

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2008-176753

(P2008-176753A)

(43) 公開日 平成20年7月31日(2008.7.31)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G06F 21/22 (2006.01)</b>	G06F 9/06 660N	5B276
<b>G06F 21/20 (2006.01)</b>	G06F 15/00 330A	5B285
<b>G06F 13/00 (2006.01)</b>	G06F 13/00 540E	5K030
<b>H04L 12/56 (2006.01)</b>	H04L 12/56 400Z	

審査請求 未請求 請求項の数 26 O L (全 24 頁)

(21) 出願番号 特願2007-12071 (P2007-12071)  
 (22) 出願日 平成19年1月22日 (2007. 1. 22)

(71) 出願人 301022471  
 独立行政法人情報通信研究機構  
 東京都小金井市貫井北町4-2-1  
 (74) 代理人 100130111  
 弁理士 新保 斉  
 (72) 発明者 吉岡 克成  
 東京都小金井市貫井北町4-2-1 独立  
 行政法人情報通信研究機構内  
 (72) 発明者 中尾 康二  
 東京都小金井市貫井北町4-2-1 独立  
 行政法人情報通信研究機構内  
 (72) 発明者 衛藤 将史  
 東京都小金井市貫井北町4-2-1 独立  
 行政法人情報通信研究機構内

最終頁に続く

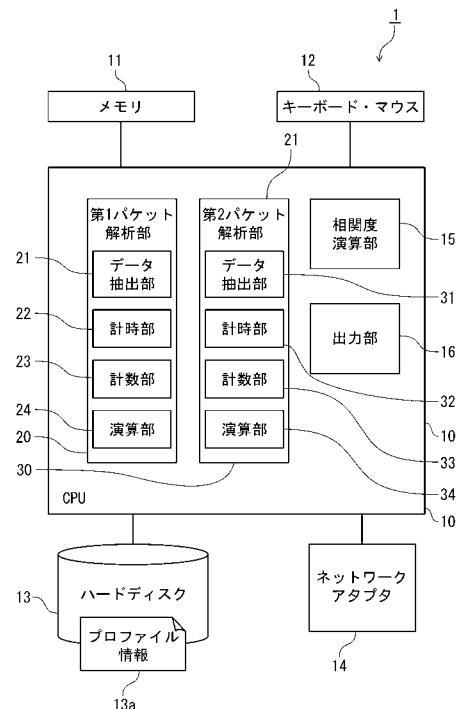
(54) 【発明の名称】 データ類似性検査方法及び装置

(57) 【要約】

【課題】 他のコンピュータに対して不正な処理を行うマルウェア等のソフトウェアと、検査するソフトウェアとの相関関係を高速、的確に自動算出する技術を提供すること。

【解決手段】 ネットワーク上で他のコンピュータに対して不正処理を行うマルウェアが送信するパケットと、検査対象ソフトウェアが送信するパケットを、所定のパラメータに着目してその相関度を算出することによりデータの類似性を検査する。パラメータの種類をプロファイル情報13aとして備えておき、第1パケット解析手段20が、それを参照してパラメータを取得する。これを第1のプロファイルデータとする。同様に第2パケット解析手段30が、第2のプロファイルデータを作成する。そして、相関度演算手段15が、両プロファイルデータから相関関係式を用いて相関度を算出する。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項 1】

ネットワーク上で他のコンピュータに対して不正処理を行う第 1 のソフトウェアが送信する第 1 の送信データと、検査対象の第 2 のソフトウェアが送信する第 2 の送信データとを比較してその類似性を検査するコンピュータのデータ類似性検査方法であって、

該送信データに関して、少なくとも送信されるパケットに含まれる情報のいずれかについて予め着目する単数又は複数のパラメータの種類をプロファイル情報として記憶手段に格納しておく、

第 1 パケット解析手段が、該プロファイル情報を参照し、第 1 のソフトウェアから送信される第 1 パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第 1 のプロファイルデータとして記憶手段に記録する第 1 パケット解析ステップ、

第 2 パケット解析手段が、該プロファイル情報を参照し、第 2 のソフトウェアから送信される第 2 パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第 2 のプロファイルデータとして記憶手段に記録する第 2 パケット解析ステップ、

コンピュータの相関度演算手段が、該第 1 のプロファイルデータと、該第 2 のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出する相関度演算ステップ、

該相関度を出力する出力ステップ

を有することを特徴とするデータ類似性検査方法。

## 【請求項 2】

前記プロファイル情報に、

パケットの送信先ポートのポート番号の一部又は全部をパラメータの種類として格納しておく、

前記第 1 及び第 2 パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出する

ことを特徴とする請求項 1 に記載のデータ類似性検査方法。

## 【請求項 3】

前記プロファイル情報に、

パケットの送信先ポートのポート番号と各ポート番号に対する送信回数をパラメータの種類として格納しておく、

前記第 1 及び第 2 パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出すると共に、所定時間における若しくはソフトウェアの所定回数の実行処理における各ポート番号に対するパケットの送信回数を計数する

ことを特徴とする請求項 1 又は 2 に記載のデータ類似性検査方法。

## 【請求項 4】

前記プロファイル情報に、

パケットの送信先ポートのポート番号間の遷移情報をパラメータの種類として格納しておく、

前記第 1 及び第 2 パケット解析手段がそれぞれ、送信先ポート番号の遷移情報をパケットから抽出する

ことを特徴とする請求項 1 ないし 3 のいずれかに記載のデータ類似性検査方法。

## 【請求項 5】

前記プロファイル情報に、

パケットの送信先アドレスが複数ある場合における該送信先アドレス間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておく、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各送信先アドレスをパケットから抽出すると共に、該送信先アドレス間の差分値又は差分値に基づく所定の統計値を算出する

ことを特徴とする請求項 1 ないし 4 のいずれかに記載のデータ類似性検査方法。

10

20

30

40

50

**【請求項 6】**

前記プロファイル情報に、

パケットの送信先アドレスと送信する送信先アドレスの数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信先アドレスの数を計数することを特徴とする請求項 1 ないし 5 のいずれかに記載のデータ類似性検査方法。

**【請求項 7】**

前記プロファイル情報に、

パケットの送信元ポートのポート番号の数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信元ポート番号の数を計数することを特徴とする請求項 1 ないし 6 のいずれかに記載のデータ類似性検査方法。

**【請求項 8】**

前記プロファイル情報に、

パケットの送信元ポートのポート番号が複数ある場合における該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各送信元ポート番号をパケットから抽出すると共に、該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を算出する

ことを特徴とする請求項 1 ないし 7 のいずれかに記載のデータ類似性検査方法。

**【請求項 9】**

前記プロファイル情報に、

パケットが準拠するプロトコル、又は、TCP プロトコルにおけるフラグの少なくともいずれかの種類又はその数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、該プロトコル又はフラグを該パケットから取得すると共に、数がパラメータの種類である場合にはその数を計数する

ことを特徴とする請求項 1 ないし 8 のいずれかに記載のデータ類似性検査方法。

**【請求項 10】**

前記プロファイル情報に、

パケットに含まれるペイロード(Payload)成分のダイジェスト(Digest)値又はその数又はそのいずれかに基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各ペイロードのダイジェスト値をパケットから抽出すると共に、数がパラメータの種類である場合にはダイジェスト値の数を計数し、該ダイジェスト値若しくはその数に基づく所定の統計値がパラメータの種類である場合にはそれを算出する

ことを特徴とする請求項 1 ないし 9 のいずれかに記載のデータ類似性検査方法。

**【請求項 11】**

前記プロファイル情報に、

パケットに含まれるペイロード(Payload)成分のデータ量又は該データ量に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、該ペイロード成分をパケットから抽出すると共に、該データ量を測定する、又は測定した該データ量に基づく所定の統計値を算出する

ことを特徴とする請求項 1 ないし 10 のいずれかに記載のデータ類似性検査方法。

**【請求項 12】**

前記コンピュータのデータ類似性検査方法において、

該送信データに関して、少なくとも送信されるパケットが所定時間に送信先ポートに到

10

20

30

40

50

達する数をパラメータの種類を含めてプロファイル情報として記憶手段に格納しておき、  
前記相関度演算ステップの前のいずれかの時点で

前記第1パケット解析手段が、該プロファイル情報を参照し、第1のソフトウェアから送信される第1パケットを所定時間に受信する数を計数して第1のプロファイルデータとして記憶手段に記録する第1パケット受信回数計数ステップ、

前記第2パケット解析手段が、該プロファイル情報を参照し、第2のソフトウェアから送信される第2パケットを所定時間に受信する数を計数して第2のプロファイルデータとして記憶手段に記録する第2パケット受信回数計数ステップ

の各ステップを行う

ことを特徴とする請求項1ないし11のいずれかに記載のデータ類似性検査方法。

10

【請求項13】

前記第1のソフトウェアが、閉じられたネットワークにおいて検査のために実行される既知のマルウェアであり、前記第2のソフトウェアが、広域ネットワークにおいて実際に実行される検査対象のソフトウェアであって、

前記請求項1ないし12に記載のデータ類似性検査方法における前記第1パケット解析ステップを単数又は複数の第1のソフトウェアについて実行処理した後、

コンピュータの検査開始指示信号送出手段が、該第2のソフトウェアからの送信データ内容又は送信データの packets に基づいて検査開始指示信号を送出する検査開始指示ステップ、

該検査開始指示信号を契機として実行処理される前記第2パケット解析ステップ、

20

コンピュータの相関度演算手段が、該単数又は複数の第1のプロファイルデータと、該第2のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出する相関度演算ステップ、

該第1のソフトウェアとの各相関度の少なくとも一部を出力するマルウェア相関度出力ステップ

を有することを特徴とするマルウェアの検査方法。

【請求項14】

ネットワーク上で他のコンピュータに対して不正処理を行う第1のソフトウェアが送信する第1の送信データと、検査対象の第2のソフトウェアが送信する第2の送信データとを比較してその類似性を検査するデータ類似性検査装置であって、

30

該送信データに関して、少なくとも送信される packets に含まれる情報のいずれかについて予め着目する単数又は複数のパラメータの種類をプロファイル情報として格納する記憶手段と、

該プロファイル情報を参照し、第1のソフトウェアから送信される第1パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第1のプロファイルデータとして記憶手段に記録する第1パケット解析手段と、

該プロファイル情報を参照し、第2のソフトウェアから送信される第2パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第2のプロファイルデータとして記憶手段に記録する第2パケット解析手段と、

40

該第1のプロファイルデータと、該第2のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出するコンピュータの相関度演算手段と、

該相関度を出力する出力手段と

を備えたことを特徴とするデータ類似性検査装置。

【請求項15】

前記プロファイル情報に、

パケットの送信先ポートのポート番号の一部又は全部をパラメータの種類として格納しておき、

前記第1及び第2パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出する

50

ことを特徴とする請求項 1 4 に記載のデータ類似性検査装置。

【請求項 1 6】

前記プロファイル情報に、

パケットの送信先ポートのポート番号と各ポート番号に対する送信回数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出すると共に、所定時間における若しくはソフトウェアの所定回数の実行処理における各ポート番号に対するパケットの送信回数を計数する

ことを特徴とする請求項 1 4 又は 1 5 に記載のデータ類似性検査装置。

【請求項 1 7】

前記プロファイル情報に、

パケットの送信先ポートのポート番号間の遷移情報をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、送信先ポート番号の遷移情報をパケットから抽出する

ことを特徴とする請求項 1 4 ないし 1 6 のいずれかに記載のデータ類似性検査装置。

【請求項 1 8】

前記プロファイル情報に、

パケットの送信先アドレスが複数ある場合における該送信先アドレス間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各送信先アドレスをパケットから抽出すると共に、該送信先アドレス間の差分値又は差分値に基づく所定の統計値を算出する

ことを特徴とする請求項 1 4 ないし 1 7 のいずれかに記載のデータ類似性検査装置。

【請求項 1 9】

前記プロファイル情報に、

パケットの送信先アドレスと送信する送信先アドレスの数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信先アドレスの数を計数する

ことを特徴とする請求項 1 4 ないし 1 8 のいずれかに記載のデータ類似性検査装置。

【請求項 2 0】

前記プロファイル情報に、

パケットの送信元ポートのポート番号の数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信元ポート番号の数を計数する

ことを特徴とする請求項 1 4 ないし 1 9 のいずれかに記載のデータ類似性検査装置。

【請求項 2 1】

前記プロファイル情報に、

パケットの送信元ポートのポート番号が複数ある場合における該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各送信元ポート番号をパケットから抽出すると共に、該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を算出する

ことを特徴とする請求項 1 4 ないし 2 0 のいずれかに記載のデータ類似性検査装置。

【請求項 2 2】

前記プロファイル情報に、

パケットが準拠するプロトコル、又は、TCP プロトコルにおけるフラグの少なくともいずれかの種類又はその数をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、該プロトコル又はフラグを該パケット

10

20

30

40

50

から取得すると共に、数がパラメータの種類である場合にはその数を計数することを特徴とする請求項 1 4 ないし 2 1 のいずれかに記載のデータ類似性検査装置。

【請求項 2 3】

前記プロファイル情報に、  
パケットに含まれるペイロード(Payload)成分のダイジェスト(Digest)値又はその数又はそのいずれかに基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、各ペイロードのダイジェスト値をパケットから抽出すると共に、数がパラメータの種類である場合にはダイジェスト値の数を計数し、該ダイジェスト値若しくはその数に基づく所定の統計値がパラメータの種類である場合にはそれを算出する

10

ことを特徴とする請求項 1 4 ないし 2 2 のいずれかに記載のデータ類似性検査装置。

【請求項 2 4】

前記プロファイル情報に、  
パケットに含まれるペイロード(Payload)成分のデータ量又は該データ量に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、

前記第 1 及び第 2 パケット解析手段がそれぞれ、該ペイロード成分をパケットから抽出すると共に、該データ量を測定する、又は測定した該データ量に基づく所定の統計値を算出する

ことを特徴とする請求項 1 4 ないし 2 3 のいずれかに記載のデータ類似性検査装置。

20

【請求項 2 5】

前記コンピュータのデータ類似性検査装置において、  
記憶手段が、該送信データに関して、少なくとも送信されるパケットが所定時間に送信先ポートに到達する数をパラメータの種類に含めてプロファイル情報として備えると共に、

前記第 1 パケット解析手段が、該プロファイル情報を参照し、第 1 のソフトウェアから送信される第 1 パケットを所定時間に受信する数を計数して第 1 のプロファイルデータとして記憶手段に記録し、

前記第 2 パケット解析手段が、該プロファイル情報を参照し、第 2 のソフトウェアから送信される第 2 パケットを所定時間に受信する数を計数して第 2 のプロファイルデータとして記憶手段に記録する

30

ことを特徴とする請求項 1 4 ないし 2 4 のいずれかに記載のデータ類似性検査装置。

【請求項 2 6】

マルウェアを検査するマルウェア検査システムであって、

前記請求項 1 4 ないし 2 5 のいずれかに記載のデータ類似性検査装置を備え、

前記第 1 のソフトウェアが、閉じられたネットワークにおいて検査のために実行される既知のマルウェアであり、前記第 2 のソフトウェアが、広域ネットワークにおいて実際に実行される検査対象のソフトウェアである構成において、

前記第 1 パケット解析手段が、単数又は複数の第 1 のソフトウェアについて各第 1 のプロファイルデータを記録すると共に、

40

該第 2 のソフトウェアからの送信データ内容又は送信データのパケットに基づいて検査開始指示信号を送出するコンピュータの検査開始指示信号送出手段を備え、

前記第 2 パケット解析手段が、該検査開始指示信号を受信すると作動し、

前記相関度演算手段が、該単数又は複数の第 1 のプロファイルデータと、該第 2 のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出し、

前記出力手段が、該第 1 のソフトウェアとの各相関度の少なくとも一部を出力する

を有することを特徴とするマルウェア検査システム。

【発明の詳細な説明】

【技術分野】

50

## 【 0 0 0 1 】

本発明はソフトウェアが出力するデータの類似性を検査する方法と装置に関し、特にネットワーク上で不正な処理を行うコンピュータから送信されるパケットに基づいてそのデータの類似性あるいは、そのソフトウェアの類似性を検査する方法と装置に係る技術である。

## 【 背景技術 】

## 【 0 0 0 2 】

インターネットにおけるインシデント対策の研究分野では、広域ネットワークでのパッシブモニタリングを行い、観測されたトラフィックを分析することで、インシデント検知を行うための研究が盛んに行われている。

10

また、本件発明者らが推進するインシデント対策のためのプロジェクトnicter（非特許文献1を参照。）では、広域観測網において観測されたトラフィックから、実時間でインシデントを検知する技術が研究されている。

広域ネットワークにおいて実際のインシデントを解析する技術をここではマクロ解析と呼ぶこととする。

## 【 0 0 0 3 】

その一方で、ウィルス(virus)、ワーム(worm)、ボット(bot)といったマルウェア(malware)検体を収集・分析し、個々のマルウェアの特徴を抽出する技術も研究が進められている。このように閉じられたネットワーク空間において、マルウェア検体の分析を行うことを、上記のマクロ解析に対して、ミクロ解析と呼ぶこととする。

20

## 【 0 0 0 4 】

マルウェアに起因するインシデントに迅速かつ的確に対処するためには、広域観測網において検出された事象(結果)に対し、その原因となったマルウェアを特定し、提示することが重要である。

このようなインシデント(結果)とマルウェア(原因)との相関関係を得るためには、それぞれの特徴を効果的に抽出した上で相関分析を行う必要がある。

## 【 0 0 0 5 】

ミクロ解析においてスキャン攻撃の特徴抽出手法としていくつかの先行研究が提案されているが、広域ネットワークでのインシデントとマルウェアとの相関分析を行うことを前提とする、個々のホストのネットワーク的挙動を分析する研究はいまだ少ない。すなわち、マクロ解析結果とミクロ解析結果との相関関係を検査して、マクロ解析において得られた特定のホストについてマルウェアの特定を行う技術はほとんど提供されていない。

30

## 【 0 0 0 6 】

マルウェアの自動解析方法としては、特許文献1のような技術が知られている。

## 【 0 0 0 7 】

【非特許文献1】中尾康二、吉岡克成、衛藤将史、井上大介、力武健次著「nicter: An Incident Analysis System using Correlation between Network Monitoring and Malware Analysis」Proceedings of The 1st Joint Workshop on Information Security, JWIS2006, Page363-377, 2006年9月

40

【特許文献1】特表2006-522395号公報

## 【 発明の開示 】

## 【 発明が解決しようとする課題 】

## 【 0 0 0 8 】

このように、従来の技術においてはミクロ的なインシデントの解析、マクロ的なマルウェアの解析が別個に行われており、ネットワーク上で発生しているインシデントが、解析済みのどのマルウェアの作用によるものなのか、特定は人手に頼らざるを得なかった。

自動的に特定する場合にも、挙動が極めて類似していれば正解が得られやすいが、実際のマルウェアは、検出されにくいように様々な挙動をするように作成されているため、受信するコンピュータによってその態様は多様に変化する。

50

## 【 0 0 0 9 】

本発明は、このように他のコンピュータに対して不正な処理を行うマルウェア等のソフトウェアと、検査するソフトウェアとの相関関係を高速、的確に自動算出する技術を提供するものである。

【課題を解決するための手段】

【0010】

本発明は、上記の課題を解決するために、次のようなデータ類似性検査方法を提供する。

すなわち、請求項1に記載の発明は、ネットワーク上で他のコンピュータに対して不正処理を行う第1のソフトウェアが送信する第1の送信データと、検査対象の第2のソフトウェアが送信する第2の送信データとを比較してその類似性を検査するコンピュータのデータ類似性検査方法を提供するものである。

10

本方法においては、送信データに関して、少なくとも送信されるパケットに含まれる情報のいずれかについて予め着目する単数又は複数のパラメータの種類をプロファイル情報として記憶手段に格納しておく。

【0011】

そして、第1パケット解析手段が、プロファイル情報を参照し、第1のソフトウェアから送信される第1パケットを受信して、そこに格納されているパラメータの種類に該当するパラメータの一部又は全部を取得する。これを第1のプロファイルデータとして記憶手段に記録する。(第1パケット解析ステップ)

【0012】

20

その上で、第2パケット解析手段が、プロファイル情報を参照し、第2のソフトウェアから送信される第2パケットを受信して、パラメータの種類に該当するパラメータの一部又は全部を取得して第2のプロファイルデータとして記憶手段に記録する。(第2パケット解析ステップ)

【0013】

続く相関度演算ステップにおいて、コンピュータの相関度演算手段が、該第1のプロファイルデータと、該第2のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出し、出力ステップにおいて相関度を出力する。

【0014】

30

また、本発明は次のようなデータ類似性検査装置を提供することもできる。

すなわち、請求項14に記載のように、ネットワーク上で他のコンピュータに対して不正処理を行う第1のソフトウェアが送信する第1の送信データと、検査対象の第2のソフトウェアが送信する第2の送信データとを比較してその類似性を検査するデータ類似性検査装置を提供する。

【0015】

本装置において、送信データに関して、少なくとも送信されるパケットに含まれる情報のいずれかについて予め着目する単数又は複数のパラメータの種類をプロファイル情報として格納する記憶手段を備える。

また、プロファイル情報を参照し、第1のソフトウェアから送信される第1パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第1のプロファイルデータとして記憶手段に記録する第1パケット解析手段と、プロファイル情報を参照し、第2のソフトウェアから送信される第2パケットを受信して、該パラメータの種類に該当するパラメータの一部又は全部を取得して第2のプロファイルデータとして記憶手段に記録する第2パケット解析手段を備える。

40

【0016】

以上の各パケット解析手段で記録された各プロファイルデータを利用し、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出するコンピュータの相関度演算手段、相関度を出力する出力手段を備える。

【0017】

50



本発明において、上記プロファイル情報に、パケットの送信先ポートのポート番号の一部又は全部をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出する構成でもよい。

【0018】

上記プロファイル情報に、パケットの送信先ポートのポート番号と各ポート番号に対する送信回数をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、送信先ポート番号をパケットから抽出すると共に、所定時間における若しくはソフトウェアの所定回数の実行処理における各ポート番号に対するパケットの送信回数を計数する構成でもよい。

【0019】

上記プロファイル情報に、パケットの送信先ポートのポート番号間の遷移情報をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、送信先ポート番号の遷移情報をパケットから抽出する構成でもよい。

【0020】

上記プロファイル情報に、パケットの送信先アドレスが複数ある場合における該送信先アドレス間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、各送信先アドレスをパケットから抽出すると共に、該送信先アドレス間の差分値又は差分値に基づく所定の統計値を算出する構成でもよい。

【0021】

上記プロファイル情報に、パケットの送信先アドレスと送信する送信先アドレスの数をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信先アドレスの数を計数する構成でもよい。

【0022】

上記プロファイル情報に、パケットの送信元ポートのポート番号の数をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、所定時間において若しくはソフトウェアの所定回数の実行処理においてパケットを送信する送信元ポート番号の数を計数する構成でもよい。

【0023】

上記プロファイル情報に、パケットの送信元ポートのポート番号が複数ある場合における該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、各送信元ポート番号をパケットから抽出すると共に、該送信元ポート番号間の差分値又は差分値に基づく所定の統計値を算出する構成でもよい。

【0024】

上記プロファイル情報に、パケットが準拠するプロトコル、又は、TCPプロトコルにおけるフラグの少なくともいずれかの種類又はその数をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、該プロトコル又はフラグを該パケットから取得すると共に、数がパラメータの種類である場合にはその数を計数する構成でもよい。

【0025】

上記プロファイル情報に、パケットに含まれるペイロード(Payload)成分のダイジェスト(Digest)値又はその数又はそのいずれかに基づく所定の統計値を導く計算をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、各ペイロードのダイジェスト値をパケットから抽出すると共に、数がパラメータの種類である場合にはダイジェスト値の数を計数し、該ダイジェスト値若しくはその数に基づく所定の統計値がパラメータの種類である場合にはそれを算出する構成でもよい。

【0026】

上記プロファイル情報に、パケットに含まれるペイロード(Payload)成分のデータ量又

10

20

30

40

50

は該データ量に基づく所定の統計値を導く計算をパラメータの種類として格納しておき、第1及び第2パケット解析手段がそれぞれ、該ペイロード成分をパケットから抽出すると共に、該データ量を測定する、又は測定した該データ量に基づく所定の統計値を算出する構成でもよい。

【0027】

また、本発明は上記のデータ類似性検査方法において、該送信データに関して、少なくとも送信されるパケットが所定時間に送信先ポートに到達する数をパラメータの種類に含めてプロファイル情報として記憶手段に格納しておくこともできる。

本構成において、上記の相関度演算ステップの前のいずれかの時点で第1パケット解析手段が、該プロファイル情報を参照し、第1のソフトウェアから送信される第1パケットを所定時間に受信する数を計数して第1のプロファイルデータとして記憶手段に記録する第1パケット受信回数計数ステップ、第2パケット解析手段が、該プロファイル情報を参照し、第2のソフトウェアから送信される第2パケットを所定時間に受信する数を計数して第2のプロファイルデータとして記憶手段に記録する第2パケット受信回数計数ステップの各ステップを行う。

本技術を実装したデータ類似性検査装置を提供することもできる。

【0028】

本発明は、上記データ類似性検査方法の技術により、次のようなマルウェアの検査方法を提供してもよい。

すなわち、請求項13に記載の発明によれば、上記第1のソフトウェアとして、閉じられたネットワークにおいて検査のために実行される既知のマルウェアを、第2のソフトウェアとして、広域ネットワークにおいて実際に実行される検査対象のソフトウェアをそれぞれ用い、上記の第1パケット解析ステップを単数又は複数の第1のソフトウェアについて実行処理する。

【0029】

その後、コンピュータの検査開始指示信号送出手段が、該第2のソフトウェアからの送信データ内容又は送信データのパケットに基づいて検査開始指示信号を送出する検査開始指示ステップ、この検査開始指示信号を契機として実行処理される前記第2パケット解析ステップ、コンピュータの相関度演算手段が、該単数又は複数の第1のプロファイルデータと、該第2のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出する相関度演算ステップ、該第1のソフトウェアとの各相関度の少なくとも一部を出力するマルウェア相関度出力ステップを順に実行する。

【0030】

さらに、マルウェアを検査するマルウェア検査システムとして提供することもできる。

すなわち、上記のデータ類似性検査装置を備え、第1のソフトウェアが、閉じられたネットワークにおいて検査のために実行される既知のマルウェアであり、前記第2のソフトウェアが、広域ネットワークにおいて実際に実行される検査対象のソフトウェアである構成において、第1パケット解析手段が、単数又は複数の第1のソフトウェアについて各第1のプロファイルデータを記録すると共に、第2のソフトウェアからの送信データ内容又は送信データのパケットに基づいて検査開始指示信号を送出するコンピュータの検査開始指示信号送出手段を備える。

【0031】

そして、第2パケット解析手段が、該検査開始指示信号を受信すると作動し、相関度演算手段が、該単数又は複数の第1のプロファイルデータと、該第2のプロファイルデータとから、予め記憶手段に格納されるパラメータの種類に応じた相関関係式を用いて相関度を算出し、出力手段が、該第1のソフトウェアとの各相関度の少なくとも一部を出力する。

【発明の効果】

【0032】

10

20

30

40

50

本発明は、以上の構成を備えることにより、次の効果を奏する。

まず、2つのソフトウェアが送信するデータの類似性を検査する際に、両ソフトウェアから出力されるパケットから、好適なパラメータの種類に着目し、これをプロファイル情報として格納することができる。

【0033】

そしてプロファイル情報に基づいてパラメータを抽出すると共に、パラメータ毎に設定してある相関関係式に基づいて両者の比較を行う。本方法によれば、パケットの送信態様が複雑で、しかも頻繁に変化するマルウェアについても複数のパラメータを融合して相関度を求めることが容易であり、その結果として高精度に相関度を算出することができる。

【0034】

またプロファイル情報を書き換えるだけで新しいマルウェアの検出にも迅速簡便に対応することが可能であり、急速に、かつ多岐にわたって進化するマルウェアの検査方法としても最適である。

【0035】

特に本発明では、パラメータの種類として特に好適なものを提供するので、マルウェアの一般的な不正処理のパケットに対し、効率よく相関度を求めることができる。

【発明を実施するための最良の形態】

【0036】

以下、本発明の実施形態を、図面に示す実施例を基に説明する。なお、実施形態は下記に限定されるものではない。

図1は本発明に係るデータ類似性検査装置(以下、本装置と呼ぶ。)(1)の全体構成図である。本装置(1)は、公知のパーソナルコンピュータやネットワークサーバによって構成するのが簡便である。

【0037】

本装置(1)には、演算処理等を司るCPU(10)を中心として、CPU(1)と協働するメモリ(11)、ユーザが入力等を行うキーボード及びマウス(12)、データを読み書き自在に格納するハードディスク(13)、インターネット等のネットワーク接続を行うネットワークアダプタ(14)などが備えられている。また、図示しないモニタを接続して画面表示を行ったり、スピーカを接続して音声出力を行うことも可能である。

これらの構成はいずれも周知の事項であって、その構造や作用については説明を省略する。

【0038】

本発明は、上記CPU(10)に、第1パケット解析部(20)と、第2パケット解析部(30)と、相関度演算部(15)と、出力部(16)とを備える。第1及び第2パケット解析部は、以下説述するさまざまなパラメータを取得するために、データ抽出部(21)(31)、計時部(22)(32)、計数部(23)(33)、演算部(24)(34)を備える。

【0039】

はじめに、本発明の最も基本的な処理を図2を用いて説明する。本実施例では不正処理を行う第1のソフトウェアとしてマルウェアを、第2のソフトウェアとして検査対象となるソフトウェアを用いる。

まず、プロファイル情報(13a)を読み出す。(プロファイル読出処理:S10)

ここでプロファイル情報(13a)とは、ソフトウェアから送信されるデータの類似性を検査するために着目するパラメータの種類とその抽出方法を予め定義したものである。パラメータの種類として本発明で用いるのは、例えばパケットに含まれる送信先ホストにおけるポート番号や、送信元ポートのポート番号などであり、それぞれについて、TCP(Transmission Control Protocol)(RFC793により規定されている)のヘッダ部分における最初の0~15ビット(送信元ポート番号)、16ビット~31ビット(送信先ポート番号)に配置されていることが記録される。

本発明ではこのパラメータの種類としていかなるものを用いるのかについても重要な意

10

20

30

40

50

義がある。

【0040】

パラメータの種類をハードディスクに格納する場合に、具体的にはコンピュータで用いられる変数の名称として記録されている。例えば、送信先ホストにおけるポート番号という種類を表す変数は「DstPort」としており、送信元ホストのポート番号は「SrcPort」としている。そして、このパラメータとしては「139」「445」などのデータである。以下、パラメータの種類とパラメータとはこのように区別して使用する。

【0041】

最初のプロファイル情報読出処理（S10）では、まずこのようにパラメータの種類（変数及びその抽出方法に関するデータ）を読み出し処理する。具体的な処理方法は、CPU（10）とハードディスク（13）の協働によるものであり公知である。

そして、読み出したプロファイル情報（13a）に基づいてマルウェアからのパケットを解析する。（第1パケット解析処理：S11）

【0042】

ここでの解析方法は、プロファイル情報によって様々であるから、具体的なパラメータの種類と共に解析方法を後述する。

そして、得られたパラメータのデータを、第1プロファイルデータ（13b）としてハードディスク（13）やメモリ（11）などの記憶手段に記録する。（プロファイル記録処理：S12）

上記処理（S10～S12）は第1パケット解析部（20）による処理である。

【0043】

本発明では以上の各処理を経て、マルウェアに関する第1のプロファイルデータを生成する。1つのマルウェアに関してのプロファイルデータを作成するだけでもよいが、多数のマルウェアのプロファイルデータを生成しておくことが好ましく、これによって次の検査対象ソフトウェアとの比較の際に、多数の候補から最も相関が認められるマルウェアを特定することができる。

【0044】

検査対象ソフトウェアについてもマルウェアと同様の処理を行う。すなわち、プロファイル情報（13a）を読み出し（プロファイル情報読出処理：S20）、検査対象ソフトウェアが送信するパケットをプロファイル情報に基づいて解析処理する。（第2パケット解析処理：S21）

そして、取得されたパラメータを第2プロファイルデータ（13c）として記録する。（プロファイル記録処理：S22）

以上の各処理（S20～S22）は第2パケット解析部（30）による処理である。

【0045】

以上の第2プロファイルデータ（13c）と、予め記録されている第1プロファイルデータ（13b）とを比較するために、相関度演算部（15）が、予め定義された相関関係式に基づいて相関度を算出する。（相関度算出処理：S30）

相関関係式は、プロファイル情報のパラメータの種類ごとに定義しておくことが好ましく、この場合にはプロファイル情報にパラメータの種類に対応して格納される。もっとも、本発明では相関関係式は1種類として、ハードディスク（13）に別に格納しておいてもよい。相関関係式についても後述する。

【0046】

そして、出力部（16）からマルウェアごとの相関度を出力する。（出力処理：S31）

出力の方法は、一覧形式にしてモニタからの画面表示、プリンタからの印刷出力、ネットワークアダプタ（14）を介して他のコンピュータへのデータ送信など、いかなる方法でもよい。CPU（10）の出力部（16）は公知のハードウェアと協働して同処理を司る。

【0047】

出力の方法は、一覧形式にしてモニタからの画面表示、プリンタからの印刷出力、ネットワークアダプタ（14）を介して他のコンピュータへのデータ送信など、いかなる方法でもよい。CPU（10）の出力部（16）は公知のハードウェアと協働して同処理を司る。

10

20

30

40

50

本発明では、さらにマルウェアの特定方法又はマルウェアの特定システムとして提供することもできる。基本的な技術としては上記データ類似性検査装置と同様であるが、本件出願人らが提案しているnicter（非特許文献1参照）に組み合わせてマルウェアの特定を行うことを提案する。

【0048】

図3において、まず広域ネットワーク（40）に複数設けたセンサー（41）でダークネットに対するパケットなどを検知し、マクロ解析器（42）に入力する。マクロ解析の結果はデータベース（43）に格納される。

【0049】

一方、ネットワーク（44）上で、キャプチャ（45）によって多数のマルウェア検体を採集し、ミクロ解析器（46）によりその静的、動的な性質を解析する。その解析結果もデータベース（47）に格納する。

【0050】

このように、実際にインシデントを発生させているマルウェアをマクロ解析器によってマクロ的に解析すると共に、検体を解析してマルウェアのミクロ的な解析を行い、それぞれのデータベースから相関分析器（48）で相関分析を行うことが考えられている。

【0051】

相関分析の結果はデータベース（49）に格納されて、さまざまな出力方法によるインシデントハンドリングシステム（50）を介してユーザ（51）に通知されたり、レポート（52）として出力されたりする。

【0052】

このシステムに対して、本発明を適用し、マクロ解析器（42）とミクロ解析器（46）とともに、相関分析器（48）に本発明のデータ類似性検査装置（1）を実装する。

もっとも、本装置（1）を分割してマクロ解析器（42）に第2パケット解析部（30）を、ミクロ解析器（46）に第1パケット解析部（20）を備えて、それぞれの挙動を検出すると共に、その結果を、相関度演算部（15）を備えた相関分析器（48）で相関分析してもよい。

【0053】

従来、マクロ解析とミクロ解析の結果を融合することが技術的に困難であったが、本発明の方法を適用することによって、これが実現され、広域ネットワークで生じているインシデントの原因を高速、的確に特定することができる。

【0054】

ミクロ解析器（46）における動作として、図4に示すようなマルウェア動的解析環境を用いた解析が好適である。

すなわち、マルウェア動的解析環境とは仮想的にローカルアドレス空間、グローバルアドレス空間を設けて、感染したホストが各空間においてどのような挙動を示すかを測定する箱庭環境である。

【0055】

そして、感染したホスト（70）から仮想グローバル空間（60）の各IPアドレスのホスト（61）に対してどのようにグローバルアドレスをスキャンするか、あるいはダミーのIRCサーバ（62）、HTTPサーバ（63）、FTPサーバ（64）、TFTPサーバ（65）に対してどのようにアクセスするかなどを測定する。

このような仮想グローバル空間（60）を用いることで、グローバルアドレス空間への挙動を測定することができる。

【0056】

また、仮想ローカル空間（66）においては、感染したホスト（70）から複数のローカルIPアドレスのホスト（67）に対してどのようにローカルアドレスをスキャンするか、あるいはダミーのSMTPサーバ（68）、DNSサーバ（69）に対してどのようにアクセスするのかなどを測定する。

このような仮想ローカル空間（66）を用いることで、ローカルアドレス空間への挙動

10

20

30

40

50

を測定することができる。

【0057】

一方、ダークネットと呼んでいる、実際には使用されていないIPアドレス領域に対して送信されるパケットを広域ネットワーク(40)上のセンサー(44)で検知する。

このようなIPアドレスに向けたパケットは規則に準じたホストに向けたものではないから、設定ミスか、ワームによるスキャン、探索、後方散乱メールなどの悪意による処理と考えられる。

【0058】

ここでセンサー(44)はダークネットにあたるIPアドレスが付与されたネットワーク端末で構成する。

さらに、本発明ではセンサー(44)又はマクロ解析器(42)に公知の不正な挙動を検出する機能を付与する。例えば非特許文献2に該方法のアルゴリズムが記載されているが、マクロ解析器の構成は非特許文献3などにも開示されるようにさまざまな技術を適用することができる。以下ではセンサー(44)に同機能を付与したものと説明する。

【0059】

【非特許文献2】Suzuki, K., Baba, S., Takakura, H.: Analyzing traffic directed to unused IP address blocks, IEICE Technical Report, vol.105, no.530, IA2005-23, pp.25-30 2006年1月

【非特許文献3】Takeuchi, J., Sato, Y., Rikitake K., Nakao K.: Development of Incident Detection System Based on Change Point Detection, SCIS2006, The 2006 Symposium on Cryptography and Information Security, Japan, 2006年1月

【0060】

相関分析器(48)に本装置(1)を実装する場合には、検査開始指示信号として受信したパケットの生データを用いるのが簡便であり、該信号に基づいて図2のプロファイル情報読出処理(S20)から相関度算出処理(S30)、出力処理(S31)までを行ってもよい。

【0061】

このとき、図5に示すような装置構成をとってもよい。

すなわち、図3の相関分析器(48)において、プロファイルデータ生成部(80)と相関度演算部(82)を設け、プロファイルデータ生成部(80)には第1パケット解析部(20)と第2パケット解析部(30)の共通する処理手段を設け、相関度演算部(82)は本装置(1)の相関度演算部(15)と同様の手段とする。

【0062】

本構成において、プロファイルデータ生成部(80)が予め相関分析結果データベース(49)にマクロ解析から得たマルウェアにつき複数のプロファイルデータ(83)~(87)・・・を作成する。

また、ダークネットからの生データが検査開始指示信号として入力されると、同じプロファイルデータ生成部(80)が該生データに関するプロファイルデータ(81)を作成し、相関度演算部(82)で相関度を演算する。

【0063】

以上、説明したように、本発明における第1パケット解析部、第2パケット解析部などは1台の装置に実装せずに複数の装置に分割して備えてもよいし、また各パケット解析部における機能に応じて複数の装置に分割して配置してもよい。さらに、各パケット解析部は同様の処理内容をなすものであるから、これを上記のように1個の処理部で構成してもよい。

【0064】

次に、本発明に係るパケット解析部の詳細な処理内容を説述する。ここではプロファイル情報(13a)におけるパラメータの種類と重要な関連があるため、好適なパラメータの種類とあわせて説明する。

まず、本実施例で用いるパラメータの種類としては次の表1に掲げる項目がある。

10

20

30

40

50

【 0 0 6 5 】

【 表 1 】

データの内容	パラメータの種類 (変数)
送信先ポート	DstPort_Count
送信先ポート	DstPort_Trans
送信元ポート	SrcPort_Unique
送信元ポート	SrcPortDif_Stats
送信先IPアドレス	DstIPDif_Stats
送信先IPアドレス	DstIP_Unique
プロトコル	Protocol_Count
フラグ	Flag_Count
時間	NumPacketRate
ペイロード	PayloadSig_Count
ペイロード	Payload_Stats
TTL	TTL_Stats
ID番号	Id_Stats
シーケンス番号	SeqNum_Stats

10

【 0 0 6 6 】

(DstPort\_Count)

不正処理における送信データの packets をどの送信先ポート番号に送るかという情報は、マルウェアを特徴づけるもっとも基本的な情報である。

そこで、最初のパラメータとしてDstPort\_Countは、パケットの送信先ポートのポート番号と各ポート番号に対する送信回数を表す。そして、第1パケット解析部(20)におけるデータ抽出部(21)が、送信されたパケットのヘッダ部分から送信先ポート番号の部分を抽出する。

20

本実施例では、送信先ポート番号だけを用いてもよいが、送信先ポート番号に何回パケットを送ったかもパラメータとして利用する。

【 0 0 6 7 】

ここで、パケットを送った回数は、単位時間あたりの回数を計数する。すなわち、データ抽出部(21)で送信先ポート番号を抽出すると共に、計時部(22)で単位時間を計時し、その間の送信回数を計数部(23)で計数する。単位時間とは例えば1秒間など、適宜定められる時間である。

30

以下、他のパラメータにおいても送信回数はこれと同様にして求める。

【 0 0 6 8 】

(DstPort\_Trans)

次に、DstPort\_Transはパケットの送信先ポートのポート番号間の遷移情報である。遷移情報とは、例えば単位時間あるいは一連の当該ホストへのパケット送信において、送信先のポートがどのような順でスキャンされるかに関する情報である。

例えば、ポート番号139、445、3127、6659の順でスキャンした場合に、139の次が445であるという情報をパラメータとする。

より好適には、遷移確率を用いて、139の次が445である確率をパラメータとしてもよい。具体的には遷移順序を記憶しておき、各ポート番号について139の次にパケットが送られた割合を遷移確率とする。全てのポート番号について格納せず、例えば上位10位の遷移確率を格納するようにしてもよい。

40

【 0 0 6 9 】

(SrcPort\_Unique)

送信元ポート番号は、送信先ポート番号に比較してマルウェアが自由に設定できる特徴がある。マルウェアによっては固定した送信元ポート番号を用いるため、どのポートをいくつ用いているかをパラメータとして利用すると好適である。また、送信元ポート番号が変動する場合でも、どのようにポート番号を変化させているかがパラメータとして利用できる。

そこで、データ抽出部(22)が送信先ポート番号を抽出し、その番号をパラメータと

50

してプロファイルデータに格納する。

【 0 0 7 0 】

(SrcPortDif\_Stats)

また、送信元ポート番号の変動については、送信元ポート番号を順にデータ抽出部(21)で抽出して、前後のポート番号の差分を演算部(24)で求めると共に、その平均値を算出する。なお、平均値以外に任意の統計計算式を用いて統計値を求めてもよい。

これは送信元ポート番号をある値ごとに变化させて送信する場合などに好適な指標となる。

【 0 0 7 1 】

(DstIPDif\_Stats)

グローバルアドレス空間及びローカルアドレス空間において、どのようなIPアドレスに対してパケットを送ってくるかも利用することができる。IPアドレスをどのようにスキャンするのはマルウェアが標的を探す上で重要であり、その検索順序もさまざまな方法がある。サブネットマスクが16ビットか、24ビットか、近いIPアドレスから検索するのかランダムに検索するのか、などマルウェアの特定に役立つ情報である。

該情報は仮想ホストに割り当てたIPアドレスから取得することができる。

【 0 0 7 2 】

送信するIPアドレスが特異な特徴を持つ場合など、IPアドレス自体をパラメータとすることもできるし、IPアドレスを順に検出して、規定の方法で数値に変換し、演算部で差分を求め、その平均値を用いてもよい。平均値を他の統計値としてもよい。

【 0 0 7 3 】

(DstIP\_Unique)

DstIP\_Unique変数はいくつの送信先IPアドレスに送信するかに関するパラメータである。そこで、単位時間における送信先IPアドレスの数を計数する。マルウェアは単位時間に大量のパケットを送信するものと、一定の時間において送信するものがあり、その間隔を指標として用いることができる。単位時間や他所定の時間における回数でもよいし、一連のパケットの送信における送信頻度でもよい。

【 0 0 7 4 】

(Protocol\_Count)

Protocol\_Count変数は、マルウェアがどのプロトコルによっていくつのパケットを送信しているか、そのプロトコル毎のパケット数をパラメータとするものである。マルウェアは様々なプロトコルを利用するため、TCPによるのかUDPによるのか、どのプロトコルでスキャンされているか、データ抽出部(21)により特定する。

【 0 0 7 5 】

(Flag\_Count)

同様に、TCPをプロトコルとして用いている場合、TCPのどのフラグにどのようなパケットがくるかの情報を用いることもできる。

フラグの種類としては、URG(urgent)フラグ、ACK(acknowledge)フラグ、PSH(push)フラグ、RST(reset)フラグ、SYN(synchronize)フラグ、FIN(finis、終了)フラグなどがある。

そして、さまざまなスキャン方法が知られており、TCP SYNスキャン、TCP FINスキャン、TCPNullスキャン、TCP Xmasスキャン、TCP Maimon スキャン、TCP ACKスキャン、UDPスキャンなど、スキャンの種類をパラメータとしてもよい。

【 0 0 7 6 】

(NumPacketRate)

本発明では、単位時間あたりに送信先ホストにどれだけのパケットを送信してくるかの情報を用いることもできる。マルウェアによって連続的に送ってくる場合や、時間をおきながらパケットを送ってくる場合もある。

本発明では計時部(22)及び計数部(23)を用いてこれを計数してパラメータとする。

10

20

30

40

50



## 【 0 0 7 7 】

(PayloadSig\_Count)

UDPによる探索を行うマルウェアなど、パケットのペイロードフィールドが攻撃用のコードを含む場合がある。そこで、ペイロードのダイジェスト値を比較することも有効である。まず実際のデータ（ダイジェスト値）を抽出して、その値自体や、ダイジェスト値の数（種類数）をパラメータとして用いることができる。

## 【 0 0 7 8 】

(Payload\_Stats)

Payload\_Statsはペイロードのサイズの統計値である。平均値など任意の統計値を用いることができる。

10

## 【 0 0 7 9 】

(TTL\_Stats, Id\_Stats, SeqNum\_Stats)

その他、本実施例では、TCPにおけるTTLフィールド、IDフィールド、シーケンス番号フィールドなどの値や、その統計値を用いることもできる。

## 【 0 0 8 0 】

本発明は、以上説述した通りのパラメータの種類を用いることを提案し、これをプロファイルデータとして第1パケット解析部（20）が、データ抽出部（21）、計時部（22）、計数部（23）、演算部（24）によって解析する。

また、これと全く同様に、第2パケット解析部（30）が、データ抽出部（31）、計時部（32）、計数部（33）、演算部（34）によって第2のソフトウェアを解析する。

20

## 【 0 0 8 1 】

さらに、相関度演算部においては第1のプロファイルデータ（13b）と第2のプロファイルデータ（13c）から、任意の相関算出式により相関度を算出する。

次に、本実施例で用いた相関関係式について説明する。

## 【 0 0 8 2 】

本発明で用いたプロファイルデータは、ポート番号ごとに値をもつ場合など、多次元データとなる場合も多く、これらの多次元データ間の相関度を演算する必要がある。

例えば、第1のソフトウェアの使用したポート番号群 $P_A$ と、第2のソフトウェアが使用したポート番号群 $P_B$ に対して、 $p \in P_A$ 、 $q \in P_B$ とする。このとき各ポート番号における第1のソフトウェアの送信回数 $C_A(p)$ と第2のソフトウェアの送信回数 $C_B(q)$ とあらわすことができる。

30

## 【 0 0 8 3 】

$C_A(p)$ と $C_B(p)$ を2つの確率変数と考える。このとき $p \in P_A \cap P_B = \Omega$ （標本空間）である。そして、 $C_A(p)$ と $C_B(p)$ の相関度は次式（数1）によって表される。なお、 $m_A$ と $m_B$ は、それぞれ $C_A(p)$ と $C_B(p)$ の平均値である。

## 【 0 0 8 4 】

## 【 数 1 】

$$\frac{\sum_{p \in \Omega} \{C_A(p) - m_A\} \{C_B(p) - m_B\}}{\sqrt{\sum_{p \in \Omega} \{C_A(p) - m_A\}^2} \sqrt{\sum_{p \in \Omega} \{C_B(p) - m_B\}^2}}$$

40

## 【 0 0 8 5 】

相関度演算部（15）において、上記式に従って、相関度を算出することにより、多次元データの相関度が得られる。該演算方法は公知のプログラム手法によることができる。なお、同様に多次元データとなるのは、上記においてDstPort\_Count、Protocol\_Count、Flag\_Count、PayloadSig\_Countなどである。

## 【 0 0 8 6 】

一方、1次元のデータとしては、SrcPort\_Unique、DstIP\_Unique、NumPacketRateなど

50

があり、これらは単純に次式（数 2）によって類似度を求めることができる。

（数 2）

$$(a/b)^k$$

なお、 $a < b$ 、 $k$ は任意に定義可能であり、ここでは 1 を用いている。

【 0 0 8 7 】

相関度演算部（15）で用いる相関関係式としては、プロファイルデータによって任意であるが、パラメータがデータの系列である場合、次のような相関関係式を用いることができる。

【 0 0 8 8 】

10

(1) ピアソンの積率相関係数

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、ピアソンの積率相関係数は次式（数 3）によって求められる。

【 0 0 8 9 】

【数 3】

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

20

ここで、

$\bar{x}, \bar{y}$

はそれぞれデータ系列  $x, y$  の平均値をあらわす。

【 0 0 9 0 】

(2) スピアマンの順位相関係数

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、スピアマンの順位相関係数は次式（数 4）によって求められる。

【 0 0 9 1 】

【数 4】

$$1 - \frac{6 \sum_{i=1}^n D_i^2}{n(n^2 - 1)}$$

30

ここで、 $D_i$  は  $x_i$  の系列  $x$  内での順位と  $y_i$  の系列  $y$  内での順位の差を表す。

【 0 0 9 2 】

(3) ケンドールの順位相関係数

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、ケンドールの順位相関係数は次式（数 5）によって求められる。

【 0 0 9 3 】

40

【数 5】

$$\frac{K - L}{n(n-1)}$$

ここで、 $K$  は、 $i=1, 2, \dots, n$  について  $x_i$  の系列  $x$  内での順位が  $y_i$  の系列  $y$  内での順位より大きい（または順位が等しい）場合の数、 $L$  は、 $y_i$  の系列  $y$  内での順位が  $x_i$  の系列  $x$  内での順位より大きい（または順位が等しい）場合の数とする。

【 0 0 9 4 】

50

## (4) 特定条件を満たす割合

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  について、 $x_i, y_i$  が共に特定の条件を満たす場合の数を  $g_{\text{condition}}$  回とする。このとき、 $x$  と  $y$  の類似度を以下のように定義する。(数6)

(数6)

$$g_{\text{condition}}/n$$

例えば、 $x_i, y_i$  が実数値である場合、 $x_i > 0, y_i > 0$  が条件として考えられる。

【0095】

10

この他、類似度と逆の概念として距離(非類似度)が考えられるが、2つのデータ系列の間の距離の定義として以下がある。距離を相関分析に用いる場合は、2つのプロファイルデータ間で距離が小さいほど類似度が高いと考える。

【0096】

## (5) ユークリッド距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、ユークリッド距離は次式(数7)によって求められる。

【0097】

【数7】

$$\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

20

【0098】

## (6) ユークリッド平方距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、ユークリッド平方距離は次式(数8)によって求められる。

【0099】

【数8】

$$\sum_{i=1}^n (x_i - y_i)^2$$

30

【0100】

## (7) チェビシエフ距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、チェビシエフ距離は次式(数9)によって求められる。

(数9)

$$\max(|x_i - y_i|)$$

40

ここで  $|x|$  は  $x$  の絶対値をあらわす。また  $\max X$  は  $i = 1, 2, \dots, n$  について  $X$  の最大値を表す。

【0101】

## (8) マンハッタン距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、マンハッタン距離は次式(数10)によって求められる。

【0102】

## 【数 1 0】

$$\sum_{i=1}^n |x_i - y_i|$$

## 【0 1 0 3】

(9) マハラノビス距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、マハラノビス距離は次式(数 1 1)によって求められる。

## 【0 1 0 4】

## 【数 1 1】

$$\sum_{i=1}^n \sum_{j=1}^n s^{ij} (x_i - y_i)(x_j - y_j)$$

但し、 $s^{ij}$ は、分散共分散行列 $s^{ij}$ の逆行列の $i, j$ 要素とする。

## 【0 1 0 5】

(10) ミンコフスキー距離

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  に対して、ミンコフスキー距離は次式(数 1 2)によって求められる。

## 【0 1 0 6】

## 【数 1 2】

$$\left( \sum_{i=1}^n |x_i - y_i|^p \right)^{1/r}$$

ここで $r$ と $p$ は重みを調整するためのパラメータであり、 $r > 0$ の範囲で任意の値を設定可能である。 $p$ は個々のデータの差分に与える重みを調節するために用い、 $r$ は全体としての距離を調節するために用いる。

## 【0 1 0 7】

(11) 不一致割合

2つのデータ列  $x=\{x_i\}, y=\{y_i\}$ ,  $i=1, 2, \dots, n$  について、データ $x_i, y_i$ が不一致の場合が $c$ 回ある場合、 $x$ と $y$ の距離を以下のように定義することが出来る。(数 1 3)

(数 1 3)

$$c/n$$

また、 $x$ と $y$ の類似度を以下のように定義できる。(数 1 4)

(数 1 4)

$$(n-c) / n$$

## 【0 1 0 8】

本発明は以上説述したように、プロフィールデータを用いて、2つのソフトウェアの送信データを比較することに最大の特徴があり、これらの比較には上記に示した様々な相関係数その他、公知の任意の相関係数を使用することができる。

【図面の簡単な説明】

10

20

30

40

50

## 【 0 1 0 9 】

【図 1】本発明に係るデータ類似性検査装置の構成図である。

【図 2】本発明に係るデータ類似性検査方法の処理フローチャートである。

【図 3】本発明に係るマルウェアの特定システムの構成図である。

【図 4】本発明におけるマルウェア動的解析環境を用いたマルウェアの解析方法の説明図である。

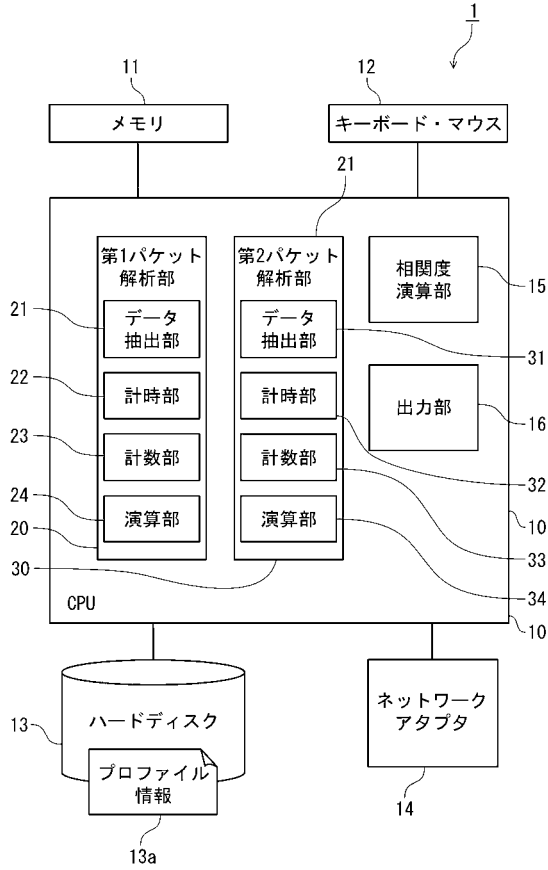
【図 5】本発明における相関解析部の別構成を示す構成図である。

## 【符号の説明】

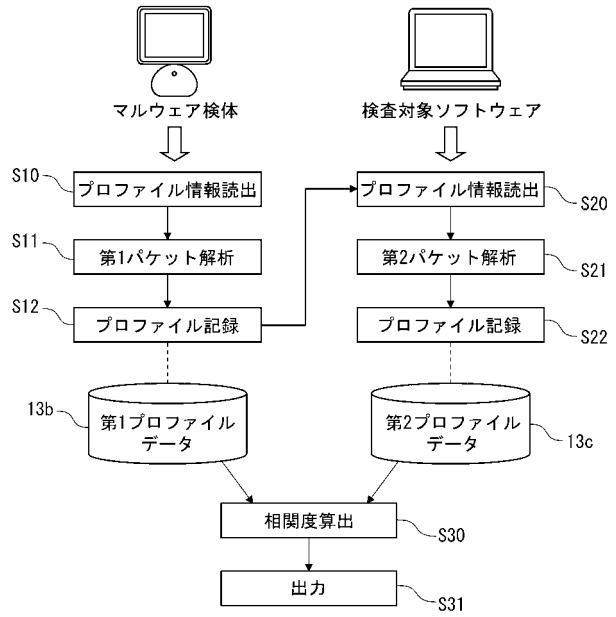
## 【 0 1 1 0 】

1	データ類似性検査装置	10
1 0	C P U	
1 1	メモリ	
1 2	キーボード・マウス	
1 3	ハードディスク	
1 3 a	プロファイル情報	
1 4	ネットワークアダプタ	
1 5	相関度演算部	
1 6	出力部	
2 0	第 1 パケット解析部	
2 1	データ抽出部	20
2 2	計時部	
2 3	計数部	
2 4	演算部	
3 0	第 2 パケット解析部	
3 1	データ抽出部	
3 2	計時部	
3 3	計数部	
3 4	演算部	

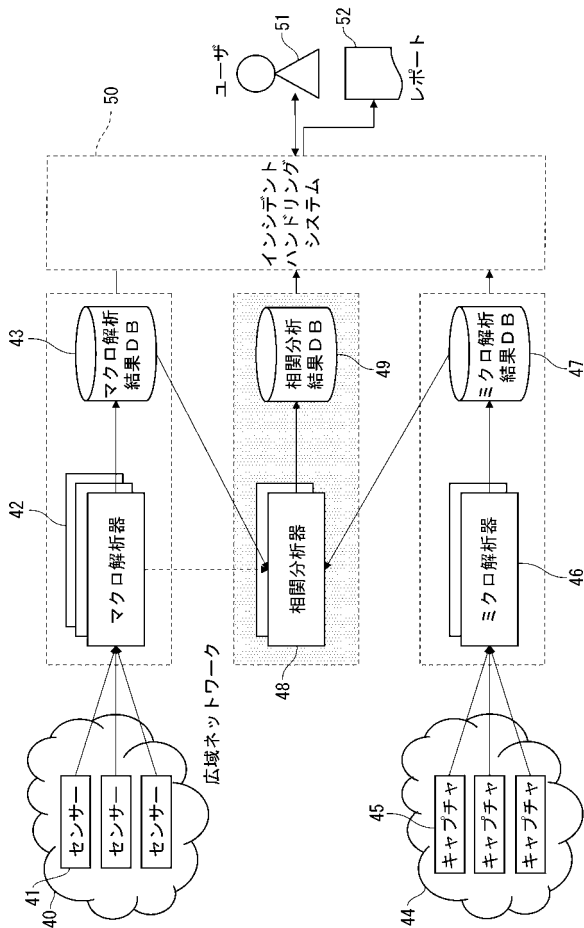
【図1】



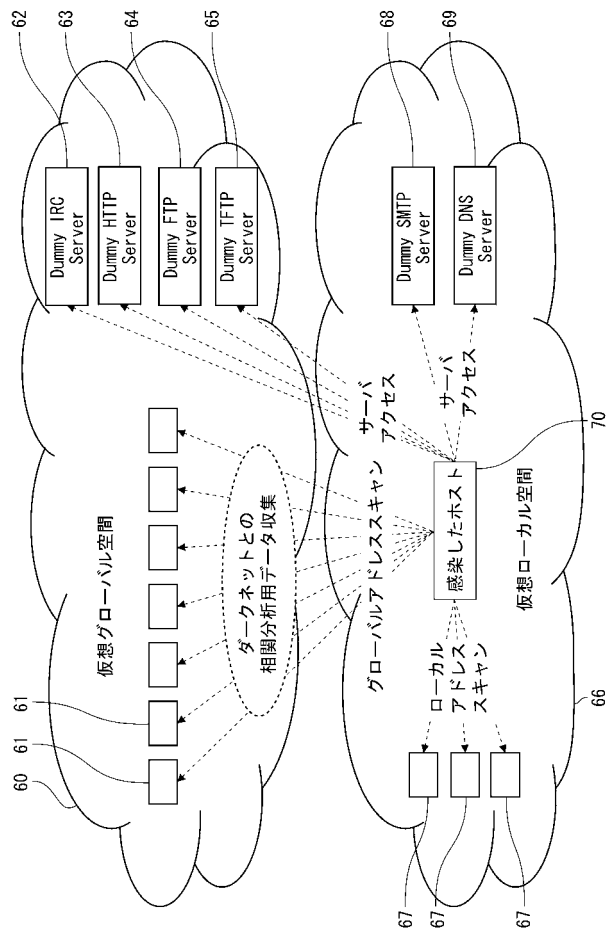
【図2】



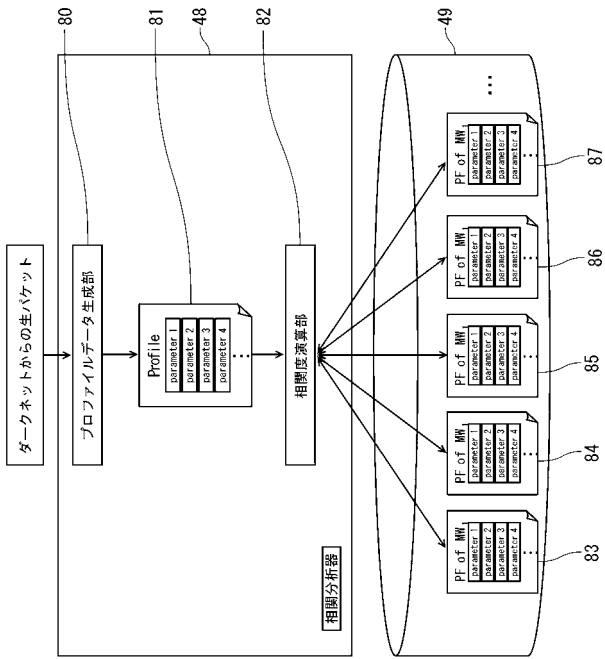
【図3】



【図4】



【 図 5 】



---

フロントページの続き

(72)発明者 井上 大介

東京都小金井市貫井北町4 - 2 - 1 独立行政法人情報通信研究機構内

Fターム(参考) 5B276 FD08

5B285 AA05 AA06 BA01 CA32 CA36 DA04

5K030 GA15 HA08 JA10 KA04 KA07 MA04 MB01 MC08 MD08