

時系列アクセスパターンのテンプレートを用いたマルウェア感染端末検知

外口 大凱*, 大下 裕一*, 芝原 俊樹†, 千葉 大紀†,
秋山 滉昭†, 村田 正幸*

* 大阪大学 大学院情報科学研究科
† NTT セキュアプラットフォーム研究所

研究背景

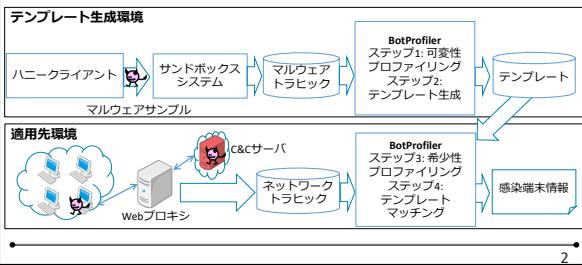
- ・マルウェア感染端末検知は重要な課題
 - ・攻撃者はC&Cサーバ経由で感染端末を制御
 - ・感染が拡大することで攻撃者は大規模な攻撃が可能
- ・2種類の感染端末検知手法が存在
 - ・ホストベース: 端末上でマルウェアが動作していないか検出するプログラムを実行
 - ・端末が攻撃者によって制御されている場合、プログラム自体が停止
 - ・ネットワークベース: 通信を検査することで感染端末を検知
- ・HTTP通信に着目した感染端末検知
 - ・C&CサーバはHTTP通信を用いて感染端末を制御
 - ・C&Cサーバが行うHTTP通信をもとにした研究が多く存在



1

HTTP通信を用いた感染端末検知

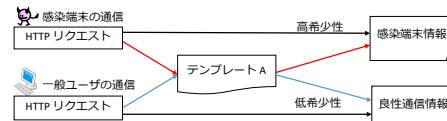
- ・テンプレート
 - ・マルウェアが送出するHTTPリクエストをもとに生成
 > 正規表現等を用いることで攻撃規則をテンプレートに反映
- ・テンプレートを用いた従来手法



2

従来手法の検知手法の課題

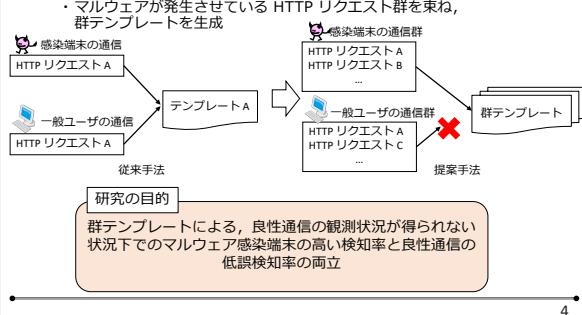
- ・正規化による良性通信の誤検知
 - ・適用先環境の観測情報を用いた誤検知の抑制が実施
 - ・観測情報を事前に収集し、各要素の頻度（希少性）を算出
 - ・通信の各要素の希少性が高い ⇒ 適用先ネットワーク内での出現が稀
 - ・感染端末の可能性大
 - ・テンプレートの類似度が高く、希少性が高いHTTPリクエストを感染端末として検知



3

研究目的

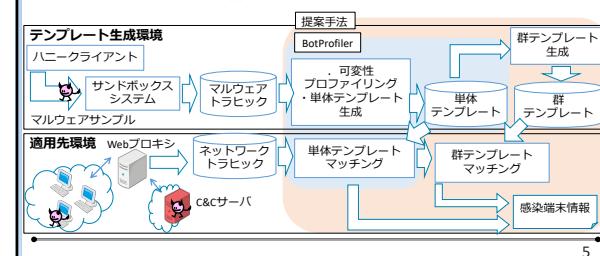
- ・良性通信が十分に得られない状況下において高い検知率と低い誤検知率を両立できる手法を提案
 - ・マルウェアが複数のHTTPリクエストを発生させていることに着目
 - ・マルウェアが発生させているHTTPリクエスト群を束ね、群テンプレートを生成



4

提案手法の概要

- ・テンプレート生成環境
 - ・生成された単体テンプレートを用いて群テンプレートを生成
- ・適用先環境
 - ・適用先で発生するHTTPリクエストを時間分割し、HTTPリクエスト群にマッチング
 - ・各HTTPリクエスト群に対しマッチングを実施

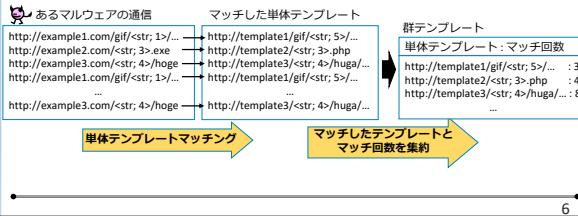


5

群テンプレートの生成

各マルウェアごとに群テンプレートを生成

- HTTP リクエストを各マルウェアごとに集約し、HTTP リクエストと単体テンプレートのマッチングを実施
- マルウェアごとにマッチした単体テンプレートとマッチ回数を集約し、「群テンプレート」として生成
- 複数通信を行う感染端末の挙動を確認



提案する感染端末検知手法

検知手順

- HTTP リクエスト群中の各 HTTP リクエストに対し 単体テンプレートマッチングを実施
 - 単体テンプレートとの類似度 Score(h, t) に応じて以下のように分岐
- ```

graph TD
 A[HTTPリクエスト群] --> B[単体テンプレートマッチング]
 B -- "Score(h, t) ≥ θ_H の HTTP リクエストが存在" --> C[感染端末情報]
 B -- "Score(h, t) < θ_L" --> D[群テンプレートマッチング]
 D -- "全てのHTTPリクエストが Score(h, t) < θ_H" --> E[群テンプレートマッチング]
 E --> F[感染端末情報]
 B -- "Score(h, t) ≤ θ_L" --> G[良性通信情報]

```

### 単体テンプレートマッチング

- HTTP リクエスト群中の HTTP リクエスト h と各テンプレート t の類似度をスコア Score(h, t) により算出
- HTTP リクエストとテンプレートの各要素 (URL パス, URL クエリ, ユーザエージェント) の類似度が高ければ Score(h, t) も高くなるよう設計

7

## 群テンプレートマッチング

### 群テンプレートとの類似度を利用したスコア計算

- HTTP リクエスト群 D と 群テンプレート T の類似度を群スコア S(D, T) により算出
- 群テンプレート内の単体テンプレートのうち、HTTP リクエスト群にてマッチした単体テンプレートの割合が高くなれば S(D, T) も高くなるよう設計

### S(D, T) をもとに感染端末を検知

- $S(D, T) \geq \theta_g$  のとき当該 HTTP リクエスト群を悪性と検知

$$S(D, T) = 1 - \frac{1}{|U_T|} \sum_{t_i \in U_T} s(d_i, t_i)$$

$U_T$ : T において 1 回以上マッチした URL の集合

$$s(d_i, t_i) = \begin{cases} \alpha & d_i = 0 \\ \beta(t_i - d_i)/t_i & 0 < d_i \leq t_i \\ 0 & d_i > t_i \end{cases}$$

$d_i, t_i : D, T$  それぞれにおいて i 番目の単体テンプレートにマッチした HTTP リクエストの個数  
 $\alpha, \beta$ : パラメータ ( $\alpha > \beta$ )

8

## 評価手順

### サンプル収集

- 悪性データ: マルウェアサンプルをサンドボックス上で実行して収集
- 良性データ: 実際に大学内で日常的なネットワーク利用の際に生じるトラヒックを収集

| ラベル | 教師                    |             | テスト                  |             |
|-----|-----------------------|-------------|----------------------|-------------|
|     | 期間                    | HTTP リクエスト数 | 期間                   | HTTP リクエスト数 |
| 悪性  | 2017/8/1 - 2017/12/31 | 656,714     | 2018/1/1 - 2018/3/31 | 442,532     |
| 良性  | -                     | -           | 2018/3/1 - 2018/3/31 | 293,120     |

### HTTP リクエスト群の生成

- ある HTTP リクエストから時間 s 秒以内に送信されている HTTP リクエストを 1 つの HTTP リクエスト群として生成

### 検知指標

- 検知率 (TPR: True Positive Rate)
  - テストデータの全感染端末中の検出された感染端末の割合
- 誤検知率 (FPR: False Positive Rate)
  - 生成された HTTP リクエスト群における誤って検出された HTTP リクエストの割合

9

## 提案手法と従来手法の比較

### 従来手法に比べて検知性能が向上

- $\theta_H$  が大きくなると検知率は低下するが、誤検知率は減少
- $\theta_L$  が大きくなると検知率は低下するが、誤検知率は減少

### $\theta_L, \theta_H$ の制御により検知性能が向上

- $\theta_H$  を高い値に設定し、 $\theta_L < \theta_H$  とすることで低い誤検知率を維持しつつ高い検知率を達成

BotProfiler without RP:  
BotProfiler で希少性を用いない手法  
Proposed: 提案手法

|     |                        | 0.87   | 0.90   | 0.93               |
|-----|------------------------|--------|--------|--------------------|
|     | BotProfiler without RP | 87.42% | 30.48% | 29.28%             |
| TPR | $\theta_L$             | 0.30   | 95.23% | 94.35%             |
|     |                        | 0.78   | 87.87% | 86.12%             |
|     |                        | 0.87   | -      | 84.23%             |
|     | BotProfiler without RP | 3.02%  | 2.79%  | 0.42%              |
| FPR | $\theta_L$             | 0.30   | 4.06%  | 3.90%              |
|     |                        | 0.78   | 3.18%  | 3.02%              |
|     |                        | 0.87   | -      | 2.75% <b>0.88%</b> |

$\theta_G = 0.85$  のときの提案手法と従来手法の比較

10

## まとめ

### 研究背景

- マルウェア感染端末検知は重要な課題
- HTTP 通信に着目した研究が多く実施

### 従来手法の課題と研究目的

- 良性通信が十分に得られない状況下では検知率と誤検知率の両立が困難
- マルウェアが複数の HTTP リクエストを送出することに着目し、群テンプレートを生成することを提案

### 提案手法

- HTTP リクエスト群に対し単体テンプレートマッチングを実施
- 類似度が中程度となった HTTP リクエスト群に対し、群テンプレートマッチングを実施

### 評価結果

- 閾値を制御することで従来手法より検知性能が向上

11