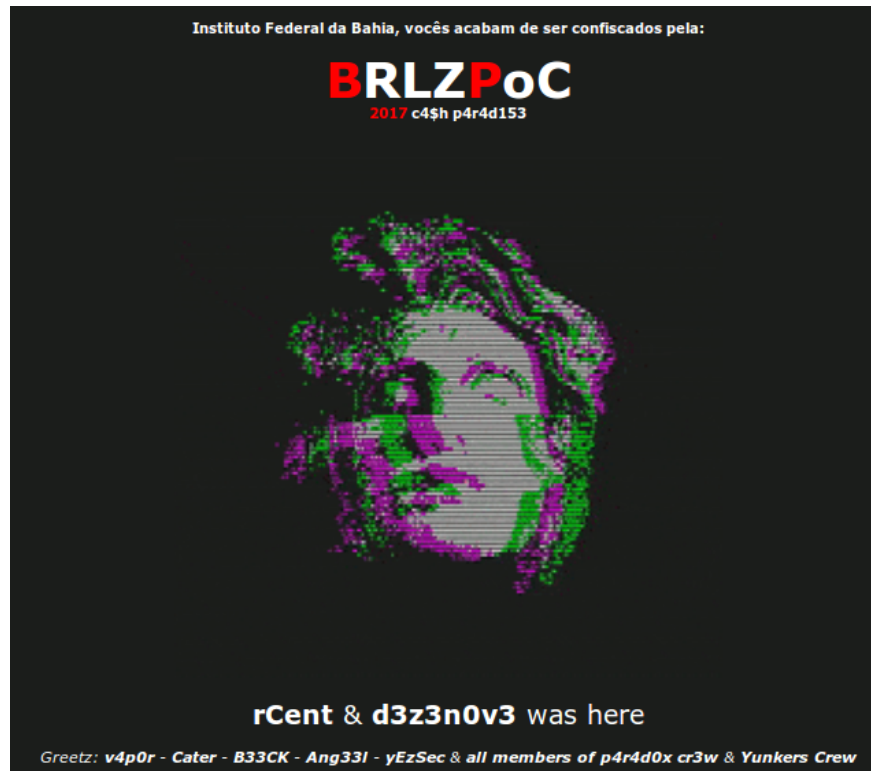
*“ ... -=[dz.pirate]=- . ba3kdor-dz security is just an illusion! we are-=[#]  
 dz.pirate [#]=-: fb : 'facebook.com/vbsdz17' | it's better to have a good security  
 no more, than fooling everyone including your selves ℰ have one. | #humbug • it  
 only makes us slower. greets:vbsdz17.i.and to all hacker groups out there. -=[  
 free ] :=-©/ vbsdz17 greetz : who am i ? -vbsdz17 ./ vendetta-dz ./imad.dz ”*

Após isso o grupo voltou a realizar ataques apenas em setembro e outubro com um padrão diferente:



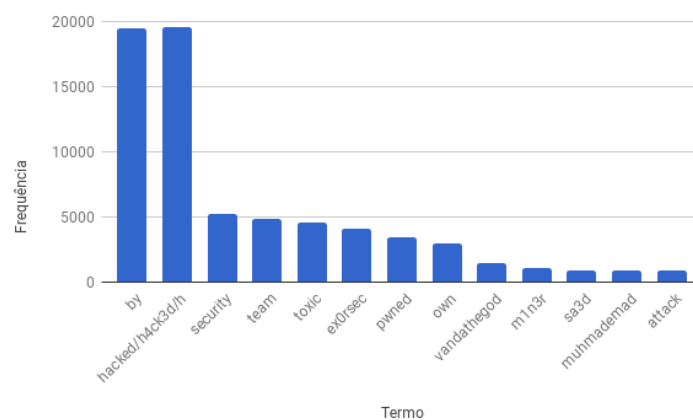
“ ...::: vbsdz17 :::... hacked by vbsdz17 ...:::my friend : imad :p::... gretz to :all  
ackers contack : hacked by vbsdz17 ”

O grupo "BRLZPoC" tem como padrão a inserção de um GIF, o uso da palavra “confiscado”, um pouco diferente da maioria, e apresentação de integrantes. A Figura 4.4 representa o padrão do grupo.



**Figura 4.4.** Padrão de desfiguração do grupo BRLZPoC

Dado os padrões individuais, obteve-se também padrões gerais. Através da análise de frequência de palavras, foi possível obter os números apresentados na Figura 4.5.



**Figura 4.5.** Termos mais frequentes nas desfigurações

Observa-se que o número de vezes que as palavras “by” e “hacked” (e derivados como: “h4ck3d e “hack3d”) aparecem é muito grande (quase 20 mil vezes cada uma), sendo isso comum, já que “hacked by” é o termo mais comum de ser visto em desfigurações ou ataque qualquer. Foi possível observar outros termos comuns de aparecerem como: “pwned” e “own”, que seria um sinônimo de “hacked”. Outros termos que aparecem, como: “ex0rsec”, “vandathegod”, “m1n3r”, “sa3d” e “muhammadmad”, são grupos/indivíduos, ilustrando um padrão de assinar a desfiguração com o *nickname* ou nome do próprio grupo.

Fazendo uma análise sobre esses principais termos, extraiu-se a frequência relativa de cada um deles em relação ao total de desfigurações (7184). A Tabela 4.5 apresenta a frequência relativa calculada.

**Tabela 4.5.** Frequência de aparição de um termo nas desfigurações

Termo	Ocorrência do Termo	%
hacked/h4ck3d/hack3d	3350	46,63
by	4296	59,80
security	517	7,20
team	977	13,60
toxic	130	1,81
ex0rsec	174	2,42
pwned	112	1,56
own	827	11,51
vandathegod	452	6,29
m1n3r	80	1,11
sa3d	190	2,64
muhammadmad	287	3,99
attack	258	3,59

Observa-se na Tabela 4.5 que alguns termos possuem uma variação em sua escrita, recebendo números no meio da palavra, por exemplo, “h4ck3d”, “ex0rsec”, “m1n3r”. Esse padrão bem comum na Internet, em ataques e em nomes de grupo/indivíduos é o **idioma leet**, que é uma alternativa ao alfabeto latino onde a forma de escrever não recebe apenas letras, mas também símbolos e números (Chaves, 2010).

Ainda observa-se na Tabela 4.5 que apesar de ser citado que o termo “hacked” e suas variações tenha aparecido muitas vezes, ao recuperar a frequência em que este ocorre nas desfigurações, nota-se que o termo aparece em menos da metade das desfigurações, sendo necessário um olhar mais criterioso para conclusões, ou seja, verificando apenas a presença desses termos não é possível concluir que ocorreu uma desfiguração.

## 4.4. Avaliação das questões de pesquisa

Essa seção apresenta as respostas para as questões de pesquisa, considerando os resultados apresentados e discutidos nas Seções 4.1, 4.2 e 4.3

### Q1. Quais os grupos/indivíduos hackers mais ativos em desfiguração no Brasil?

Os grupos/indivíduos mais ativos em desfigurações são listados na Tabela 4.2. Essa tabela apresenta os 20 principais grupos/indivíduos e o Apêndice C apresenta a classificação completa considerando todos os grupos da base coletada que realizaram pelo menos 10 desfigurações. Em complemento, observa-se na Figura 4.1 os principais alvos dos atacantes no período monitorado.

Apesar de apresentados os 20 principais grupos/indivíduos, nem todos estão realizando ataques frequentemente, por isso, além dos 20 principais, a Tabela 4.3 caracterizou dentre esses 20, aqueles que realizam desfigurações mensalmente.

A investigação dessa questão, resultou na obtenção de informações relevantes que podem ser utilizadas no desenvolvimento de mecanismos para identificação de desfigurações, além de contribuição com mecanismos de detecção automática de desfiguração, uma vez que apresentam-se os grupo que estão realizando desfigurações, aqueles que mais tem feito esses ataques, os domínios que eles tendem a atacar e os horários e dias da semanas que mais ocorrem essas ações.

### Q2. É possível relacionar os grupos/indivíduos hackers com perfis ou páginas em redes sociais?

Sim, é possível relacionar os grupos/indivíduos *hackers* com perfis ou páginas em redes sociais. Conforme foi apresentado na Seção 4.2, em mais da metade (12 grupos) dos 20 principais grupos foi identificado ao menos uma rede social. Nessas redes, verificou-se que os grupos/indivíduos divulgam suas desfigurações, por exemplo, nos seguintes perfis:

- D4RKR0N: [twitter.com/@d4rkr0n](https://twitter.com/d4rkr0n)
- ProtoWave: [facebook.com/pwave01](https://facebook.com/pwave01)
- Sillent Attack: [facebook.com/Sillent.Attack](https://facebook.com/Sillent.Attack)
- Tsunami Faction: [facebook.com/tsunamifaction](https://facebook.com/tsunamifaction)
- Tux Society: [facebook.com/tuxxsociety](https://facebook.com/tuxxsociety)
- Umbrella Gang: [twitter.com/hanneswho](https://twitter.com/hanneswho)
- VandaTheGod: <https://twitter.com/vandathegod> e [facebook.com/BrazilianCyberArmy](https://facebook.com/BrazilianCyberArmy)

- Yunkers Crew: [facebook.com/yunkers01](https://facebook.com/yunkers01)

O monitoramento desses perfis possibilitaria rastrear mais rapidamente as desfigurações com o intuito de uma reação mais rápida. Nos perfis, não foi observado a orquestração de ataque por meio da rede social, apesar da interação dos atacantes com os seguidores. Logo, como há possibilidade de organização de ataques através dos indivíduos, o monitoramento desses perfis é importante.

**Q3. É possível estabelecer um padrão de desfiguração (características na página) segundo um grupo/indivíduo?**

Sim, é possível estabelecer um padrão de desfiguração segundo um grupo/indivíduo. Observa-se na Seção 4.3 que através da aplicação de análise de frequência, associação de palavras e similaridade, resgatou-se padrões de grupos/indivíduos, os quais foram confirmados por um especialista.

Como exemplo de padrão, foi possível observar o grupo “SA3D HaCk3D” que tinha como característica alterar o título da páginas desfiguradas, inserindo o padrão “Hacked By SA3D HaCk3D”.

**Q4. É possível identificar um padrão nas páginas para caracterizar uma desfiguração?**

Sim, é possível identificar um padrão nas páginas desfiguradas. Por exemplo, A Figura 4.5 apresenta o padrão “hacked” “by” que apareceu em 46,63% e 59,80% das desfigurações respectivamente, tendo um número relevante de aparições. É importante a análise dos mesmos sempre com a associação dos termos entre eles ou com os demais termos presentes na figura ou com *nicknames* recuperados.

Apesar de apenas o termo “by” aparecer em mais da metade dos registros de desfigurações, os demais termos são muito específicos, ou seja, caso um deles esteja presente no corpo de um sítio Web, a chance dessa página estar desfigurada é grande.

## 4.5. Avaliação dos procedimentos

Essa seção discute os procedimentos utilizados na pesquisa e apresenta problemas encontrados durante o seu desenvolvimento. Os principais problemas encontrados, em específico no coletor, foram o CAPTCHA (um teste de desafio cognitivo, utilizado como

ferramenta anti-spam) e o *popup* (uma janela que abre no navegador Web). Esses problemas são apresentados a seguir, além de outras dificuldades.

O primeiro passo da pesquisa consistiu na criação da base de dados. Pretendíamos obter a base diretamente do Zone-H, porém, ao entrar em contato com o administrador, o mesmo informou que não era possível fornecer a base de dados. Com isso, foi dado início a implementação do coletor automático (Web Crawler).

Antes de utilizar o Selenium para resgatar o HTML e as informações necessárias, foi utilizado métodos padrões como GET e seus derivados, porém, todas sem sucesso, pois através da requisição que era feita para o Zone-H, a resposta obtida não era o HTML desejado do sítio (isso era esperado pois o administrador informou que o sítio continha algumas técnicas contra Web Crawler). Porém como o Selenium habilita o controle de um Navegador e “simula” um humano operando sobre o sítio, assim, foi possível a manipulação e navegação no HTML.

O CAPTCHA: quando era definido o domínio .BR (pois era desejado apenas os sítios Web brasileiros) o Zone-H solicitava autenticação de um CAPTCHA o que impedia a automatização da coleta, para isso, foi descoberto que o Zone-H possui um sítio brasileiro que contém apenas desfigurações com o domínio .BR. Porém isso continua sendo um problema, sendo que em um determinado período de tempo o CAPTCHA é solicitado novamente.

Outra dificuldade foi em relação ao domínio, nas primeiras versões do Web Crawler, o mesmo não acessava a cópia da desfiguração, acessava apenas a lista de desfiguração que já apresentava os campos citados acima, porém o domínio de alguns sítios era extenso e não exibia por completo, sendo utilizado “...”. Portanto foi necessário, a partir disso, acessarmos a cópia da desfiguração.

Haviam registros que eram feitos exatamente na mesma hora, assim, o coletor dava o registro como duplicado, pois já existia um registro com aquela data no banco de dados. Para isso foi adicionada uma comparação, para verificar o domínio, portanto, se aquele registro que continha horas iguais tivesse o domínio diferente, ele ainda não tinha sido armazenado no banco.

Como o campo `dataRegistro` é uma chave primária no banco de dados, ao corrigir o problema citado de registro duplicado, mais um problema surgiu, como já existia registro com aquela data, no momento de salvar, o registro (ou a chave) já existia e apesar do domínio ser diferente ele não permitia a inserção, pois a chave seria duplicada, para isso o campo `dataRegistro` foi modificado e deixou de ser uma chave primária.

O coletor salva todos os *mirrors* presentes em uma página e visita todos eles, passando assim para a próxima página caso não tenha nenhum registro duplicado, porém se algum registro novo for notificado, um dos registros já salvo irá para outra página e o coletor ao

compará-lo notifica que o registro já existe e interpreta que o restante também já foi coletado, porém pode existir registros que ainda não foram salvos a partir daquele que foi pra a próxima página. A solução foi, não parar a coleta com apenas o alerta de um registro duplicado, mas sim ter uma “margem de erro” para cobrir esses casos, parando a execução da coleta apenas quando houver 5 registros duplicados.

Quando o invasor insere algum *popup* o coletor não reage bem a essa ação, para isso foi verificada a exceção que era lançada e adicionado *try/exceptions* para quando a exceção do popup for lançada a coleta continua para os próximos registros. Apesar de finalizado o coletor ainda precisa de melhorias.

Os primeiros procedimentos realizados para investigação das questões de pesquisa com o intuito de encontrar os principais grupos/indivíduos foram padrões e eficientes. Padrões e eficientes pois consistiam em consultas SQL e como se tratava de uma base de dados bem estruturada, foi a melhor abordagem.

A extração de redes sociais obteve resultados desejados, porém de uma forma não tão eficiente. Tanto utilizando a biblioteca NLTK, quanto o algoritmo desenvolvido, os resultados devolvidos continham muita informação irrelevante, sendo necessário que o especialista removesse essas informações manualmente. Uma forma mais eficiente e automatizada seria a construção de uma expressão regular com mais restrições na seleção das redes sociais, além de uma verificação que comparasse as redes sociais resgatadas com o grupo/indivíduo, analisando uma similaridade entre os dois.

Os procedimentos aplicados para extração de padrões individuais e gerais, apesar de utilizar alguns métodos já utilizados na extração de redes sociais, como associação de palavras, quando somada a extração de termos mais frequentes e quando aplicado janelas de tempo, obteve-se resultados importantes para a caracterização dos grupos/indivíduos e extração de seus padrões.

A análise de similaridade foi aplicada para expressar essa semelhança entre as desfigurações através de números, porém o processo de analisar essa igualdade foi dada ao especialista, tornando o método não tão eficiente quanto podia. Era possível uma combinação de todas essas técnicas, tornando as caracterizações melhores e automatizadas.

## 4.6. Considerações do capítulo

---

## Conclusões e Trabalhos Futuros

---

Neste trabalho foram identificados os principais grupos/indivíduos que atuam em desfiguração de páginas no Brasil, como também foram extraídas informações para o monitoramento desses grupos/indivíduos, como perfis de redes sociais e padrões de desfigurações. Além disso, também foram caracterizadas outras questões interessantes, como a tendência de ataques a domínios, dias da semana e hora que comumente são publicadas as desfigurações.

Em relação à identificação dos perfis dos atacantes nas redes sociais, observou-se que há divulgação delas nas desfigurações, além de que, ataques são divulgados através dessas redes sociais. Portanto, uma vez que resgatadas, obtém-se informações de grande importância para mecanismos de monitoramento e de detecção antecipada.

Quanto a padrões de desfigurações por grupos, notou-se que esses padrões existem e podem variar conforme o integrante que realiza o ataque. Nota-se também que há desfigurações parciais e completas, muitos *blogs*, sítios de notícias ou sítios de prefeituras recebem apenas a inserção de uma postagem dentro de seus “painéis de notícias” e outros realmente sofrem alterações do sítio por completo.

O Zone-H mostrou-se uma excelente fonte de extração de informações relevantes, devido ao grande número de desfigurações que são divulgadas nele, criando assim uma rica fonte de base de dados. No entanto, foi uma limitação por conta do difícil acesso às informações.

Como uma das principais contribuições, os resultados dessa pesquisa podem ser usados em mecanismos de alerta antecipado, já que foi gerado uma base de informações inteligentes. Redes sociais já foram disponibilizadas para o projeto GT-EWS<sup>1</sup> que é um sistema de alerta antecipado e que monitora essas redes. Além das redes sociais, o coletor desenvolvido foi

---

<sup>1</sup> <https://gtews.ime.usp.br/>

adaptado e está sendo usado como um sensor no Hórus<sup>2</sup> dentro do GT-EWS.

Ainda há desafios para a pesquisa, como a questão de melhorar o processo automático de identificação de redes sociais e extração de padrões, como já comentado na Seção 4.5, existe uma dependência da análise especialista futuramente pode ser reduzido por melhoramento no uso de técnicas de Aprendizado de Máquina e novas heurísticas.

Como trabalhos futuros são propostos: (i) Automatização geral dos processos de extração de redes sociais e padrões de desfigurações. (ii) Solucionar limitações do coletor, que seria resolver o problema do CAPTCHA e dos *popups*, pois além de não ter a coleta interrompida, desfigurações não serão perdidas. (iii) Fazer utilização do HTML original em conjunto com o espelho da desfiguração, para auxiliar na criação de um mecanismo que detecte a permanência da desfiguração no sítio original, ou seja, no momento da coleta, o espelho do sítio original e o espelho da desfiguração passaria por uma análise de similaridade e em caso de uma porcentagem alta de similaridade, emite-se um alerta para o sítio, para que ocorra um reparo no mesmo.

---

<sup>2</sup> <https://horus.rnp.br/>



# Referências

---

- AMARAL, Daniela O. F. do; VIEIRA, Renata. O reconhecimento de entidades nomeadas por meio de conditional random fields para a língua portuguesa. In: , 2013. p. 59–68.
- BARTOLI, Alberto; DAVANZO, Giorgio; MEDVET, Eric. The reaction time to web site defacements. *IEEE Internet Computing*, v. 13, n. 4, p. 52–58, 2009. ISSN 10897801.
- BENJAMIN, Victor; LI, Weifeng; HOLT, Thomas; CHEN, Hsinchun. Exploring threats and vulnerabilities in hacker web: Forums, IRC and carding shops. *2015 IEEE International Conference on Intelligence and Security Informatics: Securing the World through an Alignment of Technology, Intelligence, Humans and Organizations, ISI 2015*, p. 85–90, 2015.
- BERTOGLIO, Daniel Dalalana.; ZORZO, Avelino Francisco. *Um Mapeamento Sistemático sobre Testes de Penetração*. Dissertação (Mestrado) — PONTIFÍCIA UNIVERSIDADE CATÓLICA, 2015.
- CAMPIOLO, Rodrigo. *Análise e extração de aleras antecipados sobre ameaças e incidentes de segurança em sistemas computacionais usando fontes de dados não estruturados*. Tese (Doutorado) — Instituto de Matemática e Estatística da Universidade de São Paulo, Setembro 2016.
- CERT.BR. *Incidentes Reportados ao CERT.br – Janeiro a Dezembro de 2016*. Julho 2017. <https://www.cert.br/stats/incidentes/2016-jan-dec/total.html>, disponível em 14/07/2017, acessado em 24/11/2017.
- CHAVES, Edu. *Leet (1337) - linguagem secreta dos hackers*. Abril 2010. <http://www.sequelanet.com.br/2010/04/leet-1337-linguagem-secreta-dos-hackers.html>, Disponível em 21/04/2010. Acessado em 24/11/2017.
- CONRADO, Merley da Silva. *O efeito do uso de diferentes formas de extração de termos na compreensibilidade e a representatividade dos termos em coleções textuais na língua portuguesa*. Dissertação (Mestrado) — Instituto de Ciências Matemáticas e de Computação - ICMC-USP, 2009.

- DAVANZO, Giorgio; MEDVET, Eric; BARTOLI, Alberto. A comparative study of anomaly detection techniques in web site defacement detection. In: *IFIP International Federation for Information Processing*, 2008. v. 278, p. 711–716. ISBN 9780387096988. ISSN 15715736.
- FELDMAN, Ronen; SANGER, James. *The Text Mining Handbook*. [S.l.: s.n.], 2007. 423 p. ISSN 14653133. ISBN 978-0-511-33507-5.
- HACKING, Loops. *6 Ways to Hack or Deface Websites Online*. Janeiro 2016. <https://www.hackingloops.com/6-ways-to-hack-or-deface-websites-online/>, Disponível em 04/2016. Acessado em 24/11/2017.
- HENRY, Kevin. *Penetration Testing: Protecting Networks and Systems*. [S.l.]: IT Governance Publishing, 2012.
- JOHARI, Rahul; SHARMA, Pankaj. A survey on web application vulnerabilities (SQLIA, XSS) exploitation and security engine for SQL injection. In: *Proceedings - International Conference on Communication Systems and Network Technologies, CSNT 2012*, 2012. p. 453–458. ISBN 9780769546926.
- LI, W.; CHEN, H. Identifying top sellers in underground economy using deep learning-based sentiment analysis. In: *2014 IEEE Joint Intelligence and Security Informatics Conference*, 2014. p. 64–67.
- MANNING, C. D.; RAGHAVAN, P.; SCHUTZE, H. Language models for information retrieval. In: PRESS, Cambridge University (Ed.). *An Introduction to Information Retrieval*. [S.l.: s.n.], 2008.
- MORAIS, Edison Andrade Martins; AMBRÓSIO, Ana Paula L. Mineração de textos. *Relatório Técnico-Instituto de Informática (UFG)*, 2007.
- MOTA, C.; SANTOS, D.; RANCHHOD, E. “avaliação de reconhecimento de entidades mencionadas: Princípio de harem”. In: PRESS, IST (Ed.). *Avaliação Conjunta: Um novo paradigma no processamento computacional da língua portuguesa*. [S.l.]: Diana Santos, 2007. p. capítulo 14, p. 161–176.
- NAHM, Un Yong; MOONEY, Raymond J. Text mining with information extraction. *AAAI Technical Report SS*, SS-02-06, p. 60–67, 2002.
- OWASP. *OWASP Top 10 Application Security Risks - 2010*. Abril 2010. [https://www.owasp.org/index.php/Top\\_10\\_2010-Main](https://www.owasp.org/index.php/Top_10_2010-Main), Disponível em 26/04/2010. Acessado em 24/11/2017.
- OWASP. *Cross-site Scripting (XSS)*. Abril 2016. [https://www.owasp.org/index.php/Cross-site\\_Scripting\\_\(XSS\)](https://www.owasp.org/index.php/Cross-site_Scripting_(XSS)), Disponível em 06/04/2016. Acessado em 24/11/2017.

PREATONI, Roberto. *Zone-H*. Setembro 2017. <https://en.wikipedia.org/wiki/Zone-H>, Disponível em 23/09/2017. Acessado em 24/11/2017.

SALATIEL, José Renato. *Crimes virtuais: Hackers promovem onda de ataques no Brasil*. Julho 2011. <https://vestibular.uol.com.br/resumo-das-disciplinas/atualidades/crimes-virtuais-hackers-promovem-onda-de-ataques-no-brasil.htm>, Disponível em 01/07/2011. Acessado em 24/11/2017.

SANTOS, Maria Angela Moscalewski Roveredo dos. *Extraíndo regras de associação a partir de textos*. Curitiba, 2002. 51 p. Dissertação (Mestrado) — Pontifícia Universidade Católica do Paraná, 2002.

SUREKA, A.; GOYAL, V.; CORREA, D.; MONDAL, A. Polarity classification of subjective words using common-sense knowledge-base. *Lecture Notes in Computer Science Rough Sets, Fuzzy Sets, Data Mining and Granular Computing*, p. 486–493, 2009.

WENDT, Emerson. *Os ataques crackers no Brasil e seus efeitos penais decorrentes*. Junho 2011. <http://www.emersonwendt.com.br/2011/06/os-ataques-crackers-no-brasil-e-seus.html>, Disponível em 26/06/2011. Acessado em 24/11/2017.

ZONE-H. *Estatísticas geral Anual/Mensal/diária*. Novembro 2017. <http://zone-h.com.br/stats/ynd>, Disponível em 01/11/2017. Acessado em 24/11/2017.

# Apêndices

# Apêndice A: Web Crawler

---

O Código 1, apresenta o coletor automático desenvolvido para o estudo, que realiza Web Crawling no sítio Zone-H.

```
def crawler():
    copia = [] #variavel pra salvar os mirrors
    print("Aguarde, abrindo firefox e acessando o site...")
    wd_firefox = webdriver.Firefox() # abro o firefox

    for i in range(1, 51):
        del copia[:]
        wd_firefox.get
            ("http://br.zone-h.org/archive/published=0/page=" + str(i))
        source_html = lxml.html.fromstring(wd_firefox.page_source)
        # Pegando acesso a mirror, e salvando todos os mirrors a lista copia

        print("Salvando mirrors... Pagina: " + str(i))

        for coluna in source_html.xpath("./table[@id='ldeface']/tr"):
            for linha in coluna.xpath("./td//a/@href"):
                if "/mirror" in linha:
                    copia.append("http://br.zone-h.org" + linha)

    #Navegando em todos os mirrors
    for mirror in copia:
        try:
            time.sleep(60)
            wd_firefox.get(mirror)
            html_mirror = lxml.html.fromstring(wd_firefox.page_source)

            data = html_mirror.xpath("//*[@id='propdeface']/ul/li[1]/text()")
            data = formatData(data)
            verificaD = verificaData(data)

            if(verificaD == 1):
                dominio = html_mirror.xpath
                    ("//*[@id='propdeface']/ul/li[2]/ul/li[2]/a/text()")
                invasor = html_mirror.xpath
                    ("//*[@id='propdeface']/ul/li[2]/ul/li[1]/text()")
                sistema = html_mirror.xpath
                    ("//*[@id='propdeface']/ul/li[3]/ul/li[1]/text()")
                end_ip = html_mirror.xpath
                    ("//*[@id='propdeface']/ul/li[2]/ul/li[3]/text()")
                web_service = html_mirror.xpath
                    ("//*[@id='propdeface']/ul/li[3]/ul/li[2]/text()")
                #dominio do site desfigurado no mirror
                html_deface = html_mirror.xpath
```

```

        ("//*[id='propdeface']/iframe/@src")
#Capturando html do site desfigurado no mirror
wd_firefox.get(html_deface)
content = wd_firefox.page_source
html_deface = BeautifulSoup(content, "lxml")
html_deface = str(html_deface)

#Acessando site desfigurado para captura do html
wd_firefox.get(dominio)
content = wd_firefox.page_source
html_dominio = BeautifulSoup(content, "lxml")
html_dominio = str(html_dominio)
#Observacoes:
#html_dominio = html do site desfigurado
#html_deface = html do site no mirror

saveDB(data, str(invasor), str(dominio), str(sistema), str(end_ip),
        str(web_service), html_dominio, html_deface)
else:
    dominio = html_mirror.xpath
        ("//*[id='propdeface']/ul/li[2]/ul/li[2]/a/text()")
    verificaDo = verificaDominio(str(dominio), data)
    if(verificaDo == 1):
        invasor = html_mirror.xpath
            ("//*[id='propdeface']/ul/li[2]/ul/li[1]/text()")
        sistema = html_mirror.xpath
            ("//*[id='propdeface']/ul/li[3]/ul/li[1]/text()")
        end_ip = html_mirror.xpath
            ("//*[id='propdeface']/ul/li[2]/ul/li[3]/text()")
        web_service = html_mirror.xpath
            ("//*[id='propdeface']/ul/li[3]/ul/li[2]/text()")
        html_deface = html_mirror.xpath
            ("//*[id='propdeface']/iframe/@src")
        wd_firefox.get(html_deface)
        content = wd_firefox.page_source
        html_deface = BeautifulSoup(content, "lxml")
        html_deface = str(html_deface)

#Acessando site desfigurado para captura do html
wd_firefox.get(dominio)
content = wd_firefox.page_source
html_dominio = BeautifulSoup(content, "lxml")
html_dominio = str(html_dominio)

saveDB(data, str(invasor), str(dominio), str(sistema), str(end_ip),
        str(web_service), html_dominio, html_deface)
else:
    print("Registro duplicado!")
    #break;
except ValueError:
    print("Captcha solicitado")
    exit(0);
except:
    print("Popup")
wd_firefox.close()

```

**Código 1.** Coletor automático de desfigurações no Zone-H

# Apêndice B: Estrutura da base de dados

---

A base de dados está estruturada da seguinte forma:

**Nome da base de dados:** zonehDB;

**Tabela existente na base:** dataDefacers;

- Campos presentes na tabela

- dataRegistro = Data do Registro

- dataRegistro = Data do Registro

- invasor = *nick* do invasor

- dominio = domínio atacado

- enderecoIP = endereço IP atacado

- webService = Serviço Web

- html = corresponde ao html do sítio desfigurado

- htmlMirror = corresponde ao espelho da desfiguração registrado no Zone-H

# Apêndice C: Invasores

---

A Tabela 1 apresenta os grupos/invasores presentes na base de dados da pesquisa que realizaram pelo menos 10 desfigurações.

<b>Grupo/Indivíduo</b>	<b>Quantidade de desfiguração</b>
VandaTheGod	461
MuhmadEmad	283
SA3D HaCk3D	214
Sillent_Attack	154
DARKRON	151
Tsunami Faction	145
Yunkers Crew	133
Zedan-Mrx	129
Tux Society	126
BRLZPoC	124
suliman_hacker	93
Ashiyane Digital Security Team	91
BALA SNIPER	88
XwoLfTn	83
Alarg53	82
GeNErAL	80
./CryptonKing	77
Umbrella Gang	75
Mo3Gza HaCkEr	73
ProtoWave Reloaded	73
CyberTeam	69
Mister Spy	67
4Ri3 60ndr0n9	66
HolaKo	65
Imam	63
Sons of Anarchy	63



ZoRRoKiN	62
TeaM_CC	61
Cooldsec	59
Strike King	56
Spy_Unkn0wn	53
Xinox Crew	52
MonstersDefacers-Padocas	49
AB15 Team	47
chinafans	47
RxR	46
m1n3r	46
Ex0rsec	44
w4l3XzY3	43
ayyildiz	41
Itachi.sz	41
jok3r	39
Collapse Gang	39
Mr ER	38
BILGEKULTIGIN	38
UserGhost	38
Default	35
fast	34
Kashif HaxOr	33
aDriv4	32
Brazilian Cyber Army	32
vbsdz17	29
Cater	29
Anarchy Ghost	29
@Sprek3rsSec	28
GAZA	28
Dr.S4mom	27
CyBeRiZM	27
Toxic Security Team	26
TheWayEnd	24
HighTech	23
Tr3v0r	23

VM	23
SynnX	23
Ex0rcist	23
spl0it3r	23
EvilBoyz	22
GHoST61	22
Colder	22
AnoaGhost	21
Astrologyc Hack Team	20
13CHMOD37	20
Anon Ghost Portugal	20
magelang6etar	20
Pak MONster	20
dkr	19
etownteam	19
CandySec	19
Aj4x	18
0x1999	18
SilentAngel	18
By_uMuT	18
Mr Virus Dz	18
nginxDEX	17
Fallaga Team	17
hamaminho	17
Sxtz	17
Dr.SiLnT Hill	17
BrazilObscure	17
Ayy0131ld0131z Tim	17
aPTx4869	16
black hell ahmed	16
Hitler El Maghribi	16
KingSkrupellos	16
An0n 3xPloiTeR	16
Mr.Medo	16
Mr.Vangke404	16
Unknown Al	15

cyber_hunter	15
adam tnx	15
Anonymous Indonesia	15
N3X0000S	15
Ramil Feyziyev	14
TeaMGhost	14
Mr.XaaD	14
Moroccan Revolution	14
nikkO	14
mtz13	14
LOST3R	14
h4lyz0r1337	14
xin0x w0	13
Shade	13
R3tr1ng	13
r00tkit	13
v4p0r	13
Index Php	13
b33ck	13
Mr.Rizgar.halshoy.kurdish.blackhat	12
ArchS3x	12
kefiex404	12
Panataran	12
veryhax	12
MiLwrOM_Dz	12
folps	12
E7	11
Monsters Defacers	11
Santi boy	11
Ara-C@esar	11
SNAKE2K1	11
zakiloup	11
Mr.DreamX196	11
NeT.Defacer	10
AbsenceTM	10
Hacker Khan	10

KkK1337	10
pr0s3x	10
Chucky.sh	10

**Tabela 1.** Relação de grupos/invasores presentes na base de dados