

サイバー攻撃の初期段階と推定される活動で使用されるプログラムの分類手法の提案と評価

芦野 佑樹^{†1} 中村 康弘^{†2} 矢野 由紀子^{†1} 島 成佳^{†1}

概要:サイバー攻撃による被害は深刻化しており、できるだけ早い段階でサイバー攻撃を検知し対処することが望ましいとされる。しかしながら、サイバー攻撃が開始されてから発覚までに1年以上を要するケースがある。サイバー攻撃の発覚に時間を要している原因の一つとして、攻撃者がサイバー攻撃の初期段階において攻撃目標に関する情報を収集し、検知されない攻撃手段を選択していることが考えられる。仮に攻撃者による情報収集活動を捕捉できれば、攻撃者の攻撃目標等の推定を通じてセキュリティ対策の重点箇所を決定する際の情報としての活用が期待できる。しかし、攻撃者による情報収集活動は、一般の通信と同じように見えることから捕捉は難しいことが知られている。そこで、本論文では、情報収集活動を捕捉することを目的として、情報収集活動モデルを通じて、情報収集活動で用いられているプログラムを分類する手法の提案する。併せて提案方式の有効性を評価したので報告する。

キーワード: 偵察活動, 情報収集活動, サイバークルチェーン, 大容量パケット解析

Proposal and Evaluation for a Classification Method of Cyber Attack Programs in the Initial Stage

Yuki Ashino^{†1} Yasuhiro Nakamura^{†2} Yukiko Yano^{†1} Shigeyoshi Shima^{†1}

Abstract: Cyber attacks should be detected early phase of the cyber attack in cyber security countermeasures. However according to reports, cyber attacks of some security incidents could not be detected for over one year. In these cases, attackers probably used cyber attack programs. These program methods cannot be detected by security countermeasures based on reconnaissance of attacker's targets. Target reconnaissance activities of attackers will be useful information to selectively decide enforcement points of cyber security countermeasure for the defense side. We proposes of the reconnaissance activities model and the classification of reconnaissance programs to detect reconnaissance activities. This paper also contains evaluation of the proposal classification method.

Keywords: Reconnaissance, Intelligence Activity, Cyber Kill Chain, Big Data Analysis

1. はじめに

サイバー攻撃による被害は社会問題となっている。サイバー攻撃による被害を最小限に留めるためには、サイバー攻撃をできるだけ早い段階で検知することが望ましいとされる[1]。近年では、他組織で発生したサイバー攻撃に関する知見を共有する取り組みがなされている[2]。

こうした取り組みによって過去に発生した同様のサイバー攻撃を早期に検知できるようになっている一方で、サイバー攻撃が検知されるまでに1年以上の時間を要したケースが存在する[3]。サイバー攻撃の検知に時間を要した理由として考えられることは、サイバー攻撃の目標となった組織が検知できないような攻撃手法が用いられていた可能性がある。サイバー攻撃を行う者(以下、攻撃者)は、サイバー攻撃の初期段階として徹底した情報収集活動をネットワーク経由で実施していると言われている[4]。こうした攻撃者の情報収集活動の対象の一つに攻撃目標の検知能力が含まれている可能性があり、結果として攻撃者は検知されにくい攻撃手法を選択している可能性がある。

仮に、情報収集活動の捕捉が可能となれば、攻撃者の攻撃目標や保有技術の推定を行うことができ、結果としてセキュリティ対策の重点箇所や対処方法の検討ができることが期待できる。

しかしながら、攻撃者による情報収集活動は、具体的な内容が明らかになっていない上、通常の活動に紛れた形で実施されているため捕捉は困難と言われている[5]。

本論文では、情報収集活動を捕捉するため、情報収集活動で用いられているプログラム(以下、情報収集用プログラム)の分類が必要であることを示す。併せて、情報収集用プログラムを分類する手法の提案と評価をしたので報告する。

本論文の構成は、第2章でサイバー攻撃の初期段階である情報収集活動とそれを捕捉する必要性について述べる。第3章では関連研究と本研究の位置付けについて情報収集活動モデルを用いて述べる。第4章では、筆者らのセンサーが捉えたデータを用いて情報収集用プログラムを分類する方法について検討し、第5章で情報収集用プログラムを分類する方法を提案する。第6章で提案した方法の有効性を評価し、第7章で考察を述べ、第8章でまとめる。

^{†1} NEC ナショナルセキュリティ・ソリューション事業部
サイバーセキュリティ・ファクトリー
Cyber Security Factory, National Security Solution Division, NEC Corporation

^{†2} 防衛大学校電子情報学群情報工学科
Department of Computer Science, School of Electrical and Computer Engineering,
National Defense Academy

2. サイバー攻撃の初期段階における情報収集活動

本章では、サイバー攻撃の初期段階における情報収集活動について、攻撃者とサイバー攻撃を防ぐ者(以下、防御者)の両者における位置付けを述べる。併せて、防御者が攻撃者による情報収集活動を捕捉するメリットについて述べる。

2.1 攻撃者における情報収集活動の位置付け

攻撃者によるサイバー攻撃が開始されてから検知されるまでに1年以上の時間を要したケースが存在する[3]。検知されるまでに時間がかかってしまった理由の一つとして、防御側によって検知されにくい攻撃手段を攻撃者が選択していたためと考えられる。

攻撃者がどのようにして、検知されにくい攻撃手段を選択しているのかは明らかになっていない。しかしながら、攻撃者の活動をモデル化したサイバークルチェーン[6]によると、攻撃の初期段階である第1段階目では攻撃目標に関する情報を収集する偵察フェーズが存在するとしている。さらに、JPCER/CCによれば、攻撃者は攻撃に関する情報収集を徹底的に実施していると言われている[5]。

以上のことから、攻撃者は、防御者に検知されないような攻撃手段を選択するために必要な情報を、攻撃の初期段階における情報収集活動によって集めている可能性があると言える。

2.2 防御者における攻撃者の情報収集活動の位置付け

攻撃者による情報収集活動は、一般の活動に紛れた形で行われており、通常の通信と見分けがつかないと言われている[5]。

2.1節で述べたとおり、攻撃者は徹底した情報収集活動をしているにも関わらず、それを捕捉した報告は多くないことから、それだけ情報収集活動の捕捉が困難であると言える。

攻撃者による情報収集活動は以上のような特徴があることから、情報収集活動の捕捉は防御者にとっては困難であり、結果として具体的な対策は行われていなかった。

2.3 情報収集活動を捕捉するメリット

サイバークルチェーンに基づけば第1段階である偵察フェーズからサイバー攻撃は始まっている。したがって、第3段階目である配送フェーズにおけるマルウェアの送付等を侵入検知システム(Intrusion Detection System, IDS)で検知するように、第1段階目の偵察フェーズで行われている情報収集活動を捕捉することは重要であると筆者らは考える。

仮に攻撃者による情報収集活動を捕捉することができた際のメリットについて以下に述べる。

2.3.1 攻撃目標の推定を通じた重点対策箇所の決定

セキュリティ対策へ投資できるリソースには限界がある。しかし、攻撃者が情報収集している箇所を特定できれば、その箇所は将来別のサイバー攻撃を受ける可能性が高いと言える。したがって、このように攻撃者の標的が情報収集活動の捕捉を通じて事前に把握することができるのならば、セキュリティ対策のリソースを投じるべき対象の決定に役立てることが期待できる。

2.3.2 攻撃者が保有する攻撃能力の推定と対策方法の決定

筆者らは、独特の通信パターンを持つ発信源と通信することで、インターネット全体のウェブサーバを調査していると推定

される活動を捕捉したことがある[7]。この活動で使用されたシステムは、大規模な分散システムで有り、高度な技術を有する個人でも企業でも研究機関でもない組織が開発と運用をしているのではないかと推定できた。

このような攻撃者の保有する攻撃能力の推定は、適用すべき対策方法を決定するための情報として活用できることが期待できる。

3. 関連研究と本研究の位置付け

攻撃者がネットワークを介した通信を伴う情報収集活動であれば、攻撃者によって実装されたプログラムを用いた通信は必ず存在と言える。すなわち、情報収集活動には、通信の送り手と受け手が存在する。通信の送り手と受け手の関係は、通信モデルで表現することができる。そこで、筆者らは、通信を伴う情報収集活動を通信モデルの拡張によって説明することにより、情報収集活動の補足に関する議論が進むのではないかと考えた。

本章では、攻撃者による情報収集活動において、情報の送り手と受け手の関係をモデル化し関連研究を述べる。併せて本研究の位置付けについて述べる。

3.1 情報収集活動モデル

攻撃者の情報収集活動において、情報を収集したい対象(以下、情報収集対象)と1回よりも多い回数の通信が発生する場合、攻撃者と情報収集対象は通信していると言える。つまり、情報収集活動は通信モデルで説明できる可能性がある。

本節では、通信モデルを情報収集活動に拡張した情報収集活動モデルについて述べる。

3.1.1 シュラムモデル

通信のモデルは、情報の送り手と受け手の関係を単純化したものである。代表的なものにシャノンらによって提唱されたシャノンモデル[8]の他、情報の送り手と受け手の役割を切り替えることにより複数回の通信を行うことを特徴としたシュラムモデルがある[9]。

攻撃者による情報収集活動において攻撃者と情報収集対象との間で複数回の通信が行われるならば、情報収集活動はシュラムモデルに基づいて説明できるのではないかと考えた。3.1.2ではシュラムモデルを拡張したモデルで情報収集活動を説明する。

3.1.2 シュラムモデルを拡張した情報収集活動のモデル

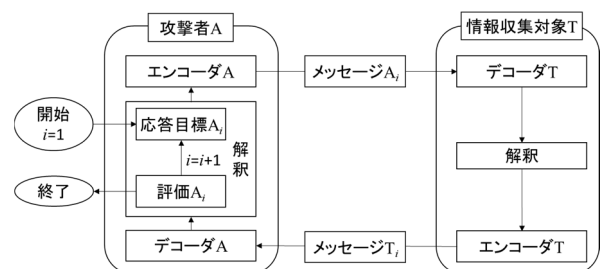


図1 シュラムモデルを情報収集活動に拡張した概念図

ここでは、攻撃者が情報収集対象にデータを送ることによって情報収集を行うシーンにおいて、シュラムモデルを拡張した情報収集活動モデルについて述べる。

情報収集活動モデルの特徴は、送り手である攻撃者が、受け

手である情報収集活動の応答するメッセージに意図的に欲しい情報を含ませるためのメッセージを送る点にある。

情報収集活動モデルを図 1 に示し、以下に概要を説明する。

(1) 攻撃者から情報収集対象に向けた通信

このモデルにおいて、通信は攻撃者 A から始められる。攻撃者 A は情報収集対象 T から情報を得るために、応答して欲しい内容を設定する。応答して欲しい内容とは、攻撃者 A にとって欲しい情報と同義であり、本モデルにおいては応答目標 T_1 である。エンコーダ A は、応答目標 A_1 をメッセージ A_1 に変換して情報収集対象 T に送る。

(2) 情報収集対象から攻撃者への応答

デコーダ T は、メッセージ A_1 を受け取り内容の解釈をする。情報収集対象 T の解釈結果はエンコーダ T に渡されメッセージ T_1 に変換する。メッセージ T_1 は攻撃者 A に送られる。

(3) 攻撃者による 2 回目以降の送信の流れ

デコーダ A は、受け取ったメッセージ T_1 を受け取り評価 A_1 に渡す。評価 A_1 では、応答目標 A_1 で設定した情報が含まれていることを評価する。評価 A_1 は、メッセージ T_1 に応答目標 A_1 が含まれていないことから情報収集が失敗したと判断した場合、もしくは、それ以上の情報収集は不要であると判断した場合はフローを終了する。

攻撃者がさらに情報収集対象 T から情報を取得したい場合は、 i の値を一つ繰り上げ、新たに応答目標 A_2 を設ける。以降のフローは(1)以降と同様である。

3.2 情報収集活動モデルに基づいた関連研究

情報収集活動を含むサイバー攻撃を捕捉することを目的とした関連研究は多く取り組まれている。

サイバー攻撃の捕捉は、サイバー攻撃に関する通信の記録や応答する装置(以下、センサー)によって行われる。センサーの種類は、攻撃者とセンサーの間で交わされる通信の方向や通信回数によって 3 種類に分類できる。

本節では、3.1 節で述べた情報収集活動モデルに基づいて 3 種類の各センサーについて述べる。

3.2.1 パッシブセンサー

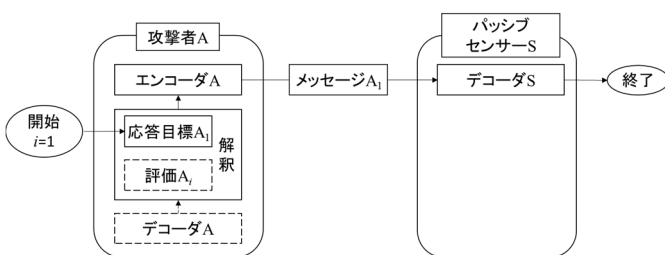


図2 パッシブセンサーのモデル

パッシブセンサーとは、通信の記録に徹しており、センサー自らは通信を発しないセンサーである。代表的なパッシブセンサーを用いた関連研究として、NICTer プロジェクトのダークネット観測がある[10]。このパッシブセンサーの情報収集活動モデルを表現した図 2 に基づいて説明する。

情報収集活動モデルにおける情報収集対象 T がパッシブセンサー S となる。パッシブセンサー S は何も応答しないことから、パッシブセンサー S は攻撃者 A の応答目標 A_1 に関係なくデコーダ S の受信にてフローが終了となる。

3.2.2 リアクティブセンサー

リアクティブセンサーとは、攻撃者からできるだけ多くの情報を得るために、攻撃者のメッセージに対して応答するセンサーである。代表的なリアクティブセンサーとして、ハニーポット[11]がある。ハニーポットとは、攻撃者によって発信された OS やアプリケーションサーバの脆弱性を突く通信を受け入れ、マルウェア等を捕獲するものである。

このリアクティブセンサーの情報収集活動モデルを図 3 に基づき説明する。

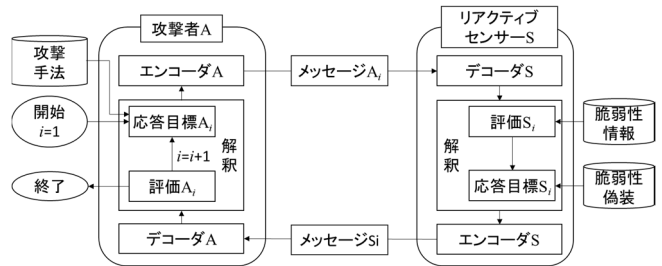


図3 リアクティブセンサーの概念図

攻撃者 A の応答目標 A_1 は、脆弱性がある場合に返してくる応答内容に相当する。エンコーダ A は、この応答目標 A_1 をメッセージ A_1 に変換しリアクティブセンサー S に送る。

リアクティブセンサー S の評価 S_1 は、受信したメッセージ A_1 から攻撃者 A の応答目標 A_1 を推定する。評価 S_1 が推定した応答目標 A_1 から応答するべきであると判断した場合は、攻撃者 A が次のメッセージ A_2 を送ってもらえるような応答目標 S_1 を設定しエンコーダ S を介してメッセージ S_1 を変換する。つまり、リアクティブセンサーに脆弱性があるように偽装したメッセージ S_1 を攻撃者 A に送り返す。

攻撃者 A は、評価 A_1 にてメッセージ T_1 の結果からリアクティブセンサー S の脆弱性の有無を評価する。攻撃手法の実行が可能であると評価した場合は、 i を一つ繰り上げ、応答目標 A_2 を新たに設ける。応答目標 A_2 は、攻撃手法に基づいている。

応答目標 A_2 に基づいた、メッセージ A_2 がリアクティブセンサー S に送られる。リアクティブセンサー S は、この攻撃手法を含むメッセージ A_2 を記録する。

このようにリアクティブセンサーは、攻撃者の応答目標に沿ったメッセージを送り返すことで、攻撃者からの多くの通信を取得する。

3.2.3 アクティブセンサー

アクティブセンサーとは、攻撃者に対して自発的に通信を發して攻撃者の応答結果から情報を得るセンサーである。

アクティブセンサーの取り組みとして、田辺らはマルウェア感染ホストに対して再侵入を行うことで感染拡大を阻止する手法を提案している[12]。

このアクティブセンサーを情報収集モデルで表現した図 4 に基づいて説明する。3.1.2 の図 1 で述べた情報収集活動モデルにおける攻撃者 A がアクティブセンサー S になり、情報収集対象 T が攻撃者 A に入れ替わった形になる。

攻撃者が特定できている状況下であれば有効なセンサーであると言える。

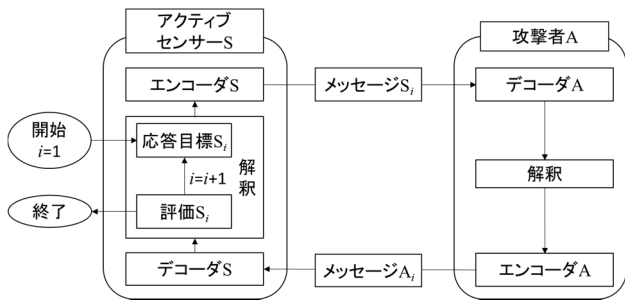


図4 アクティブセンサーの概念図

3.3 本研究の位置付け

本節では、情報収集活動モデルに基づいて、情報収集活動の捕捉に向け各センサーの有効性と課題について述べ、本研究の位置付けを述べる。

3.3.1 情報収集活動を捕捉する上での既存センサーの課題

パッシブセンサーは、攻撃者から送られる通信の受信のみに徹していることから、攻撃者からの2回目以降の通信を受信することを想定していない。2.2節で述べたとおり、情報収集活動は通常の通信と見分けがつかないとされる。そのため、攻撃者からの通信が1回のみでは、情報収集活動を捕捉するには十分であるとは言えない。以上のことから情報収集活動の捕捉としてパッシブセンサーの利用は難しいと言える。

アクティブセンサーは、攻撃者が特定できている時に有効的なセンサーであり、2.2節で述べたように特定できない攻撃者を調査する目的で使用することを想定していない。したがって、情報収集活動の捕捉にアクティブセンサーを利用することは難しいと言える。

リアクティブセンサーは、情報収集活動モデルにおける攻撃者の応答目標が判明している状況下であれば、攻撃者からの2回目以降のデータを取得することができる。したがって、2.2節の特徴を有する情報収集活動の捕捉にはリアクティブセンサーが適していると言える。

しかしながら、従来のリアクティブセンサーは、脆弱性を持っているようなサーバに偽装することによって攻撃者からの情報を多く取得しようと試みるものである。したがって、従来のリアクティブセンサーは、2.2節で述べたような、情報収集活動は通常と見分けがつかない通信に対応することを想定していない。そのため、既存のリアクティブセンサーでは攻撃者による情報収集活動の捕捉は難しい。

3.3.2 本研究の領域

筆者らの先行研究では、攻撃の初期段階と推定される情報収集活動に関する通信は送信元の環境と実装に依存している可能性が高いこと示している[13]。この結果と情報収集活動モデルに基づけば、メッセージ A_i は攻撃者が情報収集活動で使用するプログラム（以下、情報収集用プログラム）により送信されているのではないかと考えた。

したがって、この情報収集用プログラムに実装されている応答目標 A_i を推定し、適切な応答を返すリアクティブセンサーであれば、通常の通信と見分けがつかなかった情報収集活動を捕捉できるのではないかと考えた。

以上のことから、本論文では、通常の通信に紛れて行われていると考えられる情報収集活動で用いられる情報収集用プロ

ラムを分類する方式の提案と評価を行う。併せて、情報収集用プログラムの分類の結果の扱いについても検討する。具体的には、情報収集用プログラムが分類できなかった際の扱いについても述べる。

4. 予備調査

情報収集用プログラムを分類できる可能性を調査するために、筆者らは TCP レイヤにおけるリアクティブセンサーを用いて通信の観測及び分析を行うことにした。

本章では、TCP/IP レベルのリアクティブセンサーと観測データについて述べ、情報収集用プログラムの分類に活用できる指標について述べる。併せて、予備実験を通じて情報収集用プログラムを分類できる可能性について述べる。

4.1 TCP リアクティブセンサー

ドメインへの割り当てやサーバの運用をしていない IP アドレスであっても、その IP アドレスに対する通信は非常に多い。このような通信の意図を探るために、筆者らは、このような1,500個のIPアドレスに対してTCPのSYNパケットに対してSYN+ACKを応答するリアクティブセンサー(以下、TCPリアクティブセンサー)を設置して通信の観測をしている。

送信元がこのTCPリアクティブセンサー宛にSYNパケットを送るとSYN+ACKが返ってくることから、あたかもTCPコネクションが確立したように見える。そのため、送信元の一部は、ACK及びTCPコネクション確立後の最初のデータを送信してくることがある。筆者は、このデータを初期ペイロードと称して研究を進めている[14][15]。

4.2 データセット2016の諸元

表1 データセット2016の諸元

取得期間	2016/01/01～2016/12/31
観測に使用したIPアドレス数	1,500
データ形式	libpcap (ファイルは1日単位)
容量	合計約1.2TB
パケット数	約25.2億

表2 プロトコル毎のパケット数

プロトコル名	パケット数 (全体比率)
TCP/IP	約23.5億 (93.2%)
UDP/IP	約1.4億 (5.6%)
それ以外	約0.3億 (1.2%)

表3 宛先TCPポート毎のパケット数

宛先TCPポート	パケット数 (全体比率)
22	約6.4億 (27.2%)
80	約2.2億 (9.4%)
3389	約1.2億 (5.1%)

表4 宛先TCPポート毎の発信元IPアドレスの数

宛先TCPポート	アドレス数 (全体比率)
6881	約21,000 (4.5%)
80	約15,000 (3.5%)
443	約3,900 (0.9%)

4.1節で述べたTCP/IPリアクティブセンサーで観測したデータセットの諸元を表1に示す。観測期間が2016年であること

から、データセットの名称をデータセット 2016 とする。

データセットに含まれていたパケットデータの集計を示す。表 2 ではプロトコル毎のパケット数を示す。TCP/IP が全体の全体の 93.2% を占めていた。

大多数の通信が TCP/IP であったことから、宛先 TCP ポートのパケット数の集計を上位 3 位まで表 3 に示す。データセットに含まれる発信元の IP アドレスの数は、461,534 個あった。そこで、各 TCP ポートの発信元 IP アドレスの数を集計した。発信元アドレスの数が多順にして上位 3 位を表 4 に示す。

4.3 情報収集用プログラムを分類するための指標の検討

本節では、攻撃者から送られたメッセージから送信者の応答目標をするために指標を検討する。

4.3.1 IP アドレス

攻撃者は、サイバー攻撃にボットネットと呼ばれるネットワークを使用することがある[16]。ボットネットの中には長期間に渡り同一の IP アドレスを使用するケースがあるとされる。したがって、IP アドレスには攻撃者の特性が顕れる可能性があることから、情報収集用プログラムの分類に役立てられると考えた。

データセットに含まれていた全パケットを分析し、IP アドレス毎の発信回数を集計した。その結果、総じて送信回数の多い IP アドレスは少ない傾向にあった。

IP アドレスで情報収集用プログラムを分類するためには、2 回以上送信する IP アドレスである必要がある。データセットを統計した結果、1 回しか送信しない IP アドレスは全体の 31.8% であった。したがって、IP アドレスに基づいた分類では、31.8% は分類できないことになる。以上のことから、IP アドレスに基づいた情報収集用プログラムの分類は困難であると判断した。

4.3.2 ペイロードの内容

情報収集活動モデルに基づけば、同一の情報収集用プログラムが異なる IP アドレスで使用されていれば、ペイロードも同一になるはずである。そこで、ペイロードの内容によって情報収集用プログラムの分類が可能になるのではないかと考えた。各通信に含まれるペイロードの定義は通信プロトコルによって異なるため、表 5 で定義する部分をペイロードとして扱う。

表 5 プロトコル毎のペイロードの定義

プロトコル	ペイロードの定義
TCP/IP	PUSH フラグの付くパケットで TCP ヘッダ以降に続くデータ
UDP/IP	UDP ヘッダ以降に続くデータ
それ以外	IP ヘッダ以降に続くデータ

表 5 の定義に従いペイロードの種類を集計したところ、データセット 2016 中に存在するペイロードの種類は約 1 億 6400 万種類存在した。この内、同一のペイロードが 2 回以上出現したペイロードは 1.7%(約 300 万種)のみであり、残りの 98.3%(約 1 億 6100 万種類)は 1 回しか出現しなかった。そのため、ペイロードだけで分類すると通信の 98.3% は分類できないことになる。そのため、ペイロードによる応答目標の分類は不適切であると判断した。

4.3.3 抽出サイズと部分ペイロード

ペイロードはプロトコルによってデータの配置が異なること

から、プロトコルの特性に応じてペイロードを扱うべきだと考えた。

本論文で扱うペイロードの対象として、80/TCP のペイロードを選択した。80/TCP のペイロードを選択した理由は、4.2 節で述べた諸元に基づく、データセットに含まれるプロトコルは TCP が一番多く、表 3 と表 4 の両方で上位に存在したからである。80/TCP は HTTP で用いられる。HTTP は TCP コネクション確立後にクライアントからサーバに対して要求メッセージを送るプロトコルであり、データセットにはこの要求メッセージが記録されている。

筆者らは、HTTP の規約上、ウェブクライアントからウェブサーバに対する要求は GET や POST といった限られた数の種類しか存在しない特徴に着目した。具体的には、例えば、ペイロードの先頭を 1 バイトだけを取り出した場合、G や P といった限られた種類の文字列が取り出されることになる。この特徴に基づく、ペイロードの先頭から取り出すバイト数(以下、部分ペイロード)と取り出された文字列(部分ペイロード)の種類数が比例しないはずである。この特徴に基づけば、ペイロードの先頭部分の共通部分を上手く分離できるのではないかと考えた。

そこで、抽出サイズと部分ペイロードの種類数の関係を図 5 に示した。図 5 は、横軸に抽出サイズとし縦軸に部分ペイロードの種類数としてある。なお、ここでの部分ペイロードとは可読なアスキー文字のみを扱う。以降、このデータを HTTP データと称する。データセット 2016 において、HTTP データの種類数は約 70.3 万であった。

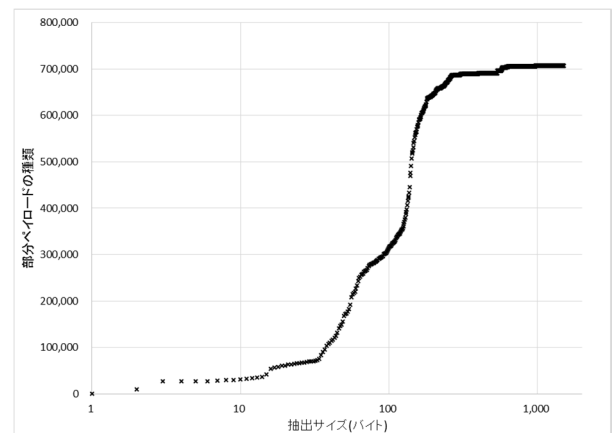


図 5 抽出サイズと部分ペイロードの種類数の傾向

図 5 のとおり、抽出サイズと部分ペイロードが比例しないのは、HTTP データは、先頭部分からある一定の長さまでは同じペイロードであったことを意味する。

この結果から、抽出サイズを変化させながら部分ペイロードに基づいて共通部分を取り出せば、少なくとも使用しているプログラムの分類に役立てられるのではないかと考えた。

4.3.4 ペイロードのサイズ

情報収集活動モデルに基づけば、同一の情報収集用のプログラムにおける初回の応答目標 1 は同一であり、同一のエンコーダ A に渡れば出力されるメッセージ A₁ も同様と言える。そして、同一の A₁ であればペイロードサイズも同一になるのではないかと考えた。そこで、データセット 2016 での HTTP データに

おけるペイロードサイズ毎の packets 数を調査した。図 6 は、横軸にペイロードサイズとし、縦軸にペイロードサイズ毎の packets 数とした。なお、今回はペイロードサイズを IP におけるペイロードサイズとした。

図 6 に示すとおり、ペイロードサイズ毎に packets 数に偏りがあることが分かった。したがって、HTTP データのサイズはランダムで決定されているものではなく、情報収集用プログラムの実装に依存している可能性が高いと言える。

以上のことから、HTTP データのサイズを応答目標の分類に用いる指標とすることにした。

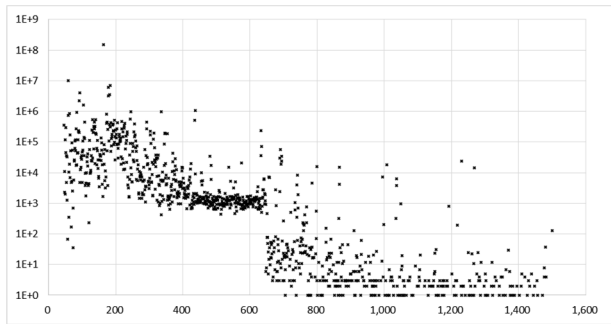


図 6 ペイロードサイズ毎の送信回数

4.3.5 情報収集用プログラム分類指標

4.4 節では、抽出サイズ・部分ペイロードと、ペイロードサイズに基づいて情報収集用プログラムを分類できる可能性について予備実験を実施したので報告する。

4.4 予備実験

抽出サイズ・部分ペイロードと、ペイロードサイズに基づいて、情報収集用プログラムの分類が可能であるのかを調査するために、特定のペイロードサイズだけを取り出し、抽出サイズ毎の部分ペイロードの種類数の変移を調査する。

ここでは、顕著な傾向が現れていたペイロードサイズ 304 バイトについて述べる。ペイロードサイズが 304 バイトとなるペイロードは 866 種類あった。抽出サイズを 1 バイトから 304 バイトに増やしていった際の部分ペイロードの種類数を図 7 に示す。

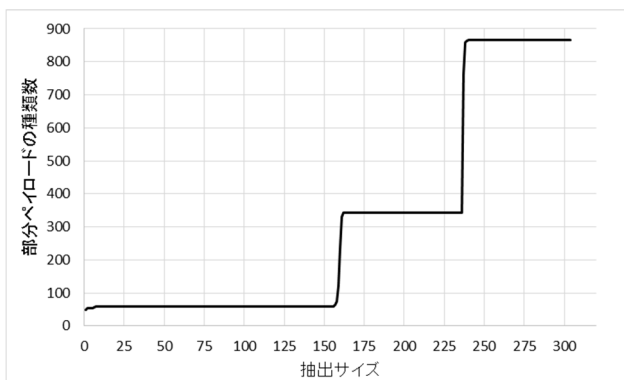


図 7 抽出サイズ毎の部分ペイロードの種類

図 7 は縦軸に部分ペイロードの種類数、横軸に抽出サイズとした。図 7 のとおり 3 段の階段状になった。

以下に、部分ペイロードの種類数が極端に変化する直前である抽出サイズ 156 の部分ペイロードを示す(表 6)。

表 6 の部分ペイロードの末尾は、Content-Length フィールド

の値を示す直前であった。抽出サイズが 156 よりも大きくすると急激に増加していた理由は、Content-Length 以降の数字が 1004 から 8754 までの複数の数字を持つ HTTP データが多く存在したことに起因する。したがって、抽出サイズによる部分ペイロードの種類数の変化は、HTTP データの共通部分を取り出しは可能である。この共通している部分ペイロードは、同一の応答目標とエンコーダの関係であると言え、つまりは同一の情報収集用プログラムであると考えられる。

以上の予備実験の結果から、抽出サイズ・部分ペイロードと、ペイロードサイズに基づいて情報収集用プログラムの分類は可能であると判断した。

表 6 抽出されたペイロード

POST /*****_action.php HTTP/1.1 (省略) Content-Length: □	※一部*文字で伏せている ※□は半角スペース
--	---------------------------

5. 提案方式

本章では、TCP リアクティブセンサーで受信したペイロードから情報収集用プログラムを分類するための方式を提案する。

5.1 構造

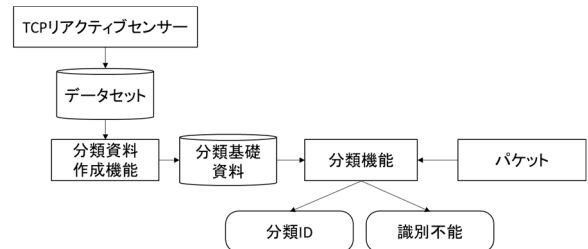


図 8 提案方式の構成

提案方式のシステム構成は、図 8 に示すような構造を持つ。

TCP リアクティブセンサーによって得たデータセットを用いて分析基礎資料を作成する。分類基礎資料には、ペイロードのサイズ毎の部分ペイロードが記録されている。

次に、パケットから情報収集用プログラムを分類するための分類機能について述べる。分類機能は、パケットを受け取るとペイロードのサイズを元に分類基礎資料からマッチする最長の部分ペイロードの ID を返す。この ID は、情報収集用プログラムの分類を意味する。マッチングする部分ペイロードがない場合は、分類不能と返す。

5.2 分類基礎資料の作成法

4.4 節の予備実験から、表 7 に示すフローでペイロードを分類するための基本データを作成する(表 7)。

データセット 2016 を用いて表 7 のフローで分類基礎資料を作成した。分類基礎資料の一部を図 9 に示す。図 9 は左側から、ID、ペイロードサイズ、抽出サイズ、部分ペイロードである。

5.3 分類機能

分類機能は、与えられたペイロードに対して、分類基礎資料中の ID を返す機能である。この ID は、情報収集用プログラムを分類する際の元となる番号となる。表 8 に示すフローで与えられたペイロードに対応する ID を返す。ただし、見つからない

場合は分類不能を返す。

表7 分類基礎資料作成フロー

Step 1. データセットからペイロードを抽出
分類基礎資料を作成する対象となるペイロードをデータセットから抽出する。
Step 2. ペイロードサイズ毎に分類
Step1で抽出されたペイロードをペイロードサイズ毎に分類する。
Step 3. 抽出サイズ毎の部分ペイロードの種類を取得
ペイロードサイズ毎に、抽出サイズを1からペイロードサイズまで変化した時の部分ペイロードの種類数を取得する。
Step 4. 部分ペイロード数が変化する直前の抽出サイズを取得
ある抽出サイズ x の部分ペイロードの種類数 y と、抽出サイズ $x+1$ における部分ペイロードの種類数 $y+1$ が異なる場合、抽出サイズ x の全ての部分ペイロードを分類基礎資料に登録する。

ID	ペイロードサイズ	抽出サイズ	部分ペイロード	ID	ペイロードサイズ	抽出サイズ	部分ペイロード
2035776	304	11	POST /icons	2035806	304	13	GET /status H
2035777	304	11	POST /index	2035807	304	13	GET /statusff
2035778	304	11	POST /login	2035808	304	13	GET /todayte
2035779	304	11	POST /nepen	2035809	304	13	GET /ubuntu/d
2035780	304	11	POST /test.	2035811	304	13	POST /admin.c
2035783	304	11	sts/precise	2035812	304	13	POST /cgi-bin
2035783	304	13	GET / HTTP/1.	2035813	304	13	POST /cgi-ovs
2035794	304	13	GET /admin.cg	2035814	304	13	POST /error H
2035795	304	13	GET /appserv/	2035815	304	13	POST /icons H
2035796	304	13	GET /cgi-bin/	2035816	304	13	POST /index.p
2035797	304	13	GET /cgi-ovs/	2035817	304	13	POST /login.p
2035798	304	13	GET /images H	2035818	304	13	POST /nepent
2035799	304	13	GET /include	2035819	304	13	POST /test.cg
2035800	304	13	GET /index.ph	2035822	304	13	sts/precise-b
2035801	304	13	GET /kali-sec	2035832	304	14	GET / HTTP/1.1
2035802	304	13	GET /login.ph	2035833	304	14	GET /admin.cgi
2035803	304	13	GET /annual H	2035834	304	14	GET /appserv/n
2035804	304	13	GET /plugins/	2035835	304	14	GET /cgi-bin/e
2035805	304	13	GET /sample H	2035836	304	14	GET /cgi-bin/h

図9 分類基礎資料の一部

表8 情報収集用プログラム分類フロー

Step1. パケットからペイロードを取り出す
分類機能に入力されたパケットからペイロードを取り出す。
Step2. ペイロードサイズ S を取得する
Step1で取り出したペイロードのサイズ S を取得する
Step3. ペイロードサイズ S の存在を確認する
分類基礎資料中に、ペイロードサイズ S と同サイズのデータが存在するか確認する。存在しない場合は、分類不能を返す。
Step4. マッチングする最長の部分ペイロードを取得する
分類基礎資料中のペイロードサイズ S と同サイズの部分ペイロードのうち、マッチした最長の部分ペイロードが存在すればその ID を返す。もし、どの部分ペイロードにもマッチしない場合は、分類不能を返す。

6. 評価実験

本章では、第5章で提案した情報収集用プログラムを分類する手法の有効性を評価する実験を行う。

6.1 データセット2017の諸元

表9 データセット2017の諸元

取得期間	2017/01/01~2017/01/31
データ形式	libpcap (ファイルは1日単位)
容量	合計約 380GB
パケット数	約 20.9 億

評価対象となるデータセットの諸元を表9と表10に示す。本データセットをデータセット2017と称する。データセット2017の80/TCPのペイロードに関する諸元を表9と表10に示

す。データセット2016に基づいて作った分類基準資料に基づいて、データセット2017のHTTPデータを分類する。

評価は、データセット2017のHTTPデータから情報収集用プログラムの分類と分類不能を出力できることとする。

表10 80/TCPに関する諸元

パケット数	2440 万 (1.2%)
ペイロードの種類数	19.0 万種
HTTP データの種類数	17.1 万種
送信元 IP アドレス	約 1110 万 (0.5%)

6.2 結果

分類機能でデータセット2017を分類した結果、HTTPデータ17.0万種の内、93.5%である15.9万種の分類ができた。この結果から、2016年で使われた情報収集用プログラムの93.5%は、2017年1月でも使われたと考えられる。一方で、分類機能は6.5%である3.1万種を分類不能と判定した。

本結果に関する考察は第7章で述べる。

7. 考察

7.1 情報収集用プログラムの分類結果と課題

本論文で提案した方式において、情報収集用プログラムを分類できない割合は6.5%であった。したがって、IPアドレスやペイロードの内容に基づいて分類した際における、分類不能の範囲が20%から6.6%に縮小できた。

以上の結果から、提案方式は、IPアドレスやペイロードの内容に基づいて分類するよりも精度が高いと言える。

しかしながら、今回の評価はあくまで、80/TCPにおける送信元からの1回目のHTTPデータから情報収集用プログラムを分類したに過ぎない。そのため、今回のHTTPデータが攻撃者による情報収集活動であるとの判断はできない。

情報収集活動の捕捉のためには、多くの情報が必要であり、送信元から2回目以降のペイロードを受信する必要がある。2回目以降のペイロードを受けるためのリアクティブセンサーについては、7.3節で述べる。

7.2 新種の情報収集用プログラムの発見と課題

本節では、6.3節の実験で分類不能となったHTTPデータの例を表11に示した上で、分類不能と判断された理由と、分類不能となったHTTPデータの扱いについて述べる。

表11 分類不能となったペイロードの一部

項番	HTTP データ (一部伏せ字)
1	HEAD / HTTP/1.1 Host: ***.***.***.*** Connection: keep-alive (以下省略)
2	POST /j***.php HTTP/1.1 (以下省略)
3	hello

項番1と項番2は、当該ペイロードサイズがデータセット2016において部分ペイロードが分類基準資料に存在していなかったため、分類不能となったものである。この結果は、項番1と項番2のHTTPデータは、データセット2017で初めて捕捉

された HTTP データであることから、2017 年 1 月に初めて使用された情報収集用プログラムであると考えられる。

このように、分類不能という結果は、新しい情報収集用プログラムの使用が確認できることを意味している。また、これらの HTTP データを 5.2 節で述べた分類基準資料作成フローで、分類基準資料に登録することで分類可能になるとも考えられる。

一方、項番 3 は、hello の 5 バイトだけを送ってきており、項番 1 や 2 と異なり、HTTP プロトコルには則っていないことから、現在の所、応答目標は不明である。筆者らは、このように HTTP に則っていない通信は、特別な情報収集用プログラムによって送信されている可能性が高く、注目するべきであると考えられる。

7.3 情報収集活動の捕捉するリアクティブセンサーに向けて

以上の議論から、情報収集活動を捕捉するためのリアクティブセンサーは、図 10 に示すような構造を持つ必要があると考えた。

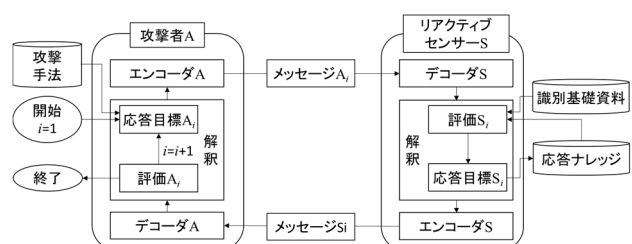


図 10 情報収集活動の捕捉するリアクティブセンサーの構造

3.2.2 で述べたリアクティブセンサー S と異なる点は、分類基礎資料、応答ナレッジがある点である。分類応答資料は、メッセージ A_i から情報収集用プログラムを分類する際に用いられる。応答ナレッジは、分類された情報収集用プログラムに対応した応答目標の方針、及び、過去に設定した応答目標が記録されている。

ハニーポットとの違いは、情報収集用プログラムの分類と過去の応答履歴に基づいてメッセージ S_i の内容を変えていく点である。このようなリアクティブセンサーであれば、通常の通信に紛れた情報収集活動を捕捉することができるのではないかと考えられる。このリアクティブセンサーを実現するに当たり、次のような課題がある。

一つ目は、応答ナレッジに入れる応答目標 S_i の方針を定める知見がない点である。例えば、GET メソッドに対しては、単純にページの存在を意味するステータスコード 200 を無条件で返すのか、特定のサイトの構造を模して必要に応じてページが存在しないことを意味するステータスコード 404 を返すべきなのかについては、議論の余地がある。

二つ目は、攻撃者の応答目標 A_1 に合致したメッセージ S_2 応答を返してきたとして、攻撃者が新たに設定した応答目標 A_2 に対応したメッセージ A_2 を識別する方法が存在しない点である。現行の提案方式のままでは、異なる情報収集用プログラムがメッセージを送ったとしか分類できない。

今後は、上記の課題を解決する手法を検討していきたい。

8. おわりに

本論文では、サイバー攻撃の初期段階である情報収集活動を

捕捉することの重要性を述べると共に、情報収集活動の捕捉には情報収集活動で用いられているプログラムを分類する必要性を述べた。

TCP の SYN パケットに対して SYN+ACK を返す TCP リアクティブセンサーによって得たデータセットを分析した結果、IP パケットにおけるペイロードサイズとペイロードの先頭から任意のサイズを取り出した部分ペイロードが統計的に偏ることを発見した。本論文では、この発見に基づいて情報収集用プログラムの分類方法を提案し、データセットを用いた実験を通じて有効性を示した。

今後は、この方式に基づいて攻撃者に応答を返すことで、情報収集活動の捕捉をするシステムの構築に向けて検討していきたい。

参考文献

- [1] “【注意喚起】攻撃の早期検知と的確な初動による深刻な被害からの回避を”。<https://www.ipa.go.jp/security/ciadr/vul/20160623-ta.html>, (参照 2017-08-25).
- [2] 伊藤大貴, 野村健太, 神菌雅紀, 白石善明, 高野泰洋, 毛利公美, 星澤裕二, 森井昌克. 脅威情報を関連付けるための攻撃活動の表現. 信学技報 116(522), 2017, p. 147-152
- [3] “M-Trends 2016”, <https://www2.fireeye.com/WEB-M-Trends-2016-JA.html>, (参照 2017-08-25).
- [4] “ネットワークビギナーのための情報セキュリティハンドブック”. <https://www.nisc.go.jp/security-site/files/handbook-all.pdf>, (参照 2017-08-25).
- [5] “高度サイバー攻撃への対処におけるログの活用と分析方法 1.0 版”. https://www.jpccert.or.jp/research/APT-loganalysis_Report_20151117.pdf, (参照 2017-08-25).
- [6] “Seven Ways to Apply the Cyber Kill Chain with a Threat Intelligence Platform”. http://lockheedmartin.com/content/dam/lockheed/data/corporate/documents/Seven_Ways_to_Apply_the_Cyber_Kill_Chain_within_a_Threat_Intelligence_Platform.PDF, (参照 2017-08-25).
- [7] 芦野佑樹, 島成佳, “インターネットノイズに対する偽装応答機能の実装と観測に基づいた意図が不明なリクエストに関する考察”, SCIS2015
- [8] Claude E. Shannon, Warren Weaver, “Mathematical Theory of Communication”, Univ of Illinois Pr, 1963.
- [9] Wilbur Lang Schramm, “The Process and Effects of Mass Communication”, Univ of Illinois Pr, Revised edition, 1971.
- [10] 中尾康二, 井上大介, 衛藤将史, 吉岡克成, 大高一弘, “ネットワーク観測とマルウェア解析の融合に向けて-インシデント分析センターnicterの研究開発-”, 情報処理学会論文誌 第 50 巻第 3 号, 2009, p.235-242.
- [11] “Dionaea”, <https://github.com/rep/dionaea>, (参照 2017-08-25).
- [12] 田辺瑠偉, 鈴木将吾, インミンパバ, 吉岡克成, 松本勉, “マルウェア感染ホストへのリモート再侵入により感染拡大を阻止する手法”, 情報処理学会論文誌 第 57 巻第 9 号, 2016, p. 2021-2033.
- [13] 芦野佑樹, 山根匡人, 矢野由紀子, 島成佳, “長期間に渡るインターネットの観測に基づいたサイバー攻撃の初期活動と推定される通信の発信源を分類する手法の提案”, CSEC, 2017.
- [14] ゴ・キムコン, 中村康弘, “走査活動観測に基づくネットワーク攻撃委との推定”, CSS2016, 2016, p.10331039.
- [15] 中村康弘, “初期ペイロードに着目したネットワーク走査活動の分析”, 情報処理学会第 79 回全国大会, 2016, 5D-92.
- [16] “Know your Enemy: Tracking Botnets”, <http://www.honeynet.org/papers/bots/>, (参照 2017-08-25).