

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6778761号  
(P6778761)

(45) 発行日 令和2年11月4日(2020.11.4)

(24) 登録日 令和2年10月14日(2020.10.14)

(51) Int.Cl. F I  
**G06F 21/56 (2013.01)** G06F 21/56 320

請求項の数 18 (全 24 頁)

|                    |                               |           |                     |
|--------------------|-------------------------------|-----------|---------------------|
| (21) 出願番号          | 特願2018-552688 (P2018-552688)  | (73) 特許権者 | 504080663           |
| (86) (22) 出願日      | 平成29年4月6日(2017.4.6)           |           | エヌイーシー ラボラトリーズ アメリカ |
| (65) 公表番号          | 特表2019-514119 (P2019-514119A) |           | インク                 |
| (43) 公表日           | 令和1年5月30日(2019.5.30)          |           | NEC Laboratories Am |
| (86) 国際出願番号        | PCT/US2017/026359             |           | erica, Inc.         |
| (87) 国際公開番号        | W02017/177003                 |           | アメリカ合衆国 08540 ニュージャ |
| (87) 国際公開日         | 平成29年10月12日(2017.10.12)       |           | ージー州 プリンストン スイート 20 |
| 審査請求日              | 平成30年11月20日(2018.11.20)       |           | 0 インディペンデンス ウェイ 4   |
| (31) 優先権主張番号       | 62/318,844                    | (74) 代理人  | 100123788           |
| (32) 優先日           | 平成28年4月6日(2016.4.6)           |           | 弁理士 宮崎 昭夫           |
| (33) 優先権主張国・地域又は機関 | 米国 (US)                       | (74) 代理人  | 100127454           |
|                    |                               |           | 弁理士 緒方 雅昭           |
| (31) 優先権主張番号       | 15/479,928                    |           |                     |
| (32) 優先日           | 平成29年4月5日(2017.4.5)           |           |                     |
| (33) 優先権主張国・地域又は機関 | 米国 (US)                       |           |                     |

最終頁に続く

(54) 【発明の名称】 ハイブリッドプログラムバイナリ特徴の抽出及び比較

(57) 【特許請求の範囲】

【請求項 1】

プログラムバイナリの類似度を特定するための方法であって、  
 1 つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成することであって、前記ハイブリッド特徴が、参照特徴、リソース特徴、抽象化制御フロー特徴及び構造特徴を含むことと、  
 前記抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成することと、  
 前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、前記バイナリ対毎の類似度スコアを決定することと、  
 入力されるハイブリッド特徴パラメータと、前記バイナリ対毎の前記類似度スコアとに基づいて、前記バイナリ対における 2 つの前記プログラムバイナリの類似度を表すハイブリッド差異スコアを生成することと、  
 既知のマルウェアのプログラムバイナリ及び前記入力プログラムバイナリ間の前記ハイブリッド差異スコアに基づいて、前記入力プログラムバイナリに関する前記マルウェアの尤度を特定することと、  
 を含む、方法。

【請求項 2】

前記複数のバイナリ対は、想定される全ての組み合わせのバイナリ対を含む、請求項 1 に記載の方法。

## 【請求項 3】

所定の閾値ハイブリッド差異スコアに到達している場合に、マルウェア対策ソフトウェアにおけるマルウェア定義ライブラリを更新することをさらに含む、請求項 1 に記載の方法。

## 【請求項 4】

前記参照特徴は、入力プログラムから参照されるプログラム及び参照される関数のリストのうちの少なくとも 1 つを含む、請求項 1 に記載の方法。

## 【請求項 5】

前記リソース特徴は、グローバルデータ、プログラムメタデータ、プログラムアイコン、文字列及びデバッグシンボルのうちの少なくとも 1 つを含む、請求項 1 に記載の方法。

10

## 【請求項 6】

入力プログラムを逆アセンブリし、命令が制御依存である場合、オペコードのみを取り出し、前記抽象化制御フロー特徴に含める、請求項 1 に記載の方法。

## 【請求項 7】

前記構造特徴は、バイナリセクションの名称及びバイナリセクションのサイズのリストのうちの少なくとも 1 つを含む、請求項 1 に記載の方法。

## 【請求項 8】

前記ハイブリッド特徴パラメータは、前記ハイブリッド差異スコアを生成するときの前記ハイブリッド特徴毎の貢献量を示す所定のレートのセットである、請求項 1 に記載の方法。

20

## 【請求項 9】

プログラムバイナリの類似度を特定するためのシステムであって、メモリと接続されたプロセッサを備え、前記プロセッサは、  
1 つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成することであって、前記ハイブリッド特徴が、参照特徴、リソース特徴、抽象化制御フロー特徴及び構造特徴を含み、  
前記抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成し、  
前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、前記バイナリ対毎の類似度スコアを決定し、

入力されるハイブリッド特徴パラメータと、前記バイナリ対毎の前記類似度スコアとに基づいて、前記バイナリ対における 2 つの前記プログラムバイナリの類似度を表すハイブリッド差異スコアを生成し、

30

既知のマルウェアのプログラムバイナリ及び前記入力プログラムバイナリ間の前記ハイブリッド差異スコアに基づいて前記入力プログラムバイナリに関する前記マルウェアの尤度を特定するように構成された、システム。

## 【請求項 10】

前記複数のバイナリ対は、想定される全ての組み合わせのバイナリ対を含む、請求項 9 に記載のシステム。

## 【請求項 11】

前記プロセッサは、所定の閾値ハイブリッド差異スコアに到達している場合に、マルウェア対策ソフトウェアにおけるマルウェア定義ライブラリを更新するようにさらに構成された、請求項 9 に記載のシステム。

40

## 【請求項 12】

前記参照特徴は、入力プログラムから参照されるプログラム及び参照される関数のリストのうちの少なくとも 1 つを含む、請求項 9 に記載のシステム。

## 【請求項 13】

前記リソース特徴は、グローバルデータ、プログラムメタデータ、プログラムアイコン、文字列及びデバッグシンボルのうちの少なくとも 1 つを含む、請求項 9 に記載のシステム。

## 【請求項 14】

50

前記プロセッサは、入力プログラムを逆アセンブリし、命令が制御依存である場合、オペコードのみを取り出し、前記抽象化制御フロー特徴に含める、請求項 9 に記載のシステム。

【請求項 15】

前記構造特徴は、バイナリセクションの名称及びバイナリセクションのサイズのリストのうち少なくとも 1 つを含む、請求項 9 に記載のシステム。

【請求項 16】

前記ハイブリッド特徴パラメータは、前記ハイブリッド差異スコアを生成するときの前記ハイブリッド特徴毎の貢献量を示す所定のレートのセットである、請求項 9 に記載のシステム。

10

【請求項 17】

プログラムバイナリの類似度を特定するためのコンピュータで読み取り可能なプログラムを備えた非一時的なコンピュータで読み取り可能な記憶媒体であって、

前記コンピュータで読み取り可能なプログラムは、

1 つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成する工程であって、前記ハイブリッド特徴が、参照特徴、リソース特徴、抽象化制御フロー特徴及び構造特徴を含み、

前記抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成する工程と、

前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、前記バイナリ対毎の類似度スコアを決定する工程と、

20

入力されるハイブリッド特徴パラメータと、前記バイナリ対毎の前記類似度スコアとに基づいて、前記バイナリ対における 2 つの前記プログラムバイナリの類似度を表すハイブリッド差異スコアを生成する工程と、

既知のマルウェアのプログラムバイナリ及び前記入力プログラムバイナリ間の前記ハイブリッド差異スコアに基づいて前記入力プログラムバイナリに関する前記マルウェアの尤度を特定する工程と、

を前記コンピュータに実行させる、非一時的なコンピュータで読み取り可能な記憶媒体。

【請求項 18】

所定の閾値ハイブリッド差異スコアに到達している場合に、マルウェア対策ソフトウェアにおけるマルウェア定義ライブラリを更新する工程をさらに含む、請求項 17 に記載の非一時的なコンピュータで読み取り可能な記憶媒体。

30

【発明の詳細な説明】

【技術分野】

【0001】

この出願は、2016年4月6日に出願された米国仮特許出願第62/318,844号を基礎とする優先権を主張し、その開示の全てをここに取り込む。

【0002】

本発明は、プログラム特徴(program feature)の抽出及び比較に関し、特にハイブリッドプログラムバイナリ特徴の抽出及び比較による、悪意のあるソフトウェア攻撃の検出及び防止に関する。

40

【背景技術】

【0003】

プログラムバイナリは、プログラムの特性を理解するためのサイバーセキュリティの重要な態様である。無害な(benign)ソフトウェア及びマルウェアはプログラムバイナリとして配信される。それらの配信及びランタイム動作を調査することは、ウィルス対策ソフトウェア等の多くのサイバーセキュリティソリューションで実行される重要なタスクである。

【0004】

従来のウィルス対策製品は、バイナリシグネチャを用いてマルウェアを特定し、プログ

50

ラムバイナリを一意に識別するためにハッシュ値が一般的に用いられてきた。これは、コンテンツ（例えば、バイナリ全体、あるいは特定の1つまたは複数のセクション）に基づく1つの特徴として分類できる。しかしながら、バイナリコンテンツのハッシュ値は、小さな変化に対して敏感であり、バイナリにおける単一ビットの差異であっても、完全に異なるハッシュ値が生じるため、そのようなシステム及び方法を用いると、正確または信頼できる結果が得られない。さらに、バイナリコンテンツにおける差異量は、ハッシュ値の差異ではうまく表せない。

#### 【0005】

従来の他の手法では、マルウェアファミリーの類似度を判定する試みにおいて、プログラムの制御フロー情報（例えば、CPU命令、システムコール）が用いられてきた。しかしながら、これらの手法は、少なくとも部分的に、マルウェアだけでなく無害なソフトウェアにおける感度に起因して、類似度比較において有効でない。したがって、当該方法は、プログラムの類似度を正確に判定するのに有効ではない。

#### 【0006】

無害なソフトウェアは、異なるプラットフォーム及びパッチのために多くのバージョンを有する。さらに、それらのソースコードが非常に類似している場合でも、一旦、コンパイルされてバイナリフォーマットになると、その命令構造はコンパイラのアルゴリズム及び最適化に起因して大幅に異なるものになる。さらに、悪意のあるソフトウェア（マルウェア）に関して、マルウェアの作成者は、コードが変形したもの（例えば、ポリモーフィックマルウェアコード）を用いる。これは事実上、上述したような従来の手法を混乱させて不正確で無効なものにする。このため、従来の手法は、例えば類似する無害なプログラムとマルウェアファミリーとを正確に判定するには十分に信頼できるものではなく、有効でもないため、プログラムバイナリの確実で有効な特性解析及び類似度比較は未解決の課題である。

#### 【発明の概要】

#### 【0007】

本原理の一態様によれば、1つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成することを含む、プログラムバイナリの類似度を特定するための方法が提供される。ハイブリッド特徴は、参照特徴（reference feature）、リソース特徴（resource feature）、抽象化制御フロー特徴（abstract control flow feature）及び構造特徴（structural feature）を含む。抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成し、前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、バイナリ対毎の類似度スコアを決定する。入力されるハイブリッド特徴パラメータと、バイナリ対毎の類似度スコアとに基づいてバイナリ対における2つのプログラムバイナリの類似度を表すハイブリッド差異スコアを生成する。既知のマルウェアのプログラムバイナリ及び入力プログラムバイナリ間のハイブリッド差異スコアに基づいて入力プログラムに関するマルウェアの尤度を特定する。

#### 【0008】

本原理の他の態様によれば、プログラムバイナリの類似度を特定するためのシステムが提供される。システムは、1つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成するように構成された、メモリと接続されたプロセッサを備える。ハイブリッド特徴は、参照特徴、リソース特徴、抽象化制御フロー特徴及び構造特徴を含む。抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成し、前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、バイナリ対毎の類似度スコアを決定する。入力されるハイブリッド特徴パラメータと、バイナリ対毎の類似度スコアとに基づいてバイナリ対における2つのプログラムバイナリの類似度を表すハイブリッド差異スコアを生成する。既知のマルウェアのプログラムバイナリ及び入力プログラムバイナリ間のハイブリッド差異スコアに基づいて入力プログラムバイナリに関するマルウェアの尤度を特定する。

#### 【0009】

本原理の他の態様によれば、1つまたは複数の入力プログラムバイナリからプログラムバイナリ特徴を抽出し、対応するハイブリッド特徴を生成することを含む、プログラムバイナリの類似度を特定するための非一時的なコンピュータで読み取り可能な媒体が提供される。ハイブリッド特徴は、参照特徴、リソース特徴、抽象化制御フロー特徴及び構造特徴を含む。抽出されたハイブリッド特徴を用いた複数のバイナリ対の組み合わせを生成し、前記ハイブリッド特徴のうち同じ種類の特徴どうしを比較することで、バイナリ対毎の類似度スコアを決定する。入力されるハイブリッド特徴パラメータと、バイナリ対毎の類似度スコアとに基づいて**バイナリ対における2つのプログラムバイナリの類似度を表すハイブリッド差異スコアを生成する。既知のマルウェアのプログラムバイナリ及び入力プログラムバイナリ間のハイブリッド差異スコアに基づいて入力プログラムに関するマルウェアの尤度を特定する。**

10

【0010】

これら及び他の特徴並びに利点は、以下の典型的な実施形態の詳細な説明を添付の図面と併せて読むことで明らかになるであろう。

【図面の簡単な説明】

【0011】

本開示では、後述するように、以下の図面を参照しながら好ましい実施形態について詳細に説明する。

【0012】

【図1】図1は、本原理による、本原理を適用できる典型的な処理システムを示すブロック/フロー図である。

20

【0013】

【図2A】図2Aは、本原理による、プログラムバイナリ特徴抽出のための高レベルのシステム/方法を示すブロック/フロー図である。

【0014】

【図2B】図2Bは、本原理による、ハイブリッド特徴類似度解析のための高レベルのシステム/方法を示すブロック/フロー図である。

【0015】

【図3】図3は、本原理による、プログラムバイナリ特徴抽出のための方法を示すブロック/フロー図である。

30

【0016】

【図4】図4は、本原理による、プログラムバイナリ特徴抽出のための典型的な参照特徴を示すブロック/フロー図である。

【0017】

【図5】図5は、本原理による、抽象化制御特徴の生成のための方法を示すブロック/フロー図である。

【0018】

【図6】図6は、本原理による、ハイブリッド特徴類似度比較のための方法を示すブロック/フロー図である。

【0019】

【図7】図7は、本原理による、2つのバイナリの類似度比較のための方法を示すブロック/フロー図である。

40

【0020】

【図8】図8は、本原理による、ハイブリッドプログラムバイナリ特徴の抽出及び比較のためのシステムを示すブロック/フロー図である。

【発明を実施するための形態】

【0021】

本原理によれば、ハイブリッドバイナリプログラム特徴の抽出及び比較のためのシステム及び方法が提供される。

【0022】

50

特に有用な実施形態において、本原理によるハイブリッドバイナリプログラム特徴の抽出及び比較により、1つまたは複数のコンピュータシステムに対する悪意のあるソフトウェア（マルウェア）攻撃を検出及び/または防止するためのシステム及び方法が提供される。

#### 【0023】

一実施形態において、本原理は、例えばウイルス対策保護システムのための総合的なウイルス定義アップデートとして、マルウェア攻撃から1つまたは複数のコンピューティングシステムを保護する実用的な解決策として用いることができる。プログラムバイナリは、プログラムの特性を理解するためのサイバーセキュリティの重要な態様である。無害なソフトウェア及びマルウェアは、プログラムバイナリとして配信される。それらの配信及びランタイム動作を調査することは、多くのサイバーセキュリティソリューション（例えば、ウイルス対策ソフトウェア）で実行される重要なタスクである。

10

#### 【0024】

現代のコンピューティングシステムは、多くのプログラムを使用し、コンピュータシステムの複雑性は、しばしば複数の機能を実装するために、一緒に動作させる複数のプログラムが必要であることがよく知られている。さらに、例えば多様なオペレーティングシステム及びプラットフォームに起因して、同じソフトウェアが多くの異なるバージョンとしてパッケージ化されて配信されることがある。各オペレーティングシステムは、複数のバージョン（例えば、WindowsにおけるService Pack、Linuxにおける配信バージョン等）を有することがある。さらに、プログラムは、多数のエラー及びセキュリティ脆弱性に起因して、頻繁に再コンパイルされる、またはパッチが適用される。例えば、現代のオペレーティングシステムは、複数の理由のうちのいずれかで、プログラムの新たな配信バージョン及び更新バージョンがしばしば発行される。

20

#### 【0025】

これらの理由により、企業環境において日々多数のプログラムバイナリが配備及び更新される。これは、サイバーセキュリティシステム及び方法にとって、プログラムの監視及び解析が困難な課題になる。ウイルス対策会社は、新たなマルウェアを解析し、それらのバイナリシグネチャデータベースを更新するために多くの労力及びリソースを集中させている。しかしながら、バイナリベースのシグネチャを解析して分類する既存の手法は、例えばそれらの低い信頼性及びバイナリの規模に起因して、多数のプログラムバイナリの取

30

#### 【0026】

上述したように、従来のシステム及び方法は、プログラム間の類似度または無害なプログラムとマルウェアファミリーとの正確な類似度を有効にまたは確実に判定しない。無害なソフトウェアは、異なるプラットフォーム及びパッチのために多くのバージョンを有する。それらのソースコードが非常に類似している場合でも、一旦、コンパイルされてバイナリフォーマットになると、その命令構造は、例えばコンパイラのアルゴリズム及び最適化に起因して大幅に異なるものになる。

#### 【0027】

さらに、悪意のあるソフトウェア（マルウェア）に関して、マルウェアの作成者は、コードが変形したもの（例えば、ポリモーフィックマルウェアコード）を用いる。これは事実上、上述したような従来の手法を混乱させて不正確で無効なものにする。このため、上述したような従来の手法は、例えば類似する無害なプログラムとマルウェアファミリーとを正確にかつ確実に判定するために、十分に信頼できるものではなく、効果があるものでも無いため、プログラムバイナリの確実に有効な特性解析及び類似度比較は未解決の課題である。

40

#### 【0028】

様々な実施形態によれば、本原理は、プログラムバイナリから複数の特徴を抽出し、プログラムの特性を数値化して、それらの類似度をブラックボックス手法で（例えば、ソースコードまたはデバッグ情報を何も用いず）比較することに適用できる。特に、本原理に

50

よれば、抽出される複数の特徴は、(1)参照特徴、(2)リソース特徴、(3)抽象化制御フロー特徴及び/または(4)構造特徴を含む。これらの特徴は、参照されるバイナリ、リソース、制御フロー及びバイナリ構造の観点から見て、バイナリの複数の態様を表している。さらに、これらの特徴は、バイナリの他の特徴よりも関連する特性の適用範囲(coverage)がより豊富であり、相補的にプログラムの類似度を数値化するためにより有効である。

【0029】

いくつかの実施形態において、本原理は、サイバーセキュリティシステムにおけるプログラムのホワイトリストへの登録及びマルウェアのクラスタリングのための重要な特徴である、プログラムバイナリの有効で正確な比較を可能にするために適用できる。

10

【0030】

プログラムのホワイトリストへの登録に関して、本原理による既知の無害なプログラムとの類似度を判定することは、無害なプログラムの多様な変形に対する、複雑度を低減して、マルウェア検出の正確性及び有効性を増大させるために有用である。以下でさらに詳細に示すように、本発明の複数の特徴は、類似度計算のプロセスを改善する。

【0031】

マルウェアのクラスタリングに関して、大量の新たなマルウェアが毎日のように新たにリリースされ、発見されている。このマルウェアの全てを手作業で調べることは、人力のコストが高く、人力及び/または従来のウィルス対策システムを用いるタイムリーな方法(例えば、絶え間なく変化するマルウェアに対処し、マルウェア変形からの攻撃をリアルタイムで防止するために十分にタイムリーに)では、膨大な量の潜在的なマルウェアを処理するには能力が欠如しており、不的確であるため、非常に困難である。本発明は、以下でさらに詳細に説明する様々な実施形態による、マルウェアの比較及びクラスタリングの品質、有効性及び正確性を改善するのに有利である。

20

【0032】

本明細書に記載される実施形態は、全てハードウェアで実現してもよく、全てソフトウェアで実現してもよく、ハードウェアとソフトウェアの両方の要素を含んでもよい。好ましい実施形態において、本発明は、ファームウェア、常駐ソフトウェア、マイクロコード等を含むが、これらに限定されないソフトウェアでも実現可能である。

【0033】

実施形態には、コンピュータもしくは任意の命令実行システムによって使用される、または関連して使用されるプログラムコードを提供する、コンピュータで利用可能な、またはコンピュータで読み取り可能な媒体からアクセスできる、コンピュータプログラム製品を含んでもよい。コンピュータで利用可能な、またはコンピュータで読み取り可能な媒体には、命令実行システム、機器、もしくは装置によって使用される、または関連して使用されるプログラムを格納、伝達、伝搬または転送する任意の機器を含んでもよい。該媒体は、磁気媒体、光学媒体、電子媒体、電磁気媒体、赤外線媒体、または半導体システム(または機器もしくは装置)、あるいは伝搬媒体であってよい。該媒体には、半導体または固体メモリ、磁気テープ、取り外し可能なコンピュータディスク、ランダムアクセスメモリ(RAM)、リードオンリメモリ(ROM)、リジッド磁気ディスク及び光ディスク等のコンピュータで読み取り可能な媒体を含んでもよい。

30

40

【0034】

各コンピュータプログラムは、汎用または特別な目的を持つプログラム可能なコンピュータで読み取ることができる、機械で読み取り可能なストレージメディアまたは装置(例えば、プログラムメモリまたは磁気ディスク)に格納される。該コンピュータプログラムは、ストレージメディアまたは装置から本明細書に記載された手順を実行するコンピュータで読み出される、該コンピュータの設定及び制御動作のためのものである。本発明のシステムには、本明細書に記載した機能を実行する、特定の及び事前に定義された方法でコンピュータを動作させるように構成されたコンピュータプログラムを含む、コンピュータで読み取り可能なストレージメディアも考慮される。

50

## 【0035】

プログラムコードを記憶及び/または実行するのに適したデータ処理システムは、システムバスを介してメモリ要素に直接または間接的に接続された少なくとも1つのプロセッサを備えていてもよい。このメモリ要素には、処理の実行中にバルク記憶装置からコードが検索される回数を減らすために、プログラムコードの実際の実行中に用いられるローカルメモリ、バルク記憶装置及び少なくともいくつかのプログラムコードを一時的に記憶するキャッシュメモリを備えていてもよい。入力/出力すなわちI/O装置(限定されるものではないが、キーボード、ディスプレイ、ポインティング装置等を含む)は、直接またはI/Oコントローラを介してシステムに接続されてもよい。

## 【0036】

ネットワークアダプタは、データ処理システムが、プライベートネットワークまたは公衆ネットワークを介して、他のデータ処理システムまたは遠隔プリンタもしくは記憶装置に接続されることを可能にするために、上記システムと接続されていてもよい。モデム、ケーブルモデム及びイーサネット(登録商標)カードは、現在利用可能なタイプのネットワークアダプタのほんの一握りのものである。

## 【0037】

ここで、同じ数字が同一または同様の要素を表す図面、まず図1を詳細に参照すると、図1には本原理の一実施形態による、本原理が適用できる典型的な処理システム100が示されている。

## 【0038】

処理システム100は、システムバス102を介して他のコンポーネントと動作可能に接続された、少なくとも1つのプロセッサ(CPU)104を含む。システムバス102には、キャッシュ106、リードオンリメモリ(ROM)108、ランダムアクセスメモリ(RAM)110、入力/出力(I/O)アダプタ120、サウンドアダプタ130、ネットワークアダプタ140、ユーザインタフェースアダプタ150及びディスプレイアダプタ160が動作可能に接続されている。

## 【0039】

第1の記憶装置122及び第2の記憶装置124は、I/Oアダプタ120によってシステムバス102に動作可能に接続されている。記憶装置122及び124は、ディスク記憶装置(例えば、磁気ディスク記憶装置または光ディスク記憶装置)、固体磁気装置等のいずれであってもよい。記憶装置122及び124は、同じタイプの記憶装置であってもよく、異なるタイプの記憶装置であってもよい。

## 【0040】

スピーカ132は、サウンドアダプタ130によってシステムバス102に動作可能に接続されている。トランシーバ142は、ネットワークアダプタ140によってシステムバス102に動作可能に接続されている。ディスプレイ装置162は、ディスプレイアダプタ160によってシステムバス102に動作可能に接続されている。

## 【0041】

第1のユーザ入力装置152、第2のユーザ入力装置154及び第3のユーザ入力装置156は、ユーザインタフェースアダプタ150によってシステムバス102に動作可能に接続されている。ユーザ入力装置152、154及び156は、キーボード、マウス、キーパッド、イメージキャプチャ装置、モーション感知装置、マイクロフォン、あるいはこれらの装置のうちの少なくとも2つの装置の機能を組み込んだ装置等のいずれであってもよい。本原理の趣旨を維持する限りにおいて、他のタイプの入力装置を使用することも可能である。ユーザ入力装置152、154及び156は、同じタイプのユーザ入力装置であってもよく、異なるタイプのユーザ入力装置であってもよい。ユーザ入力装置152、154及び156は、システム100に情報を入力し、システム100から情報を入力するために使用される。

## 【0042】

処理システム100は、当業者であれば容易に思いつくような他の要素(図示せず)を

10

20

30

40

50

含んでいてもよく、特定の要素を省略することも可能である。例えば、当業者であれば容易に理解できるが、処理システム100には、その詳細な実装に応じて他の様々な入力装置及び/または出力装置を使用できる。例えば、無線及び/または有線による様々なタイプの入力装置及び/または出力装置を使用できる。さらに、当業者であれば容易に理解できるが、様々な構成において追加のプロセッサ、コントローラ、メモリ等を使用することも可能である。処理システム100の上記及び他の変形例は、本明細書で提供される本原理の教示によって当業者であれば容易に考えられるであろう。

【0043】

さらに、図1、図2A、図2B及び図8に関して説明するシステム100、200、210及び800は、本原理の各実施形態を実現するためのシステムであることを理解されたい。本原理の様々な実施形態によれば、処理システム100の一部または全ては、システム200、210及び800の要素の1つまたは複数を実行してもよい。

10

さらに、処理システム100は、例えば、図2A、図2B、図3、図4、図5、図6及び図7の方法200、210、300、400、500、600及び700の少なくとも一部を含む、本明細書で説明する方法の少なくとも一部を実行してもよいことを理解されたい。同様に、本原理の様々な実施形態によれば、システム800の一部または全ては、図2A、図2B、図3、図4、図5、図6及び図7の方法200、210、300、400、500、600及び700の少なくとも一部を実行するために使用されてもよい。

【0044】

次に図2Aを参照すると、図2Aには本原理の実施形態によるプログラムバイナリ特徴抽出のための高レベルの方法200が例示的に示されている。

20

【0045】

一実施形態において、ブロック202にて、1つまたは複数のプログラムバイナリが入力される。本原理によれば、ブロック204において、ハイブリッドバイナリ特徴を生成するためにプログラムバイナリ特徴抽出が実行され、ブロック206において、ハイブリッドバイナリ特徴が出力される。

【0046】

次に図2Bを参照すると、図2Bには本原理の実施形態によるハイブリッド特徴類似度解析のための高レベルの方法210が例示的に示されている。

【0047】

30

一実施形態において、ブロック212、214及び/または216において、1つまたは複数の生成されたハイブリッドバイナリ特徴が入力され、ブロック218において、1つまたは複数のハイブリッド特徴パラメータが入力される。本原理によれば、ブロック220において、2つ以上のプログラムバイナリの類似度を判定するために、特徴212、214及び216、並びに特徴パラメータ218を用いることができる。ブロック222において、ハイブリッド特徴類似度解析220に基づく類似度ベクトルが出力される。様々な実施形態によれば、ブロック222で出力される類似度ベクトルは、例えばマルウェア攻撃の検出及び防止のために、リアルタイムな悪意のあるソフトウェア(マルウェア)定義の比較及び更新を提供するために用いることができる。

【0048】

40

次に図3を参照すると、図3には本原理の実施形態によるプログラムバイナリ特徴抽出のための方法300が例示的に示されている。

【0049】

一実施形態において、本原理によれば、ブロック302において、1つまたは複数のプログラムバイナリが入力され、ブロック304において、プログラムバイナリ特徴が抽出される。様々な実施形態によれば、プログラムバイナリ314で抽出されるハイブリッドバイナリ特徴には、参照特徴316、リソース特徴318、抽象化制御フロー特徴320及び構造特徴322のうちの1つまたは複数が含まれる。本原理によれば、参照特徴抽出306、リソース特徴抽出308、抽象化制御フロー特徴抽出310及び構造特徴抽出312を含む、対応する抽出機能を用いてハイブリッド特徴が抽出される。様々な実施形態

50

によれば、出力は、個々の特徴としてハイブリッド特徴で保存される。

【0050】

説明を容易にするため、本システム及び方法の様々な特徴に関して、以下の表現を用いる。

P: 入力プログラムバイナリ 302

F\_F(P): プログラムバイナリ P の参照特徴 316

F\_R(P): プログラムバイナリ P のリソース特徴 318

F\_C(P): プログラムバイナリ P の抽象化制御特徴 320

F\_S(P): プログラムバイナリ P の構造特徴 322

【0051】

—実施形態において、プログラム P のためのハイブリッドバイナリ特徴 314 は、以下で示すように、プログラム P のための参照特徴、リソース特徴、抽象化制御特徴及び構造特徴の 4 タプルである。

$HF(P) = (F\_F(P), F\_R(P), F\_C(P), F\_S(P))$

【0052】

ハイブリッド特徴の類似度は、以下のように表せばよい。

【0053】

プログラム P1、P2 の参照特徴 316 の類似度:

$Sim\_F(F\_F(P1), F\_F(P2))$

その入力、プログラム P1 及び P2 の 2 つの参照特徴である。

この値は、0 (0%) と 1 (100%) との間である。

【0054】

プログラム P1、P2 のリソース特徴 318 の類似度:

$Sim\_R(F\_R(P1), F\_R(P2))$

その入力、プログラム P1 及び P2 の 2 つのリソース特徴である。

この値は、0 (0%) と 1 (100%) との間である。

【0055】

プログラム P1、P2 の抽象化制御フロー特徴 320 の類似度:

$Sim\_C(F\_C(P1), F\_C(P2))$

その入力、プログラム P1 及び P2 の 2 つの抽象化制御フロー特徴である。

この値は、0 (0%) と 1 (100%) との間である。

【0056】

プログラム P1、P2 の構造特徴 322 の類似度:

$Sim\_S(F\_S(P1), F\_S(P2))$

その入力、プログラム P1 及び P2 の 2 つの構造特徴である。

この値は、0 (0%) と 1 (100%) との間である。

【0057】

様々な実施形態によれば、プログラムバイナリ特徴抽出 304 は、以下のように実行すればよい。

【数 1】

Given the following:

P: an input program binary

F\_F(P): Reference feature for a program P

F\_R(P): Resource feature for a program P

F\_C(P): Abstract control feature for a program P

F\_S(P): Structural feature for a program P

10

20

30

40

50

## 【 0 0 5 8 】

－実施形態において、参照特徴抽出 3 0 6 は、以下のように実行すればよい。

## 【 数 2 】

Extract\_F\_F(P) // Extract Reference Feature from a Program binary P

ListOfReferences = [] // empty set

Sections = Get the list of binary sections of Program P

For each section in Sections:

    If section is a reference table (For example, import Table, GOT in ELF) 10

        For each entry in section:

            program = getReferredProgram(entry)

            function = getReferredFunctionName(entry)

            ListOfReferences.add(pair(program, function))

Return ListOfReferences // This is F\_F(P)

## 【 0 0 5 9 】

－実施形態において、リソース特徴抽出 3 0 8 は、以下のように実行すればよい。

## 【 数 3 】

Extract\_F\_R(P) // Extract Resource Feature from a Program binary P

ListOfResource = [] // empty set

Sections = Get the list of binary sections of Program P

For each section in Sections:

    If section has resource (E. g., String, Symbol, Global Data, Icon, etc.) 30

        kind = getKind(section)

        For each entry in section:

            value = getValue(entry)

            ListOfResource.add(pair(kind, value))

Return ListOfResource // This is F\_R(P)

## 【 0 0 6 0 】

－実施形態において、抽象化制御特徴抽出 3 1 0 は、以下のように実行すればよい。

40

## 【数4】

```
Extract_F_C(P) // Extract Abstract Control Flow Feature from a Program binary P
```

```
ListOfACF = [] // empty set
```

```
Sections = Get the list of binary sections of Program P
```

```
For each section in Sections:
```

```
    If section is a code section
```

```
        DisassembledInstructions = Disassemble(section)
```

```
        For each instruction in DisassembledInstructions:
```

```
            If instruction is control dependent
```

```
                ACF = getOpCode(instruction)
```

```
                ListOfACF.add(ACF)
```

```
Return ListOfACF // This is F_C(P)
```

10

20

## 【0061】

一実施形態において、構造特徴抽出312は、以下のように実行すればよい。

## 【数5】

```
Extract_F_S(P) // Extract Structural Feature from a Program binary P
```

```
ListOfSections = [] // empty set
```

```
Sections = Get the list of binary sections of Program P
```

```
For each section in Sections:
```

```
    name = getSectionName(section)
```

```
    size = getSectionSize(section)
```

```
    ListOfSections.add(pair(name, size))
```

```
Return ListOfSections // This is F_S(P)
```

30

## 【0062】

本原理によれば、上記特徴の4タプルは、少なくとも参照特徴がライブラリとの関係を示し、リソース特徴が同じプログラムの異なる複数のバージョンまたは類似する複数のプログラムで共有できる共通の文字列またはメッセージを示し、抽象化制御フロー特徴がより少ないノイズを有するプログラム命令の類似度を示し、構造特徴がセクション情報を比較することでプログラム構造全体の類似度を示すという理由で用いられるため、例えば本原理によるマルウェア攻撃を防止するためのプログラム間の類似度の検出に有用である。

40

## 【0063】

次に図4を参照すると、図4には本原理によるプログラムバイナリ特徴抽出のための典型的な参照特徴のダイアグラム400が例示的に示されている。

## 【0064】

ほとんどのプログラムは、個別のプログラムバイナリに格納されるライブラリコードを用いる。この特徴は、関連するプログラムバイナリに対する機能的な依存性を示している

50

。プログラムは、他のどのバイナリが、より正確に言えば、どの機能を見つけて、適切に動作するように結びつけるべきかを記述するバイナリセクションを有する。

【0065】

例えば、Linuxプラットフォームで普及しているELF (Executable and Linkable Format) バイナリフォーマットにおいて、この情報を収集するために、インポートテーブル及びGOT (Global Offset Table) を用いることができる。他のバイナリフォーマットは、類似のバイナリセクションまたはテーブルを有する。

【0066】

典型的な一実施形態において、プログラムA 402は、プログラムB 404の関数B1 405とプログラムC 406の関数C1 407とを用いる。この典型的なプログラムからの参照特徴は、ブロック410において、より詳細に示される。

10

【0067】

再び図3を参照すると、典型的なリソース特徴318が以下の表1で示されている。

【0068】

表1：リソース特徴

【表1】

リソース特徴

| 種類       | 値             |
|----------|---------------|
| 文字列      | 「バージョン 1.0.1」 |
| 文字列      | 使用:プログラム      |
| シンボル     | 関数 A          |
| グローバルデータ | 192.168.1.0   |
| グローバルデータ | 255           |
| アイコン     | バイナリデータ       |

20

30

【0069】

プログラムは、プログラムバイナリに埋め込まれた様々なリソースを使用する。これらのリソースのうちいくつかは、プログラムで用いられるデータである。例えば、グローバルデータ、プログラムメタデータ、プログラムアイコン、文字列 (strings)、デバッグシンボル等は、リソース特徴318のこのカテゴリに属する。そのような情報は、通常、個別のバイナリセクションで保存される。例えば、ELFバイナリフォーマットにおいて、読み出し専用データセクション及びシンボルテーブルセクションがそのような情報のために用いられる。

【0070】

表1で示したように、リソース特徴のこの例において、このテーブルは、種類 (Kind) の列を有し、データは実際に測定及び/または受信した値 (Value) である。この例は、いくつかの文字列と、プログラム関数シンボルと、読み出し専用セクションからのグローバルデータと、メタデータに属するアイコンデータとを示している。

40

【0071】

次に図5を参照すると、図5には本原理による抽象化制御フロー特徴の生成及び抽出のための方法500が例示的に示されている。

【0072】

一実施形態において、与えられたプログラムが、ブロック502において、逆アセンブルされ、アルゴリズムは各命令を反復する。ブロック504において、命令が制御依存命令 (例えば、算術命令) であると判定されない場合、ブロック508において、この命令

50

が廃棄される。本原理の様々な実施形態によれば、以下の典型的な擬似コード 1 で示すように、ブロック 504 において、命令が制御依存命令であると判定された場合、オペコードのみを取り出し、抽象化制御フロー特徴に含める。

【 0073 】

擬似コード 1 : 抽象化制御フロー特徴の生成

【表 2】

```

4028b0: mov  $0x413c98,%edi
4028bf: callq 402450
4028ca: xor  %eax,%eax
402905: mov  $0x413800,%edi
40290a: callq 402340
402975: je   402203
402994: test %eax,%eax
4029a4: jmp  4029e5
402bc4: mov  $0x0,0x218582(%rip)
4048dc: retq

```

10

【 0074 】

制御フロー情報（例えば、ファンクションコール、リターン、ジャンプ及びシステムコール）は、それらの動作を表す重要な記述である。しかしながら、それらの十分な情報（full information）を用いると、特定の細部（certain details）がマイナーチェンジのみで敏感に変化し、ノイズが過度に多くなることがある。例えば、プログラム命令は、他のサブルーチンに対する複数のジャンプ命令を用いるため、バイナリにおけるそれらの位置は、小さなコードパッチで変化する対象（subject）となる。したがって、本原理の様々な実施形態によれば、制御フロー情報のサブセットを用いる。このため、十分な情報を用いるよりも、敏感な変化に対してより回復力が持つ。様々な実施形態によれば、オペコードを含むが、制御依存命令のための命令パラメータを有しない命令情報のサブセット（例えば、ジャンプ、コール及びリターン命令）を用いてもよい。上述した実施形態では、非制御依存命令を用いていない。

20

30

【 0075 】

再び図 3 を参照すると、典型的な構造特徴 322 が以下の表 2 で示される。

【 0076 】

表 2 : 構造特徴

【表 3】

| セクション名    | セクションサイズ |
|-----------|----------|
| .gnu.hash | 104      |
| .dynsym   | 3096     |
| .init     | 26       |
| .plt      | 1808     |
| .text     | 63066    |
| .fini     | 9        |
| .got.plt  | 920      |
| .data     | 596      |
| .bss      | 3424     |
| .rodata   | 20732    |

10

## 【0077】

20

バイナリの他の特徴は、バイナリの構造情報である。本原理によれば、バイナリセクションの特性（例えば、名称、サイズ及びバイナリセクションの数）が、構造特徴322の1つとして用いられる。表2で詳細に記載したように、構造特徴322の典型的なテーブルにおいて、1つの列はバイナリセクションの名称を示し、他の列はバイナリセクションのサイズを示している。

## 【0078】

次に図6を参照すると、図6には本原理によるハイブリッド特徴類似度比較のためのシステム及び方法600が例示的に示されている。

## 【0079】

様々な実施形態によれば、（例えば、マルウェア攻撃を検出及び/または防止するための）バイナリの複数の特性の類似度を比較及び判定するために、本原理によるプログラムバイナリのハイブリッド特徴を用いてもよい。ブロック602、604及び606において、複数の（例えばN個の）プログラムバイナリから生成される1組のN個のハイブリッド特徴と、1つまたは複数のハイブリッド特徴パラメータ608（例えば、比較における各特徴の貢献を判定するレートの設定）とが、ブロック610におけるハイブリッド特徴類似度比較のために入力される。

30

## 【0080】

ハイブリッド特徴パラメータ608は、以下のように表すことができる。

C\_\_F：参照特徴のためのパラメータ

C\_\_R：リソース特徴のためのパラメータ

C\_\_C：抽象化制御フロー特徴のためのパラメータ

C\_\_S：構造特徴のためのパラメータ

C\_\_F、C\_\_R、C\_\_C、C\_\_Sは、0と1との間の比率である。

$C\_F + C\_R + C\_C + C\_S = 1$

C\_\_F\_\_P：類似したプログラム参照と照合するための閾値

C\_\_F\_\_F：類似した関数参照と照合するための閾値

40

## 【0081】

一実施形態において、組み合わせ生成器612は、想定される全てのバイナリ対の組み合わせを生成するように構成されている。ブロック616において、各2つのバイナリ（例えば、対）の類似度比較が実行され、ブロック618において、ハイブリッド差異スコ

50

アを生成する。本原理によれば、ブロック 6 1 4 において、類似度比較（例えば、特徴比較）が 1 つまたは複数のバイナリ対について反復され、ブロック 6 2 0 において、例えばマルウェア攻撃の検出及び防止で用いるために、類似度ベクトルが生成及び出力される。

【 0 0 8 2 】

本原理によれば、ブロック 6 1 0 におけるハイブリッド特徴類似度比較は以下のように実行すればよい。

【 数 6 】

Hybrid\_Feature\_Similarity\_Comparison(HFList, HFIndex, C\_F, C\_R, C\_C, C\_S, C\_F\_P, C\_F\_F)

Explanation of input parameters:

10

(1) A list of hybrid features of binary P1, P2, ..., PN

HFList = [HF(P1), HF(P2), ..., HF(PN)]

(2) Index of binaries

HFIndex = [P1, P2, ..., PN]

(3) Hybrid Feature Parameters C\_F, C\_R, C\_C, C\_S, C\_F\_P, C\_F\_F

【 0 0 8 3 】

類似度ベクトル 6 2 0 :

20

【 数 7 】

Similarity Vector (Component 620)

= an empty N \* N vector having N rows and N columns

【 0 0 8 4 】

組み合わせ生成器 6 1 2 :

【 数 8 】

// Make N \* N combinations of P1, P2, ..., PN excluding the comparison  
to itself and duplicates

30

CombinationList = GenerateCombination(HFIndex)

【 0 0 8 5 】

特徴比較 6 1 4 の反復 :

【 数 9 】

For each (PX, PY) in CombinationList:

// (Component 616, 618), Note: Definition of Component 618 follows.

40

HybridDiffScore (component 618) = HDS(HF(PX), HF(PY), C\_F,

C\_R, C\_C, C\_S, C\_F\_P, C\_F\_F)

Similarity Vector [PX, PY] = HybridDiffScore

return Similarity Vector (Component 620)

【 0 0 8 6 】

上述したように、ハイブリッド差異スコア 6 1 8 は 2 つのバイナリ間の類似度を表し、全ての組み合わせのスコアは類似度ベクトル 6 2 0 にて保存される。いくつかの実施形態

50

において、本原理によれば、バイナリのクラスタは、例えばクラスタリングアルゴリズムを、類似度ベクトルで保存されたデータに適用することで生成される。

【 0 0 8 7 】

次に図 7 を参照すると、図 7 には本原理による 2 つのバイナリの類似度比較のための方法 7 0 0 が例示的に示されている。

【 0 0 8 8 】

一実施形態において、本原理によれば、ブロック 7 2 2 において、2 つのハイブリッドバイナリ特徴 7 0 2、7 1 2 の類似度比較が実行される。比較は、ブロック 7 2 4、7 2 6、7 2 8 及び 7 3 0 において同じ種類の特徴間の差異スコア値を決定するために、同じ種類の 2 つの特徴（例えば、参照特徴 7 0 4、7 1 4；リソース特徴 7 0 6、7 1 6；抽象化制御フロー特徴 7 0 8、7 1 8；及び構造特徴 7 1 0、7 2 0）間でそれぞれ実行される。

10

【 0 0 8 9 】

本原理の様々な実施形態によれば、差異スコア値 7 2 4、7 2 6、7 2 8 及び 7 3 0 は、例えばハイブリッド特徴パラメータ 7 3 2 の値を差異スコア値 7 2 4、7 2 6、7 2 8 及び 7 3 0 と乗算することで、ブロック 7 3 4 にて決定ハイブリッド差異スコア (determined hybrid difference score) を決定するために用いられる。

【 0 0 9 0 】

一実施形態において、プログラムバイナリ P 1 及び P 2 間のハイブリッド差異スコア 7 3 4 は以下のように決定すればよい。

20

【数 1 0】

$$\begin{aligned} \text{HDS}(\text{HF}(\text{P1}), \text{HF}(\text{P2}), \text{C\_F}, \text{C\_R}, \text{C\_C}, \text{C\_S}, \text{C\_F\_P}, \text{C\_F\_F}) = \\ \text{C\_F} * \text{Sim\_F}(\text{F\_F}(\text{P1}), \text{F\_F}(\text{P2}), \text{C\_F\_P}, \text{C\_F\_F}) + \text{C\_R} * \text{Sim\_R}(\text{F\_R}(\text{P1}), \\ \text{F\_R}(\text{P2})) + \text{C\_C} * \text{Sim\_C}(\text{F\_C}(\text{P1}), \text{F\_C}(\text{P2})) + \text{C\_S} * \text{Sim\_S}(\text{F\_S}(\text{P1}), \\ \text{F\_S}(\text{P2})) \end{aligned}$$

【 0 0 9 1 】

Sim\_F、Sim\_R、Sim\_C 及び Sim\_S の定義は以下で説明する。

30

【 0 0 9 2 】

一実施形態において、プログラム P 1、P 2 の参照特徴 7 0 4、7 1 4 の類似度比較 7 2 4 は、以下のように実行すればよい。

【数 1 1】

```

Sim_F(F_F(P1), F_F(P2), C_F_P, C_F_F)
  Total = Min(|F_F(P1)|, |F_F(P2)|)
  Count = 0
  For (RP1, RF1) from F_F(P1)
    For (RP2, RF2) from F_F(P2)
      If (Difference(RP1, RP2) <= C_F_P and Difference(RF1,
        RF2) <= C_F_F)
        Count += 1
  return Count / Total

```

40

ここで、RP はプログラム基準を表し、RF は関数基準を表す。

【 0 0 9 3 】

一実施形態において、プログラム P 1、P 2 のリソース特徴 7 0 6、7 1 6 の類似度比較 7 2 6 は、以下のように実行すればよい。

50

【 0 0 9 4 】

【 数 1 2 】

```
Sim_R(F_R(P1), F_R(P2))
```

```
Total = Min(|F_R(P1)|, |F_R(P2)|)
```

```
Count = 0
```

```
For (K1, V1) from F_R(P1)
```

```
    For (K2, V2) from F_R(P2)
```

```
        If (K1 = K2 and V1 = V2)
```

```
            Count += 1
```

```
return Count / Total
```

10

ここで、Kは種類 (Kind) を表し、Vは値 (Value) を表す。

【 0 0 9 5 】

一実施形態において、プログラム P 1、P 2 の抽象化制御フロー特徴 7 0 8、7 1 8 の類似度比較 7 2 8 は、以下のように実行すればよい。

【 数 1 3 】

```
Sim_C(F_C(P1), F_C(P2))
```

```
Score = |LongestCommonSubsequence(F_C(P1), F_C(P2))| /
```

```
min(|F_C(P1)|, |F_C(P2)|)
```

```
return Score
```

30

【 0 0 9 6 】

一実施形態において、プログラム P 1、P 2 の構造特徴 7 1 0、7 2 0 の類似度比較 7 3 0 は、以下のように実行すればよい。

## 【数 1 4】

```
Sim_S(F_S(P1), F_S(P2))
```

```
Scores = []
```

```
For (Sec_1, size_1) from F_S(P1)
```

```
For (Sec_2, size_2) from F_S(P2)
```

```
if (Sec_1 = Sec_2)
```

```
Score = |size1 - size2| / max(size1, size2)
```

```
Score.push(Score)
```

```
Sum = 0, Count = 0
```

```
For score from Scores:
```

```
Sum += score
```

```
Count += 1
```

```
return Sum / Count
```

10

20

## 【0097】

様々な実施形態によれば、参照特徴724の比較において、例えば類似するプログラム名及び関数名を照合するために閾値を用いてもよい。例えば、同じライブラリの異なるバージョン（例えば、LibX\_V1及びLibX\_V2）は、閾値を用いて照合できる。対応する2つの閾値は、C\_F\_P及びC\_F\_Fとして追加される。

## 【0098】

次に図8を参照すると、図8には本原理によるハイブリッドプログラムバイナリ特徴の抽出及び比較のための典型的なシステム800が例示的に示されている。

30

## 【0099】

システム800の多くの態様が、例示及び明確さのために単数で記載されているが、システム800の説明に関して言及される要素は、複数の場合にも同じことが適用可能である。例えば、単一のコントローラ816が記載されているが、本原理の趣旨を維持する限り、本原理の教示による2つ以上のコントローラ816を用いることができる。さらに、ストレージデバイス818は、システム800に含まれる単なる1つの態様のすぎず、本原理の趣旨を維持する限り、複数の場合にも拡張できることを理解されたい。

## 【0100】

本原理の様々な実施形態によれば、本原理は、例えばプログラムのホワイトリスト登録、未知のバイナリの特性の判定、悪意のある機能の尤度判定及びマルウェアクラスタリングのために用いてもよい。

40

## 【0101】

プログラムのホワイトリスト登録に関して、これはソフトウェアの変形または異なるバージョンを判定することで達成してもよい。ソフトウェア会社及び開発者は、バグ修正、セキュリティ更新及び新たな特徴のために、ソフトウェアの様々なバージョンを作成する。例えば、ある会社が数週間または数ヶ月に1度バイナリを更新する場合、全てのバイナリ情報を企業内で取得すると、数十から数百の異なるバージョンのプログラムが存在する可能性がある。

## 【0102】

本発明は、そのような複数のプログラムの類似度を判定することができる。無害なソフ

50

トウェアの異なるバージョンを知ることは、これらを悪意のあるソフトウェアとの比較から除外して、マルウェア検出の複雑度を低減するために有用である。

【0103】

例えば、10,000のバイナリのうち、4,000のバイナリが既知で無害な500のバイナリの変形であると判定する。この場合、4,000のバイナリをホワイトリスト登録した後、マルウェアに関して6,000のバイナリのみを調べればよい。

【0104】

未知のバイナリの特性を判定することに関して、未知のソフトウェアバイナリに遭遇したとき、そのバイナリの特性が何であるか（例えば、そのバイナリがユーティリティプログラムであるか、ワードプロセッサであるか、またはマルウェアであるか）は分からない。本発明を用いて未知のバイナリを既知のソフトウェアのリストと比較することで、該バイナリの特性を正確に判定することができる。例えば、他のソフトウェアとの比較において、以下の類似度を得た場合：

- 未知のプログラムとファイル検索ユーティリティとの類似度：20%、
- 未知のプログラムとネットワークユーティリティとの類似度：80%、

本原理によれば、このプログラムは、ネットワーク機能及びファイル検索機能を有する可能性があると判定される。

【0105】

悪意のある機能の尤度を推定することに関して、未知のバイナリを悪意のあるソフトウェアバイナリのリストと比較してもよく、該バイナリが悪意のある機能を有する可能性を判定してもよい。例えば：

- 未知のプログラムとマルウェアXとの類似度：50%、
- 未知のプログラムとファイル検索ユーティリティとの類似度：20%、
- 未知のプログラムとネットワークユーティリティとの類似度：30%。

このプログラムは、マルウェアXで見られる悪意のある機能を有する可能性（例えば、尤度百分率）があると判定される。

【0106】

マルウェアクラスタリングにおいて、未知のバイナリが複数種類の悪意のあるソフトウェアとの類似度を示すとき、本発明は、ウィルス対策アプリケーションで用いるマルウェアクラスタリングによりマルウェアのカテゴリを理解するために役に立つ。例えば：

- 未知のプログラムとマルウェアファミリー1との類似度：70%、
- 未知のプログラムとマルウェアファミリー2との類似度：27%、
- 未知のプログラムとマルウェアファミリー3との類似度：3%

本原理によれば、このバイナリは、マルウェアファミリー2または3よりもマルウェアファミリー1により密接に関連していると判定され、それに応じて、ウィルス対策/マルウェア対策アプリケーションが悪意のあるソフトウェア攻撃を検出及び防止するために更新/適用される。

【0107】

システム800は、1つまたは複数のコンピューティングネットワーク及び/またはストレージデバイス818と接続可能なバス801を含むことができる。バイナリ特徴の抽出のためにプログラムバイナリ特徴抽出器802を用いてもよく、ハイブリッドバイナリ特徴生成器804を用いてハイブリッド特徴を生成してもよい。ハイブリッドバイナリ特徴は、例えばハイブリッド特徴パラメータ判定装置808によって提供されるハイブリッド特徴パラメータを入力としてさらに取得できる、ハイブリッドバイナリ特徴類似度解析器806を用いて解析してもよい。

【0108】

様々な実施形態によれば、バイナリ対どうしの差異を判定するために類似度判定装置812を用いてもよく、類似度ベクトル生成器810は、類似度比較に基づいて、類似度スコアベクトルを生成して出力する。生成された類似度ベクトルは、マルウェア攻撃解析器、検出器及びプリベント814において、マルウェア攻撃保護のために（例えば、リアル

10

20

30

40

50

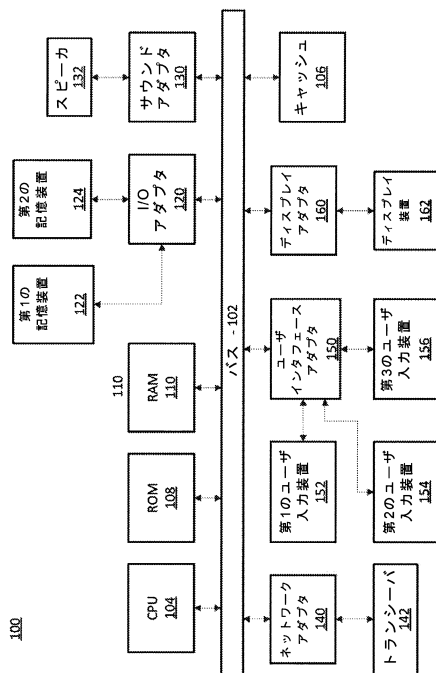
タイムで、将来において)用いてもよい。それは、例えば類似度ベクトルに基づいてマルウェア定義を更新すること、更新されたマルウェア定義によって検出されたマルウェアを隔離すること等を(例えば、手動または自動で)ウィルス対策ソフトウェアに命令するためにコントローラ816によって制御できる。本原理の様々な実施形態によれば、例えば、マルウェア攻撃の検出及び防止に使用するため、更新されたマルウェア定義、類似度比較の結果等を保存するためにストレージデバイス818を用いてもよい。

【0109】

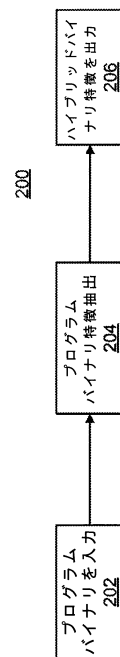
上記は、あらゆる観点において例示的(illustrative)かつ典型的(exemplary)であって限定的でないものと理解されるべきであり、本明細書で開示する本発明の範囲は、詳細な説明から決定されるべきではなく、特許法で認められた最大限の広さに基づいて解釈される特許請求の範囲から決定されるべきである。本明細書中に図示及び記載されている実施形態は、本発明の原理を説明するものにすぎず、本発明の範囲及び主旨から逸脱することなく当業者は様々な変更を実施することができることを理解されたい。当業者は、本発明の範囲及び精神から逸脱することなく、様々な他の特徴の組み合わせを実施できる。以上、本発明の態様について、特許法で要求される詳細及び特殊性と共に説明したが、特許証で保護されることを要求する特許請求の範囲は、添付の特許請求の範囲に示されている。

10

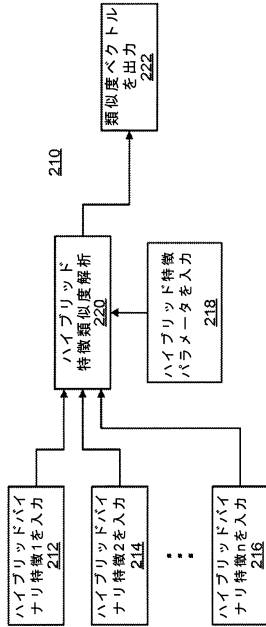
【図1】



【図2A】

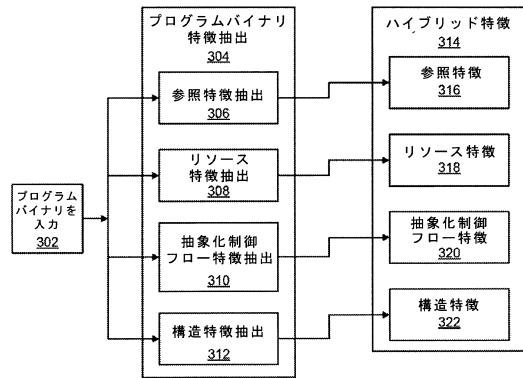


【図2B】



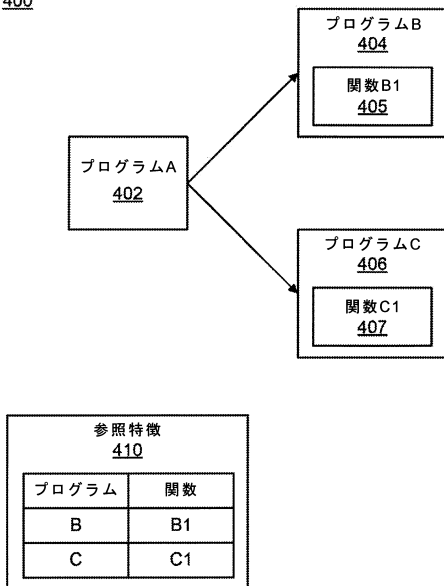
【図3】

300



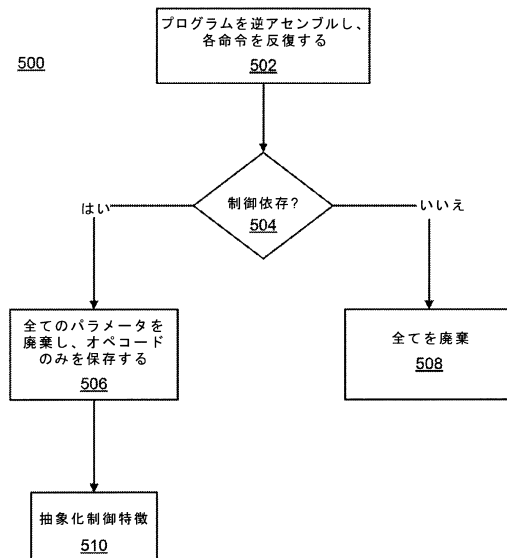
【図4】

400

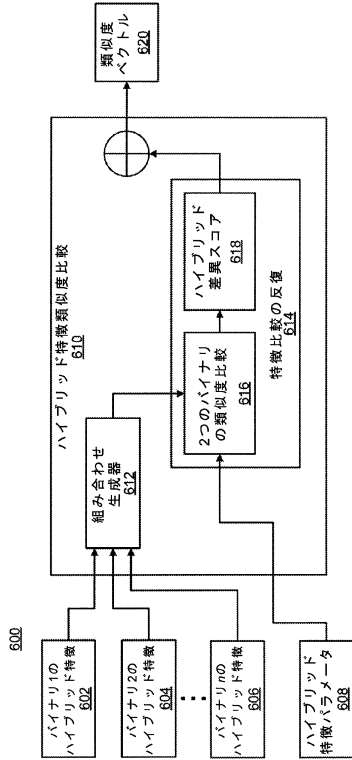


【図5】

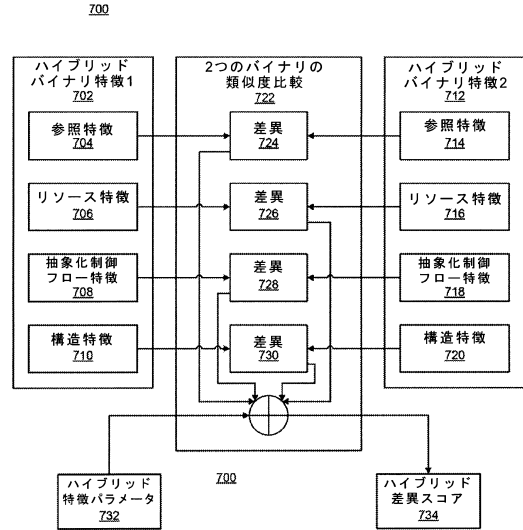
500



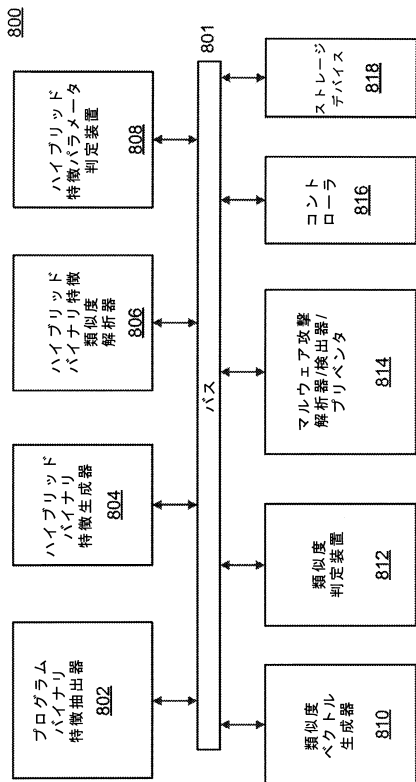
【図6】



【図7】



【図8】



## フロントページの続き

- (72)発明者 リー、 ジュンワン  
アメリカ合衆国 08540 ニュージャージー州 プリンストン クロムウェル コート 6
- (72)発明者 リ、 ジチュン  
アメリカ合衆国 08540 ニュージャージー州 プリンストン セイヤー ドライヴ 306
- (72)発明者 ウ、 ジェンユ  
アメリカ合衆国 08536 ニュージャージー州 プレインスポロ マーション レーン 4
- (72)発明者 ジー、 カンクック  
アメリカ合衆国 08540 ニュージャージー州 プリンストン ジョナソン デイトン コー  
ト 192
- (72)発明者 ジアン、 グオフェイ  
アメリカ合衆国 08540 ニュージャージー州 プリンストン ダンビー コート 5

審査官 田名網 忠雄

- (56)参考文献 特表2014-534531(JP, A)  
特開2013-077154(JP, A)  
特開2012-027710(JP, A)  
特開2010-198565(JP, A)  
国際公開第2015/099780(WO, A1)  
米国特許出願公開第2012/0159625(US, A1)

- (58)調査した分野(Int.Cl., DB名)  
G06F 21/50 - 21/57