

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5571035号  
(P5571035)

(45) 発行日 平成26年8月13日(2014.8.13)

(24) 登録日 平成26年7月4日(2014.7.4)

(51) Int.Cl. F 1  
**G 0 6 F 21/56 (2013.01)**  
 G 0 6 F 21/00 1 5 6 A  
 G 0 6 F 21/00 1 5 6 D

請求項の数 5 (全 20 頁)

(21) 出願番号	特願2011-120293 (P2011-120293)	(73) 特許権者	000004226
(22) 出願日	平成23年5月30日 (2011.5.30)		日本電信電話株式会社
(65) 公開番号	特開2012-248058 (P2012-248058A)		東京都千代田区大手町一丁目5番1号
(43) 公開日	平成24年12月13日 (2012.12.13)	(74) 代理人	100089118
審査請求日	平成25年7月17日 (2013.7.17)		弁理士 酒井 宏明
		(74) 代理人	100112656
			弁理士 宮田 英毅
		(72) 発明者	岩村 誠
			東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
		(72) 発明者	伊藤 光恭
			東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
		審査官	戸島 弘詩

最終頁に続く

(54) 【発明の名称】 特定装置、特定方法及び特定プログラム

(57) 【特許請求の範囲】

【請求項1】

パックされたプログラムコードをアンパックすることで得られるオリジナルコードから、少なくとも間接分岐命令と解釈可能な箇所を示すアドレスである間接分岐箇所と、直接分岐命令と解釈可能な箇所である直接分岐箇所を抽出する抽出部と、

前記抽出部により抽出された箇所間の分岐関係を解析することで、前記間接分岐箇所が直接的又は間接的に分岐先となる前記直接分岐箇所を識別する識別部と、

前記抽出部により抽出された箇所各々について分岐命令となる第1の確率を算出し、該箇所各々について算出した該第1の確率に基づいて、前記間接分岐箇所の先頭と、該間接分岐箇所が分岐先となる前記識別部により識別された前記直接分岐箇所の先頭のうちの、少なくとも一つの箇所の先頭が機械語命令の先頭となる第2の確率を算出する算出部と、

前記算出部により算出された第2の確率が閾値以上である場合に、前記間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する特定部と

を備えたことを特徴とする特定装置。

【請求項2】

前記分岐先アドレス格納場所が同一となる前記間接分岐箇所が複数抽出された場合に、前記算出部は、前記抽出部により抽出された箇所各々について前記第1の確率を算出し、該箇所各々について算出した該分岐命令となる該第1の確率に基づいて、前記分岐先アドレス格納場所が同一となる前記間接分岐箇所各々の先頭と、該間接分岐箇所各々のうちい

ずれかが分岐先となる前記識別部により識別された前記直接分岐箇所の先頭とのうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる前記第2の確率を算出することを特徴とする請求項1に記載の特定装置。

【請求項3】

前記算出部は、各箇所について、前記抽出部により抽出された前記箇所各々に該当する前記第1の確率を用いて、該当箇所を含む、他の命令との分岐元又は分岐先の分岐関係を示す分岐ツリー内に存在するすべての箇所が分岐命令とならない第3の確率を1から減算することで、前記第2の確率をそれぞれ算出することを特徴とする請求項1または2に記載の特定装置。

【請求項4】

特定装置で実行される特定方法であって、  
前記特定装置の抽出部が、バックされたプログラムコードをアンバックすることで得られるオリジナルコードから、少なくとも間接分岐命令と解釈可能な箇所を示すアドレスである間接分岐箇所と、直接分岐命令と解釈可能な箇所である直接分岐箇所を抽出して記憶部に格納する抽出工程と、

前記特定装置の識別部が、前記記憶部に記憶された箇所であって、前記抽出工程により抽出された箇所間の分岐関係を解析することで、前記間接分岐箇所が直接的又は間接的に分岐先となる前記直接分岐箇所を識別する識別工程と、

前記特定装置の算出部が、前記記憶部に記憶された箇所であって、前記抽出工程により抽出された箇所各々について分岐命令となる第1の確率を算出し、該箇所各々について算出した該第1の確率に基づいて、前記間接分岐箇所の先頭と、該間接分岐箇所が分岐先となる前記識別工程により識別された前記直接分岐箇所の先頭とのうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第2の確率を算出する算出工程と、

前記特定装置の特定部が、前記算出工程により算出された第2の確率が閾値以上である場合に、前記間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する特定工程と

を含んだことを特徴とする特定方法。

【請求項5】

コンピュータを請求項1～3のいずれか一つに記載の特定装置として機能させるための特定プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、特定装置、特定方法及び特定プログラムに関する。

【背景技術】

【0002】

近時、コンピュータウイルス等の悪意あるソフトウェアに対する対策が不可欠となっている。コンピュータウイルス等の悪意あるソフトウェアとは、例えば、マルウェアが該当する。マルウェアには、パッカーと呼ばれるツールにより、機械語のオリジナルコードに対して解読を困難にするため隠蔽処理であるパッキングが施される。オリジナルコードにパッキングが施されることで、アンチウイルスソフトのパターンマッチング機構を回避され、マルウェアの解析が困難となる。また、プログラムの解析を困難にする機構を持つパッカーもある。なお、オリジナルコードを復元する処理を「アンパッキング」とも称する。プログラムの解析を困難にする機構とは、例えば、Anti-DebugやAnti-VMなどがある。

【0003】

また、実行可能形式のファイルを受け付け、オリジナルコードを隠蔽しつつも実行可能形式を保ったファイルを出力するパッカーであるランタイムパッカーがある。ランタイムパッカーによりパッキングされたプログラムは、オリジナルコードを復元し、通常はローダが行う動的ライブラリのリンク処理等を実施した後に、オリジナルコードのエントリポ

10

20

30

40

50

イントへ処理を渡す。

【0004】

ここで、マルウェアの脅威を把握する際には、オリジナルコードを抽出し、抽出したオリジナルコードから外部関数を呼び出すときに使用される外部関数のアドレスが格納される格納場所であるインポートアドレス格納場所を特定する必要がある。なお、外部関数とは、例えば、Win32 API (Application Program Interface)、DLL (Dynamic Link Library) などがある。

【0005】

なお、インポートアドレス格納場所を特定する手法として、逆アセンブル結果に基づき、オリジナルコード領域から間接call命令を抽出し、抽出した間接call命令により示されるメモリ領域をインポートアドレス格納場所として特定する特定手法がある。

10

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2010-92179号公報

【非特許文献】

【0007】

【非特許文献1】Daniel Quist, Lorie Liebrock, Joshua Neil, Improving antivirus accuracy with hypervisor assisted analysis, Journal in Computer Virology (Published online: 6 April 2010) [online]、[2011年5月16日検索]、[インターネット] (URL: <http://csr.lanl.gov/vera/jcv-quist-liebrock-neil.pdf>)

20

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかしながら、上述の従来の特特定手法では、逆アセンブル結果が不正確である場合に、誤ったインポートアドレス格納場所が特定されるという課題がある。例えば、データ部分を誤って間接call命令として抽出されると、誤ったインポートアドレス格納場所が特定される。

30

【0009】

すなわち、オリジナルコードから間接call命令を抽出し、抽出した間接call命令により示されるメモリ領域をインポートアドレス格納場所とする上述の従来の特特定手法では、オリジナルコードに対応する正確な逆アセンブル結果が必要となる。ここで、マルウェア作者は、アンチウィルスソフトベンダ等による解析や対策から逃れることを目的として、ソースコードや逆アセンブルに要するシンボル情報などを非公開とすることがある。この場合、正確な逆アセンブル結果が得られず、誤ったインポートアドレス格納場所を特定されることがある。

【0010】

開示の実施形態は、上述に鑑みてなされたものであって、適切なインポートアドレス格納場所を特定可能となる特定装置、特定方法及び特定プログラムを提供することを目的とする。

40

【課題を解決するための手段】

【0011】

開示する特定装置は、一つの態様において、抽出部と、算出部と、特定部とを有する。抽出部は、バックされたプログラムコードをアンパックすることで得られるオリジナルコードから、少なくとも間接分岐命令と解釈可能な箇所を示すアドレスである間接分岐箇所を抽出する。算出部は、前記抽出部により抽出された箇所が分岐命令となる第1の確率を算出し、算出した第1の確率に基づいて、抽出された該箇所の先頭が機械語命令の先頭と

50

なる第2の確率を算出する。特定部は、前記算出部により算出された第2の確率が閾値以上である場合に、前記間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する。

【発明の効果】

【0012】

開示する特定装置の一つの態様によれば、適切なインポートアドレス格納場所を特定可能となるという効果を奏する。

【図面の簡単な説明】

【0013】

【図1】図1は、実施例1における特定装置の構成の一例を示すブロック図である。 10

【図2】図2は、実施例1における特定装置に入力されるオリジナルコードの一部を示す図である。

【図3】図3は、実施例1における暫定分岐情報格納テーブルに記憶された情報の一例を示す図である。

【図4】図4は、実施例1における暫定分岐ツリー情報テーブルに記憶された情報の一例を示す図である。

【図5】図5は、隠れマルコフモデルの一例を示す図である。

【図6】図6は、隠れマルコフモデルの一例を示す図である。

【図7】図7は、実施例1における分岐命令解析部による処理の一例について示す図である。 20

【図8-1】図8-1は、実施例1における分岐命令解析部による処理の一例について示す図である。

【図8-2】図8-2は、実施例1における分岐命令解析部による処理の一例について示す図である。

【図9】図9は、実施例1における分岐命令解析部による処理の一例について示す図である。

【図10】図10は、実施例1における特定装置による処理の流れの一例を示すフローチャートである。

【図11】図11は、特定装置による一連の処理を実行するための特定プログラムによる情報処理がコンピュータを用いて具体的に実現されることを示す図である。 30

【発明を実施するための形態】

【0014】

以下に、開示する特定装置、特定方法及び特定プログラムの実施例について、図面に基づいて詳細に説明する。なお、本実施例により開示する発明が限定されるものではない。各実施例は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

【実施例1】

【0015】

図1は、実施例1における特定装置の構成の一例を示すブロック図である。図1に示す例では、特定装置100は、記憶部110と、制御部120とを有する。また、特定装置100は、以下に詳細に説明するように、オリジナルコード10の入力を受け付け、インポートアドレス格納場所を出力する。 40

【0016】

なお、以下では、インポートアドレス格納場所20を「分岐先アドレス格納場所」とも記載する。言い換えると、以下では、特定装置100が、パッキングされたプログラムコードをアンパックすることで得られるオリジナルコード10の入力を受け付ける場合を例に説明する。言い換えると、パッキング前のプログラムコードの入力を受け付ける場合を用いて説明する。

【0017】

ただし、これに限定されるものではなく、例えば、特定装置100は、パッキングされたファイルを受け付け、受け付けたファイルからオリジナルコード10を抽出した上で、 50

後述する一連の処理を実行することでインポートアドレス格納場所情報を出力しても良い。

【0018】

図2は、実施例1における特定装置に入力されるオリジナルコードの一部を示す図である。図2に示す例では、説明の便宜上、オリジナルコード10に記載された一連の機械語のうちの機械語命令と解釈される部分の先頭を示す「アドレス」と対応付けて、「アドレス」から始まる機械語が機械語命令であると解釈した場合に得られる「ニーモニック」を示した。

【0019】

図2に示す例では、アドレス「0x0005」とニーモニック「call 0x000F」とを対応付けて記憶する。すなわち、オリジナルコード10のうち、「0x0005」が機械語命令であると解釈した場合には、ニーモニック「call 0x000F」が得られることを示す。なお、図2に示す例では、ニーモニックが「call」で始まる場合には、call命令を示す。また、ニーモニックが「jump」で始まる場合には、jump命令を示す。また、図2において、[ ]がある場合には、間接分岐命令を示し、[ ]がない場合には、直接分岐命令を示す。例えば、ニーモニック「call 0x000F」は、「0x000F」に分岐する直接call命令を示す。また、ニーモニック「call [0x2004]」は、アドレス「0x2004」に格納されている値を分岐先アドレスとして分岐する間接call命令を示す。

【0020】

記憶部110は、制御部120と接続される。記憶部110は、制御部120による各種処理に用いるデータを記憶する。記憶部110は、例えば、RAM(Random Access Memory)やROM(Read Only Memory)、フラッシュメモリ(Flash Memory)などの半導体メモリ素子、又は、ハードディスクや光ディスクなどである。図1に示す例では、記憶部110は、暫定分岐情報格納テーブル111と、暫定分岐ツリー情報テーブル112とを有する。

【0021】

暫定分岐情報格納テーブル111は、後述する制御部120の分岐命令解析部121による処理結果を記憶する。図3は、実施例1における暫定分岐情報格納テーブルに記憶された情報の一例を示す図である。

【0022】

図3に示す例では、暫定分岐情報格納テーブル111は、「アドレス」と、「種別」と、「命令確率」と、「暫定分岐先アドレス」と、「暫定分岐先アドレス格納場所」とを対応付けて記憶する。ここで、「アドレス」は、オリジナルコード10に記載された一連の機械語のうち機械語命令と解釈される部分の先頭を示す。「種別」は、分岐命令の種別を示す。例えば、「種別」は、「アドレス」から始まる機械語が機械語命令であると解釈した場合に、「直接call命令」か「間接call命令」か「直接jump命令」か「間接jump命令」かを示す。また、「命令確率」は、「アドレス」から始まる機械語が機械語命令である確率を示す。「命令確率」は「第1の確率」とも記載する。「暫定分岐先アドレス」と「暫定分岐先アドレス格納場所」とは、それぞれ、「アドレス」から始まる機械語が機械語命令についての分岐先アドレス、又は、分岐先アドレスが格納された格納場所を示す。

【0023】

例えば、暫定分岐情報格納テーブル111は、アドレス「0x0005」と、種別「直接call」と、命令確率「0.2」と、暫定分岐先アドレス「0x000F」とを含むレコードを記憶する。すなわち、暫定分岐情報格納テーブル111は、アドレス「0x0005」から始まる機械語が機械語命令であると解釈した場合に、分岐先アドレスが「0x000F」となる直接call命令である確率が「0.2」であることを記憶する。なお、暫定分岐情報格納テーブル111は、直接分岐命令についてのレコードでは、「暫定分岐先アドレス」を記憶する。また、暫定分岐情報格納テーブル111は、間接分岐命令

10

20

30

40

50

についてのレコードでは、「暫定分岐先アドレス格納場所」を記憶する。

【0024】

なお、暫定分岐情報格納テーブル111に記憶された情報は、分岐命令解析部121により格納され、分岐関係解析部122により用いられる。

【0025】

暫定分岐ツリー情報テーブル112は、後述する制御部120の分岐関係解析部122による処理結果を記憶する。具体的には、暫定分岐ツリー情報テーブル112は、分岐命令解析部121により抽出された分岐命令と解釈可能な部分各々について、他の命令との分岐元又は分岐先を示す分岐関係を示す暫定分岐関係ツリーを記憶する。

【0026】

図4は、実施例1における暫定分岐ツリー情報テーブルに記憶された情報の一例を示す図である。図4に示す例では、図3における「アドレス」を用いて、分岐命令解析部121により抽出された部分各々を示した。図3に示すように、暫定分岐ツリー情報テーブル112は、例えば、「0x0005」が分岐元となり「0x000F」が分岐先となる関係を記憶し、「0x000F」が分岐元となり「0x1004」が分岐先となる関係を記憶する。

【0027】

なお、暫定分岐ツリー情報テーブル112に記憶された情報は、分岐関係解析部122により格納され、インポートアドレス格納場所特定部123により用いられる。

【0028】

制御部120は、記憶部110と接続される。制御部120は、各種の処理手順などを規定したプログラムを記憶する内部メモリを有し、種々の処理を制御する。制御部120は、例えば、ASIC(Application Specific Integrated Circuit)、FPGA(Field Programmable Gate Array)、CPU(Central Processing Unit)、MPU(Micro Processing Unit)などの電子回路である。図1に示す例では、制御部120は、分岐命令解析部121と、分岐関係解析部122と、インポートアドレス格納場所特定部123とを有する。

【0029】

分岐命令解析部121は、パックされたプログラムコードをアンパックすることで得られるオリジナルコード10から、少なくとも間接分岐命令と解釈可能な箇所を示すアドレスである間接分岐箇所を抽出する。また、分岐命令解析部121は、間接分岐箇所に加えて、直接分岐命令と解釈可能な箇所である直接分岐箇所を抽出する。すなわち、分岐命令解析部121は、例えば、オリジナルコード10の入力を受け付けると、分岐命令と解釈可能な部分をすべて抽出する。

【0030】

より詳細な一例をあげて説明すると、分岐命令解析部121は、オリジナルコード10から、直接call命令と間接call命令と直接jump命令と間接jump命令とのうち、いずれかであると解釈可能な部分を抽出する。

【0031】

なお、分岐命令解析部121により抽出された箇所は、分岐命令と解釈可能である箇所であり、実際には分岐命令ではない箇所も含まれることがあり得る。

【0032】

ここで、分岐命令について簡単に説明した上で、直接分岐命令と間接分岐命令との違いについて簡単に説明する。分岐命令とは、次に実行される命令を切り替える命令である。なお、一連の命令は、分岐命令でない場合には、逐次順番に実行される。分岐命令には、例えば、jump命令やcall命令などがある。jump命令やcall命令は、次に実行する命令を変更する。言い換えると、jump命令やcall命令の次に実行される命令は、jump命令やcall命令の次に記載された命令ではなく、jump命令やcall命令のオペランドとして指定された命令となる。例えば、jump命令やcall

10

20

30

40

50

命令は、オペランドとしてメモリアドレスを指定しており、オペランドとして指定されたメモリアドレスに格納されている値が分岐先アドレスとなる。つまり、j u m p 命令や c a l l 命令の次に、オペランドとして指定されたメモリアドレスに格納されている値が取得されたり処理されたりする。

#### 【 0 0 3 3 】

直接分岐命令と間接分岐命令との違いについて簡単に説明する。直接分岐命令には、例えば、直接 c a l l 命令と直接 j u m p 命令とがある。また、間接分岐命令には、例えば、間接 c a l l 命令と間接 j u m p 命令とがある。直接分岐命令では、命令内に分岐先となる分岐先アドレスが明示されており、命令解読時に分岐先が決定される。間接分岐命令では、命令内に分岐先アドレスが明示されておらず、実際に間接分岐命令を実行する際に初めて分岐先アドレスが決定される。例えば、間接分岐命令では、命令内には分岐先アドレスが明示されておらず、分岐先アドレス格納場所から分岐先アドレスを読み出すことで、分岐先アドレスが決定される。例えば、間接分岐命令は、実行時にロードされた D L L が分岐先となり、D L L がロードされるメモリアドレスが動的に決定される場合に利用される。

10

#### 【 0 0 3 4 】

また、分岐先アドレス格納場所と、間接分岐箇所と、直接分岐箇所との関係について簡単に説明する。分岐先アドレス格納場所を参照する分岐命令は、間接分岐箇所のみとなる。また、間接分岐箇所を分岐先とする分岐命令は、直接分岐箇所のみとなる。間接分岐箇所を分岐先とする直接分岐箇所は、ない場合もある。直接分岐箇所を分岐先とする分岐命令は、直接分岐箇所のみとなる。直接分岐箇所を分岐先とする直接分岐箇所は、ない場合もある。

20

#### 【 0 0 3 5 】

分岐命令解析部 1 2 1 の説明に戻る。分岐命令解析部 1 2 1 は、公知の手法を用いて、オリジナルコード 1 0 から分岐命令と解釈可能な箇所を抽出する。例えば、分岐命令解析部 1 2 1 は、機械語にて記載された命令である機械語命令として解釈可能なバイト列を識別することで、直接分岐箇所や間接分岐箇所を抽出する。詳細な一例をあげて説明すると、I A - 3 2 命令セットの場合、分岐命令解析部 1 2 1 は、1 バイト目が「 0 x E B 」であれば、直接 j u m p 命令であると解釈可能な部分として抽出する。

#### 【 0 0 3 6 】

また、分岐命令解析部 1 2 1 は、抽出した間接分岐箇所が分岐命令となる第 1 の確率を算出する。また、同様に、分岐命令解析部 1 2 1 は、抽出した間接分岐箇所が分岐命令となる第 1 の確率に加えて、抽出した直接分岐箇所が分岐命令となる第 1 の確率を算出する。なお、分岐命令解析部 1 2 1 による第 1 の確率を算出する算出手法の一例については、後述する。

30

#### 【 0 0 3 7 】

また、分岐命令解析部 1 2 1 は、抽出した間接分岐箇所と直接分岐箇所とについて、算出した第 1 の確率と対応付けて暫定分岐情報格納テーブル 1 1 1 に格納する。また、分岐命令解析部 1 2 1 は、抽出した直接分岐箇所について格納する際には、抽出した直接分岐箇所が実際に直接分岐命令である場合に分岐先と解釈される分岐先アドレスを「暫定分岐先アドレス」として併せて格納する。また、分岐命令解析部 1 2 1 は、抽出した間接分岐箇所を格納する際には、抽出した間接分岐箇所が実際に間接分岐命令である場合に分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所と解釈される位置を「暫定分岐先アドレス格納場所」として併せて格納する。

40

#### 【 0 0 3 8 】

例えば、分岐命令解析部 1 2 1 は、先頭が「 0 x 0 0 0 5 」から始まる箇所が直接 c a l l 命令である第 1 の確率が「 0 . 2 」であると算出した場合を用いて説明する。また、先頭が「 0 x 0 0 0 5 」から始まる箇所が直接 c a l l 命令であり、分岐先の先頭のアドレスが「 0 x 0 0 0 F 」である場合を用いて説明する。この場合、分岐命令解析部 1 2 1 は、アドレス「 0 x 0 0 0 5 」と、種別「直接 c a l l 」と、命令確率「 0 . 2 」と、暫

50

定分岐先アドレス「0x000F」とを含むレコードを暫定分岐情報格納テーブル111に格納する。

【0039】

分岐関係解析部122は、分岐命令解析部121により抽出された箇所間の分岐関係を解析することで、間接分岐箇所が直接的又は間接的に分岐先となる直接分岐箇所を識別する。そして、分岐関係解析部122は、分岐命令解析部121により抽出された箇所間の分岐関係を示すツリーである暫定分岐関係ツリーを生成する。なお、分岐関係解析部122は、「識別部」とも称する。

【0040】

例えば、分岐関係解析部122は、分岐命令解析部121により生成された暫定分岐情報格納テーブル111内に記憶されたレコードのうち、種別が「直接call」又は「直接jump」となるレコードを識別する。そして、分岐関係解析部122は、種別が「直接call」又は「直接jump」となるレコードの「先頭」から抽出された箇所が、他のレコードの「暫定分岐先アドレス」と一致するか否かを判定する。すなわち、直接分岐命令として抽出された箇所の分岐元が、直接分岐命令として抽出された箇所と一致するか否かを判定する。なお、レコードの「先頭」から抽出された箇所とは、例えば、図3の「アドレス」に示す箇所が該当する。ここで、分岐関係解析部122は、一致すると判定した場合には、一致すると判定したレコードに対応する箇所各々が分岐元と分岐先とが関係あると判定する。また、分岐関係解析部122は、種別が「直接call」又は「直接jump」となるレコードについての判定結果に加えて、種別が「間接call」又は「間接jump」となるレコードの「先頭」と「暫定分岐先アドレス格納場所」とを加えることで、暫定分岐関係ツリーを生成する。

【0041】

例えば、分岐関係解析部122は、図3に示す暫定分岐情報格納テーブル111に基づいて処理を実行する場合には、図4に示すような暫定分岐ツリーを生成し、暫定分岐ツリー情報テーブル112に格納する。図4は、実施例1における分岐関係解析部により生成される暫定分岐関係ツリーの一例を示す。図4に示す例では、説明の便宜上、直接分岐箇所については4角形で表し、間接分岐箇所については3角形で表し、暫定分岐先アドレス格納場所については丸で表した。例えば、分岐関係解析部122は、直接分岐箇所についての分岐関係の判定結果からツリーを生成する。その後、分岐関係解析部122は、直接分岐箇所の分岐先となる間接分岐箇所と、間接分岐箇所の暫定分岐先アドレス格納場所とを加えることで、暫定分岐ツリーを生成する。なお、図4に示す例では、「0x0005」が分岐元となり「0x000F」が分岐先となり、「0x000F」が分岐元となり「0x1004」が分岐先となる。

【0042】

インポートアドレス格納場所特定部123は、間接分岐箇所の先頭と、間接分岐箇所が分岐先となる直接分岐箇所の先頭とのうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第2の確率を算出する。なお、間接分岐箇所の「先頭」や直接分岐箇所の「先頭」とは、間接分岐箇所や直接分岐箇所としてレコードから抽出された箇所のうち、先頭部分にある箇所を示す。すなわち、図3に示す例では、「アドレス」に示された箇所が該当する。

【0043】

また、インポートアドレス格納場所特定部123は、暫定分岐先アドレス格納場所が同一となる間接分岐箇所が複数抽出された場合に、暫定分岐先アドレス格納場所が同一となる間接分岐箇所各々の先頭と、暫定分岐先アドレス格納場所が同一となる間接分岐箇所各々のうちいずれかが分岐先となる直接分岐箇所の先頭とのうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第2の確率を算出する。つまり、インポートアドレス格納場所特定部123は、抽出された間接分岐箇所の分岐先アドレスが格納されている暫定分岐先アドレス格納場所が実際に分岐先アドレス格納場所となる確率を算出する。

【0044】

10

20

30

40

50



ここで、計算の単純化を目的として、分岐命令解析部 1 2 1 により抽出された箇所各々が機械語命令として解釈される事象がお互いに独立であると仮定して説明する。また、暫定分岐先アドレス格納場所「X」が実際に分岐先アドレス格納場所となる第 2 の確率を「PX」として説明する。また、「X」と同一の暫定分岐ツリー内に存在する分岐命令「Y<sub>i</sub> ( 1 ≤ i ≤ N, N は X と同一ツリー内に存在する分岐命令の数を示し、i はツリー内の分岐命令各々を表すインデックスを示す。 ) 」として説明する。また、分岐命令「Y<sub>i</sub>」が実際に分岐命令となる第 1 の確率を「PY<sub>i</sub>」とした上で説明する。この場合、「PX」は、下記の式 ( 1 ) にて算出可能となる。

【 0 0 4 5 】

【数 1】

$$PX = 1 - \prod(1 - PY_i) \quad \dots (1)$$

10

【 0 0 4 6 】

すなわち、式 ( 1 ) に示すように、「X」と同一の暫定分岐ツリー内に存在する分岐命令のうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第 2 の確率は、「1」から、「X」と同一の暫定分岐ツリー内に存在する分岐命令のすべてが機械語命令の先頭とはならない確率を減算することで得られる。なお、式 ( 1 ) の「 $\prod$ 」は、総積を示す。言い換えると、式 ( 1 ) に示す例では、「 $\prod(1 - PY_i)$ 」は、「X」と同一の暫定分岐ツリー内に存在するすべての分岐命令「Y<sub>i</sub>」について「1 - PY<sub>i</sub>」を乗算した結果を示す。

20

【 0 0 4 7 】

すなわち、インポートアドレス格納場所特定部 1 2 3 は、抽出された箇所各々に該当する第 1 の確率を用いて、抽出された同一ツリー内にある箇所すべてが分岐命令とならない第 3 の確率を 1 から減算することで、第 2 の確率を算出する。ここで、第 3 の確率とは、例えば、式 ( 1 ) における「 $\prod(1 - PY_i)$ 」を示す。

【 0 0 4 8 】

より詳細な一例をあげて説明すると、インポートアドレス格納場所特定部 1 2 3 は、処理対象となる箇所について算出された第 1 の確率を 1 から減算した値を処理対象となる箇所ごとに算出し、処理対象となる箇所ごとに算出された値各々を積算し、積算することで得られた値を 1 から減算することで、第 2 の確率を算出する。処理対象となる箇所とは、同一の暫定分岐ツリー内に存在する分岐命令となり、例えば、間接分岐箇所の先頭一つであったり、間接分岐箇所の先頭、及び間接分岐箇所が直接的又は間接的に分岐先となる直接分岐箇所の先頭であったり、分岐先アドレス格納場所が同一となる間接分岐箇所各々の先頭及び間接分岐箇所各々のうちいずれかが分岐先となる直接分岐箇所の先頭であったりする。

30

【 0 0 4 9 】

例えば、図 3 における暫定分岐先アドレス格納場所「0x2004」がインポートアドレス格納場所 2 0 となる第 2 の確率は下記の式 ( 2 ) にて算出される。

【 0 0 5 0 】

【数 2】

$$P(0x2004) = 1 - (1 - 0.2) \times (1 - 0.4) \times (1 - 0.1) \times (1 - 0.6) \times (1 - 0.8) = 0.96544 \quad \dots (2)$$

40

【 0 0 5 1 】

同様に、図 3 における暫定分岐先アドレス格納場所「0x20AA」がインポートアドレス格納場所 2 0 である第 2 の確率は下記の式 ( 3 ) にて算出される。

【 0 0 5 2 】

【数 3】

$$P(0x20AA) = 1 - (1 - 0.4) = 0.4 \quad \dots (3)$$

【 0 0 5 3 】

50

また、インポートアドレス格納場所特定部 1 2 3 は、算出された第 2 の確率が閾値以上である場合に、間接分岐箇所が間接分岐命令であるとした場合における間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する。例えば、閾値として「0.5」を用いる場合には、インポートアドレス格納場所特定部 1 2 3 は、アドレス「0x2004」をインポートアドレス格納場所 20 として特定し、アドレス「0x20AA」をインポートアドレス格納場所 20 として特定しない。

【0054】

なお、インポートアドレス格納場所特定部 1 2 3 は、間接分岐箇所について算出された第 1 の確率が閾値以上となる場合に分岐先アドレス格納場所を特定しても良い。言い換えると、同一の暫定分岐ツリー内に存在する分岐命令のうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第 2 の確率を用いることなく、間接分岐箇所について算出された第 1 の確率のみを用いても良い。

10

【0055】

ここで、同一の暫定分岐ツリー内に存在する分岐命令のうち、少なくとも一つの箇所の先頭が機械語命令の先頭となる第 2 の確率が閾値以上となる場合に、分岐先アドレス格納場所を特定する点について補足する。暫定分岐ツリー内に分岐命令が存在する場合、その分岐先も機械語命令と考えられることを踏まえ、インポートアドレス格納場所特定部 1 2 3 は、少なくとも一つの箇所の先頭が機械語命令の先頭となる第 2 の確率を算出して用いても良い。

【0056】

20

ここで、インポートアドレス格納場所特定部 1 2 3 により分岐先アドレス格納場所として特定されるのは、間接分岐箇所として抽出された箇所のうち、算出された第 2 の確率が閾値以上となった間接分岐箇所についての暫定分岐先アドレス格納場所となる。言い換えると、インポートアドレス格納場所特定部 1 2 3 により分岐先アドレス格納場所として特定されるのは、直接分岐箇所の暫定分岐先アドレスではない。この点について補足する。上述したように、例えば、ランタイムパッカーによりパッキングされたプログラムは、通常はロードが行う動的ライブラリのリンク処理等を実施した後に、オリジナルコードのエントリーポイントへ処理が渡される。ここで、Win32 API や DLL に含まれる外部関数は、通常、実行時にメモリ上に一度ロードされた上で、実行時に決定されたメモリアドレスが参照されて実行される。このことを踏まえ、インポートアドレス格納場所特定部 1 2 3 は、直接分岐箇所の分岐先アドレスではなく、間接分岐箇所の分岐先アドレス格納場所を特定する。

30

【0057】

ここで、機械語命令として解釈可能な部分が実際に機械語命令である第 1 の確率を算出する分岐命令解析部 1 2 1 による処理の一例について説明する。以下では、隠れマルコフモデルに基づき箇所各々について命令確率を算出する場合を用いて説明するが、これに限定されるものではない。

【0058】

分岐命令解析部 1 2 1 は、例えば、隠れマルコフモデルによりオリジナルコード 10 をモデル化し、Forward / Backward アルゴリズムを用いることで、オリジナルコード 10 の各バイトが機械語命令の先頭である確率を算出する。なお、以下に説明する処理は一例であり、これに限定されるものではない。

40

【0059】

隠れマルコフモデルによりオリジナルコード 10 をモデル化する点について簡単に説明する。なお、以下では、オリジナルコード 10 に対して逆アセンブルが行われることで、オリジナルコード 10 を構成する複数のバイナリ値が複数の単語に分割され、分割された複数の単語それぞれに「命令部」又は「データ部」のいずれかの状態であることを示す「タグ」が割り当てられ、「命令部」としての「タグ」が割り当てられた単語の命令長に基づいて、ニーモニック（アセンブルコード）を当てはめられた場合を用いて説明する。すなわち、オリジナルコード 10 を構成する複数のバイナリ値が複数の単語に分割されてお

50

り、複数の単語それぞれにタグが付されている場合を用いて説明する。

【0060】

また、以下では、「入力バイナリ列：X」とは、「逆アセンブル」の対象となる「オリジナルコード10」のバイナリ列を示し、「N」個のバイナリ値であるとする。式(4)に示すように、「逆アセンブル」の対象となる「オリジナルコード」を構成するN個のバイナリ値は、「 $x_1 \sim x_N$ 」として表される。

【0061】

【数4】

$$\text{入力バイナリ列 } X = x_1^N = x_1, x_2, \dots, x_N \quad \dots (4)$$

10

【0062】

また、「単語列：w」とは、「入力バイナリ列：X」を1命令の「命令部」もしくは1データの「データ部」としての単語として分割したものである。式(5)に示すように、「入力バイナリ列：X」を分割したM個の単語それぞれは、「 $w_1 \sim w_M$ 」として表される。「 $w_i$ 」は1命令もしくは1データを表す。なお、「命令部」は、複数のバイナリ値から構成される場合もあるため、「単語数：M」「入力バイナリ数：N」となる。

【0063】

【数5】

$$\text{単語列 } W = w_1^M = w_1, w_2, \dots, w_M \quad \dots (5)$$

20

【0064】

また、「タグ列：T」とは、単語「 $w_1 \sim w_M$ 」それぞれに対して、「命令部」か「データ部」であるかの「タグ」が割り当てられたものである。式(6)に示すように、単語「 $w_1 \sim w_M$ 」に対応付けてタグ「 $t_1 \sim t_M$ 」として表される。

【0065】

【数6】

$$\text{タグ列 } T = t_1^M = t_1, t_2, \dots, t_M \quad \dots (6)$$

【0066】

また、「命令タグ集合：I」は、「命令部」としての状態を表す「タグ」の集合であり、「データタグ集合：D」は、「データ部」としての状態を表す「タグ」の集合である。ここで、タグ「 $t_i (1 \leq i \leq M)$ 」は、命令部かデータ部かのいずれかとなる。この結果、式(7)に示すように、「 $t_i$ 」は、「命令タグ集合：I」あるいは「データタグ集合：D」のいずれかに属する。

30

【0067】

【数7】

$$\begin{aligned} &\text{命令タグ集合 } I, \text{ データタグ集合 } D \\ &t_i \in I \cup D \quad \dots (7) \end{aligned}$$

【0068】

図5及び図6は、隠れマルコフモデルの一例を示す図である。すなわち、図5に示すように、「命令タグ集合：I」に属するタグを「継続命令状態：S」及び「データ直前命令状態：T」の2種類に更に分割し、「データタグ集合：D」に属する「データ状態：U」と合わせて3種類の状態から構成される隠れマルコフモデルを前提とする。

40

【0069】

「継続命令状態：S」は、1命令を出力したのち、引き続き「継続命令状態：S」に留まる場合と、「データ直前命令状態：T」に遷移する場合とがある。

【0070】

「データ直前命令状態：T」は、「継続命令状態：S」と同様に、1命令を出力するが、その遷移先は、「データ状態：U」のみとなる。一般的に、後方にデータが続く命令は、無条件分岐であることが多い。この結果、命令状態を、継続命令状態と、データ直前命

50

令状態に分割することで、逆アセンブルの精度を向上することが期待できる。

【0071】

ここで、「継続命令状態：S」、「データ直前命令状態：T」、又は、「データ状態：U」のいずれかの「状態*i*」から始まる確率（初期確率）を「 $p_i$ 」とし、「状態*i*」から「状態*j*」へ遷移する確率（遷移確率）を「 $a_{ij}$ 」とし、「状態*i*」におけるシンボルとしての「単語*w*」が出力される確率（シンボル出力確率）を「 $b_i(w)$ 」とする。

【0072】

このような隠れマルコフモデルの一例において、「データ状態：U」で出力されるシンボルをデータ1バイトとすると、これにより、「データ状態：U」におけるシンボル出力確率「 $b_U(w)$ 」において、「*w*」は、「0以上255以下の範囲にある整数」とすることができる。

10

【0073】

これに対して、「命令タグ集合：I」に属する「状態*i*」において出力されるシンボルの長さ（シンボル長）は、1命令の長さとなる。ここで、複合命令セットコンピュータ（CISC：Complex Instruction Set Computer）アーキテクチャの代表的な例であるIntel社の「x86命令」の場合、1命令の長さは最大で16バイトにも及ぶため、そのまま統計的に信頼できるシンボル出力確率「 $b_i(w)$ 」を学習することは容易ではない。こうした状況に対応するため、図6を用いて近似的にシンボル出力確率「 $b_i(w)$ 」（*i*はIに属する）を算出する方法について述べる。「x86命令」は、「PREFIX（命令長：0～4バイト）」、「OPCODE（命令長：1～2バイト）」、「ModRM（命令長：0～1バイト）」、「SIB（命令長：0～1バイト）」、「DISPLACEMENT（命令長：0～4バイト）」、「IMMEDIATE（命令長：0～4バイト）」といった命令部から構成される。また、これらの命令部間の遷移パターンは、図6に示すパターンとなる。

20

【0074】

ここで、図6に示す遷移パターンによって遷移する各命令部を「状態」とし、「命令開始状態」と「命令終了状態」とを除いた各状態（PREFIX, OPCODE, ModRM, SIB, DISPLACEMENT, IMMEDIATE）では、1バイトの命令部を出力するとする。

【0075】

また、「単語*w*」を1バイトごとに分解した結果を、式（8）によって表し（「 $x_a \sim x_b$ 」）、対応する命令部の種別を、式（9）によって表すとする。

30

【0076】

【数8】

$$W \text{ を } 1 \text{ バイトごとに分解した結果: } x_a^b \quad \dots \quad (8)$$

【数9】

$$\text{対応する命令部種別: } v_a^b \in \{\text{PREFIX, OPCODE, ModRM, SIB, DISPLACEMENT, IMMEDIATE}\}$$

$$\dots \quad (9)$$

40

【0077】

このとき、「命令部1バイトを出力する確率は、その時点での命令部の状態によってのみ決まる」と仮定し、更に、「命令部の状態へ遷移する確率は、一つ前の命令部の状態によって決まる」と仮定すると、「命令タグ集合：I」に属する「状態*i*」におけるシンボルとしての「単語*w*」のシンボル出力確率「 $b_i(w)$ 」は、式（10）に示すように、近似することができる。なお、ここで示す「命令部の状態」とは、図6のOPCODEやModRMなどを示す。

【0078】

【数 1 0】

$$b_i(w) \approx \prod_{i=a}^b P(x_i | v_i) P(v_i | v_{i-1}) \quad \dots (10)$$

【0079】

これにより、隠れマルコフモデルにおける最尤状態系列算出の問題として、オリジナルコード10から命令部とデータ部とを識別するために用いるモデルパラメータは、命令部間の状態遷移確率と、命令部ごとの1バイトの出現確率のみとすることができる。この命令部に関するモデルパラメータは、「継続命令状態：S」と「データ直前命令状態：T」とで個別に持たせる。

10

【0080】

次に、分岐命令解析部121によるモデルパラメータの学習について説明する。分岐命令解析部121は、『命令部間の遷移確率「 $P(v_i | v_{i-1})$ 」及び各命令部における1バイト値のシンボル出力確率「 $P(x_i | v_i)$ 」』を、各状態(タグ)間での遷移回数及び各状態(タグ)におけるシンボル出現回数をカウントすることで算出する。ここで、『命令部間の遷移確率「 $P(v_i | v_{i-1})$ 」及び各命令部における1バイト値のシンボル出力確率「 $P(x_i | v_i)$ 」』は、式(11)~式(14)を算出するためのモデルパラメータである。

【0081】

【数 1 1】

$$\pi_i (i \in I \cup D) \quad \dots (11)$$

20

【数 1 2】

$$a_{ij} (i, j \in I \cup D) \quad \dots (12)$$

【数 1 3】

$$b_i(w) (i \in D, 0 \leq w \leq 255) \quad \dots (13)$$

【数 1 4】

$$b_i(w) (i \in I) \quad \dots (14)$$

30

【0082】

ここで、式(11)は、「命令タグ集合：I」又は「データタグ集合：D」のいずれかに属する「状態i」の初期確率「 $\pi_i$ 」を示す。式(12)は、「命令タグ集合：I」又は「データタグ集合：D」のいずれかに属する「状態i」から「命令タグ集合：I」又は「データタグ集合：D」のいずれかに属する「状態j」への遷移確率「 $a_{ij}$ 」を示す。式(13)は、「状態i」が「データタグ集合：D」に属する場合のシンボル出力確率「 $b_i(w)$ 」を示す。式(14)は、「状態i」が「命令タグ集合：I」に属する場合のシンボル出力確率「 $b_i(w)$ 」を示す。

【0083】

例えば、分岐命令解析部121は、「初期状態」、「継続命令状態：S」、「データ直前命令状態：T」及び「データ状態：U」の間での遷移確率を、図7に示すように、算出する。なお、図7は、実施例1における分岐命令解析部による処理の一例について示す図である。

40

【0084】

なお、逆アセンブルされていないタグなしのオリジナルコード10について処理を実行する場合の一例について簡単に説明する。分岐命令解析部121は、タグ付きのオリジナルコード10に基づいて決定されたモデルパラメータと、タグなしのオリジナルコード10とを用いて、バウム・ウェルチアルゴリズムによって新たなモデルパラメータを決定して用いる。

50

## 【 0 0 8 5 】

次に、Forward / Backwardアルゴリズムについて簡単に説明する。分岐命令解析部 1 2 1 は、Forwardアルゴリズムに基づき、モデルパラメータ が与えられたときのオリジナルコード 1 0 候補 X の出力確率  $P(X | )$  を算出する。図 8 - 1 及び図 8 - 2、図 9 は、実施例 1 における分岐命令解析部による処理の一例について示す図である。

## 【 0 0 8 6 】

図 8 - 1 に示すような 1 6 進数表記の「入力バイナリ列」がオリジナルコード 1 0 として特定装置 1 0 0 に入力された場合を用いて説明する。また、「入力バイナリ列」を先頭から 1 バイトずつずらしながら、命令として解釈した場合の命令長を取得することで、図 8 - 2 に示すデータが得られた場合を用いて説明する。例えば、図 8 - 2 に示すように、「入力バイナリ列」が「5 5」である場合に、「命令長：1」が取得され、これに対応するニーモニックが「PUSH EBP」となる場合を用いて説明する。

## 【 0 0 8 7 】

図 9 に示す行列は、横軸に「入力バイナリ列」が配置され、縦軸に「継続命令状態：S」、「データ直前命令状態：T」及び「データ状態：U」が配置される。j 行目 i 列目の要素には、「 $x_1, \dots, x_{i-1}$ 」を出力し且つ「状態 j」で「 $x_i$  (状態 j が命令状態の場合は、 $x_i$  を命令の先頭としたときの命令全体)」を出力する「累積最大確率値」が格納される。また、各要素には、「累積最大確率値」以外にも、「遷移元要素リスト」と「累積最大確率値算出の元になった遷移元要素」が格納される。

## 【 0 0 8 8 】

ここで、各要素における「遷移元要素リスト」は、図 8 - 2 に示す命令長と、式 ( 1 1 ) ~ 式 ( 1 4 ) に示される遷移状態相関関係とを利用することで算出可能である。具体的には、図 9 に示す行列における 1 行目 1 列目 ( 継続命令状態：S ) の場合、「5 5」は、1 バイト命令であり、遷移先は、1 行目 2 列目 ( 継続命令状態：S ) と、2 行目 2 列目 ( データ直前命令状態：T ) となる。つまり、1 行目 2 列目と、2 行目 2 列目の「遷移元要素リスト」へ、1 行目 1 列目を追加する。分岐命令解析部 1 2 1 は、全要素について同様の処理を繰り返すことで、各要素における「遷移元要素リスト」が算出する。

## 【 0 0 8 9 】

なお、分岐命令解析部 1 2 1 は、すべての入力バイナリ列を出力し終える場合には、図 9 に示す行列における終了状態 ( 出力確率は「1」 ) の列に遷移するとする。なお、例外として、1 列目の要素の遷移元は、図 9 に示す行列における初期状態 ( 累積最大確率値は「1」 ) としておく。

## 【 0 0 9 0 】

なお、累積最大確率値の算出手法についても簡単に補足する。例えば、分岐命令解析部 1 2 1 は、j 行目 i 列目の遷移元要素が、n 行目 m 列目であり、n 行目 m 列目の累積最大確率値を「 $P_{nm}$ 」、「遷移元状態：n」から「現状態：j」に遷移する確率 ( 図 7 のモデルパラメータを参照 ) を「 $a_{nj}$ 」とすると、「最大確率値算出の元となった遷移元要素」は、式 ( 1 5 ) に示すように、「 $P_{nm} \times a_{nj}$ 」が最大となる「m」及び「n」を探することで算出される。そして、「 $P_{nm} \times a_{nj}$ 」の最大値に、「 $x_i$ 」 ( 現状態が命令状態の場合は、 $x_i$  を命令の先頭としたときの命令全体 ) のシンボル出力確率を乗算した値を、j 行目 i 列目の累積最大確率値として算出し、対応する要素に格納する。

## 【 0 0 9 1 】

## 【 数 1 5 】

$$\arg \max_{nm} (P_{nm} \times a_{nj}) \quad \dots (15)$$

## 【 0 0 9 2 】

図 9 に示すように、分岐命令解析部 1 2 1 は、図 9 に示す行列の要素間の遷移において、すべてのタグ系列の総和を計算する。例えば、分岐命令解析部 1 2 1 は、n 行目 m 列目の要素から j 行目 i 列目の要素への遷移 (  $m < i$  とする。 ) に対して、式 ( 1 6 ) に示す

ように、 $n$  行目  $m$  列目の確率値  $P_{nm}$  に、状態「 $n$ 」から状態「 $j$ 」に遷移する確率「 $a_{nj}$ 」（例えば、図 7 のモデルパラメータを参照）を乗算し、遷移元となるすべての  $m$ 、 $n$  について和をとる。そして、式 (16) の値に「 $x_i$ 」のシンボル出力確率を乗算した値を  $j$  行目  $i$  列目の確率値として算出し、このような計算を初期状態から終了状態まで算出して確率  $P(X|)$  を得る。

【0093】

【数16】

$$\sum_{nm} P_{nm} \times a_{nj} \quad \dots (16)$$

【0094】

分岐命令解析部 121 は、上述したようなモデルを利用し、Forward/Backward アルゴリズムを用いることで、各バイトが機械語命令である第 1 の確率を算出する。

【0095】

なお、特定装置 100 は、既知のパーソナルコンピュータ、ワークステーション、携帯電話、PHS (Personal Handyphone System) 端末、移動体通信端末又は PDA (Personal Digital Assistant) などの情報処理装置を利用して実現しても良い。例えば、PDA などの情報処理装置に、図 1 に示した記憶部 110 や制御部 120 の各機能を搭載することによって実現しても良い。

【0096】

[ 特定装置による処理 ]

図 10 は、実施例 1 における特定装置による処理の流れの一例を示すフローチャートである。

【0097】

図 10 に示すように、特定装置 100 では、オリジナルコード 10 の入力があると (ステップ S101 肯定)、分岐命令解析部 121 が、分岐命令と解釈可能な箇所を抽出する (ステップ S102)。例えば、分岐命令解析部 121 は、直接分岐命令又は間接分岐命令と解釈可能な箇所を抽出する。また、分岐命令解析部 121 は、抽出した箇所の先頭が機械語命令の先頭となる確率を算出する (ステップ S103)。

【0098】

そして、分岐関係解析部 122 は、分岐命令解析部 121 により抽出された箇所間の分岐関係を解析することで、間接分岐箇所が直接的又は間接的に分岐先となる直接分岐箇所を識別し (ステップ S104)、分岐命令解析部 121 により抽出された箇所間の分岐関係を示すツリーである暫定分岐関係ツリーを生成する (ステップ S105)。

【0099】

そして、インポートアドレス格納場所特定部 123 は、間接分岐箇所の先頭と、間接分岐箇所が分岐先となる直接分岐箇所の先頭のうちの、少なくとも一つの箇所の先頭が機械語命令の先頭となる確率を算出する (ステップ S106)。図 4 における暫定分岐先アドレス格納場所「 $0 \times 2004$ 」を例に説明すると、インポートアドレス格納場所特定部 123 は、「 $1 - (1 - 0.2) \times (1 - 0.4) \times (1 - 0.1) \times (1 - 0.6) \times (1 - 0.8) = 0.96544$ 」であると算出する。

【0100】

そして、インポートアドレス格納場所特定部 123 は、閾値以上の確率が算出されると (ステップ S107 肯定)、間接分岐箇所が間接分岐命令である場合に間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する (ステップ S108)。例えば、所定の閾値として「 $0.5$ 」を用いる場合には、インポートアドレス格納場所特定部 123 は、アドレス「 $0 \times 2004$ 」をインポートアドレス格納場所 20 として特定する。一方、インポートアドレス格納場所特定部 123 は、閾値以上の確率が算出されないと (ステップ S107 否定)、そのまま処理を終了する。

【0101】

10

20

30

40

50

なお、上記の処理手順は、上記の順番に限定されるものではなく、処理内容を矛盾させない範囲で適宜変更しても良い。例えば、上記のステップ S 1 0 3 を S 1 0 5 の後に実行しても良い。

#### 【 0 1 0 2 】

##### [ 実施例 1 の効果 ]

上述したように、実施例 1 によれば、パックされたプログラムコードをアンパックすることで得られるオリジナルコード 1 0 から、少なくとも間接分岐命令と解釈可能な箇所を示すアドレスである間接分岐箇所を抽出し、抽出された箇所が分岐命令となる第 1 の確率を算出し、算出した第 1 の確率に基づいて、抽出された該箇所の先頭が機械語命令の先頭となる第 2 の確率を算出する。そして、閾値以上の第 2 の確率が算出されると、間接分岐箇所が間接分岐命令である場合に間接分岐箇所からの分岐先を示す分岐先アドレスが格納される分岐先アドレス格納場所を特定する。この結果、適切な分岐先アドレス格納場所を特定可能となる。すなわち、逆アセンブル結果が不正確であったとしても、機械語命令の先頭となる第 1 の確率に基づいて処理を実行することで、誤ったインポートアドレス格納場所 2 0 を特定することを減らすことが可能となる。

10

#### 【 0 1 0 3 】

すなわち、実施例 1 によれば、プログラムコードにおいて、外部関数を呼び出すときに使用される分岐先アドレス格納場所を適切に特定可能となる。従来手法では、プログラムコードから間接 `call` 命令を探し出し、探し出した間接 `call` 命令が示すメモリ領域を分岐先アドレス格納場所として特定していた。この従来手法では、オリジナルコード 1 0 がマルウェアである場合など、逆アセンブルが困難で正確な逆アセンブル結果が得られない場合には、誤ったインポートアドレス格納場所 2 0 を特定することになる。これに対して、実施例 1 によれば、プログラムコード内の間接分岐命令と解釈できるすべての箇所について、機械語命令と解釈できる第 1 の確率に基づき、分岐先アドレス格納場所である確率が高い箇所を特定可能となる。この結果、従来手法よりも高精度に分岐先アドレス格納場所を特定可能となる。また、分岐先アドレス格納場所を特定できれば、例えばマルウェアが利用しようとしている外部 API を特定でき、マルウェアがどんな機能を持つか等の脅威把握や、脅威に基づくマルウェアの分類作業が可能となる。

20

#### 【 実施例 2 】

#### 【 0 1 0 4 】

さて、これまで本発明の実施例について説明したが、本発明は上述した実施例以外にも、その他の実施例にて実施されても良い。そこで、以下では、その他の実施例を示す。

30

#### 【 0 1 0 5 】

##### [ システム構成 ]

また、本実施例において説明した各処理のうち、自動的に行われるものとして説明した処理の全部又は一部を手動的に行うこともでき、あるいは、手動的に行われるものとして説明した処理の全部又は一部を公知の方法で自動的に行うこともできる。この他、上述文書中や図面中で示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報（図 1 ~ 図 1 0）については、特記する場合を除いて任意に変更することができる。

40

#### 【 0 1 0 6 】

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部又は一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的又は物理的に分散・統合して構成することができる。例えば、図 1 に示す例では、記憶部 1 1 0 を特定装置 1 0 0 の外部装置としてネットワーク経由で接続するようにしても良い。

#### 【 0 1 0 7 】

##### [ プログラム ]

図 1 1 は、特定装置による一連の処理を実行するための特定プログラムによる情報処理

50



がコンピュータを用いて具体的に実現されることを示す図である。図 11 に例示するように、コンピュータ 3000 は、例えば、メモリ 3010 と、CPU (Central Processing Unit) 3020 と、ネットワークインタフェース 3070 とを有する。コンピュータ 3000 の各部はバス 3100 によって接続される。

【0108】

メモリ 3010 は、図 11 に例示するように、ROM 3011 及び RAM 3012 を含む。ROM 3011 は、例えば、BIOS (Basic Input Output System) 等のブートプログラムを記憶する。

【0109】

ここで、図 11 に例示するように、ハードディスクドライブ 3080 は、例えば、OS 3081、アプリケーションプログラム 3082、プログラムモジュール 3083、プログラムデータ 3084 を記憶する。すなわち、開示の技術に係る特定プログラムは、コンピュータによって実行される指令が記述されたプログラムモジュール 3083 として、例えばハードディスクドライブ 3080 に記憶される。具体的には、上記実施例で説明した記憶部 110 と同様の情報処理を実行する手順各々が記述されたプログラムモジュールが、ハードディスクドライブ 3080 に記憶される。

【0110】

また、上記実施例で説明した記憶部 110 に記憶されるデータのように、特定プログラムによる情報処理に用いられるデータは、プログラムデータ 3084 として、例えばハードディスクドライブ 3080 に記憶される。そして、CPU 3020 が、ハードディスクドライブ 3080 に記憶されたプログラムモジュール 3083 やプログラムデータ 3084 を必要に応じて RAM 3012 に読み出し、各種の手順を実行する。

【0111】

なお、特定プログラムに係るプログラムモジュール 3083 やプログラムデータ 3084 は、ハードディスクドライブ 3080 に記憶される場合に限られない。例えば、プログラムモジュール 3083 やプログラムデータ 3084 は、着脱可能な記憶媒体に記憶されても良い。この場合、CPU 3020 は、ディスクドライブなどの着脱可能な記憶媒体を介してデータを読み出す。また、同様に、更新プログラムに係るプログラムモジュール 3083 やプログラムデータ 3084 は、ネットワーク (LAN (Local Area Network)、WAN (Wide Area Network) 等) を介して接続された他のコンピュータに記憶されても良い。この場合、CPU 3020 は、ネットワークインタフェースを介して他のコンピュータにアクセスすることで各種データを読み出す。

【0112】

[その他]

なお、本実施例で説明した特定プログラムは、インターネットなどのネットワークを介して配布することができる。また、特定プログラムは、ハードディスク、フレキシブルディスク (FD)、CD-ROM、MO、DVD などのコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行することもできる。

【符号の説明】

【0113】

- 100 特定装置
- 110 記憶部
- 111 暫定分岐情報格納テーブル
- 112 暫定分岐ツリー情報テーブル
- 120 制御部
- 121 分岐命令解析部
- 122 分岐関係解析部
- 123 インポートアドレス格納場所特定部

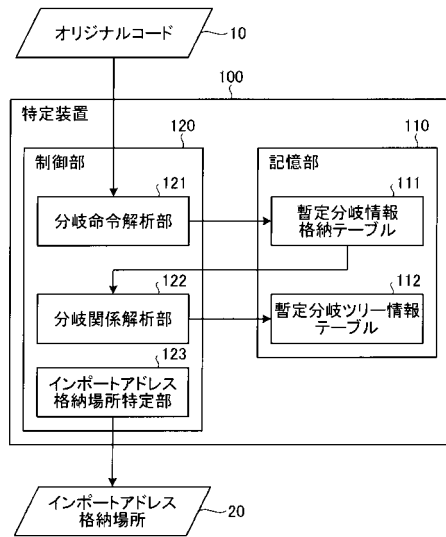
10

20

30

40

【図1】



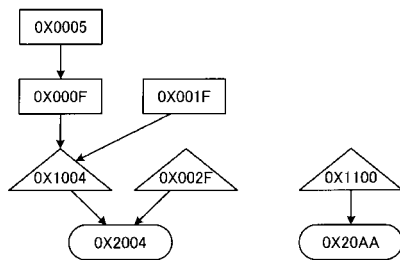
【図2】

アドレス	機械語命令として解釈したときのニーモニック
...	...
0x0005	call 0x000F
0x000F	call 0x1004
...	...
0x001F	jmp 0x1004
...	...
0x002F	call [0x2004]
...	...
0x1004	jmp [0x2004]
...	...
0x1100	call [0x20AA]
...	...

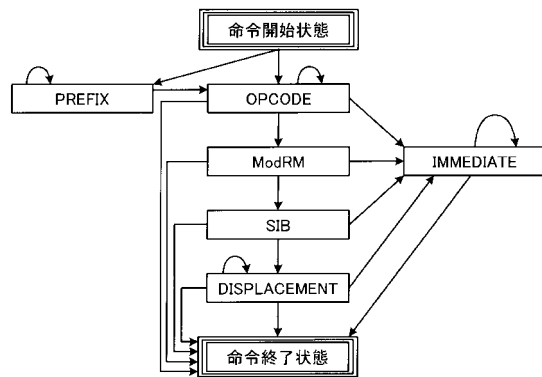
【図3】

アドレス	種別	命令確率	暫定分岐先アドレス	暫定分岐先アドレス格納場所
0x0005	直接call	0.2	0x000F	—
0x000F	直接call	0.4	0x1004	—
0x001F	直接jmp	0.6	0x1004	—
0x002F	間接call	0.8	—	0x2004
0x1004	間接jmp	0.1	—	0x2004
0x1100	間接call	0.4	—	0x20AA

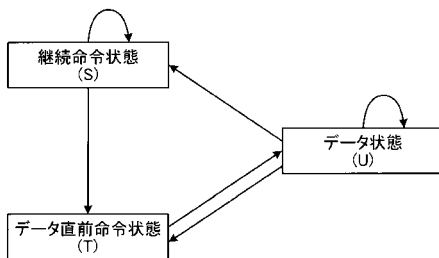
【図4】



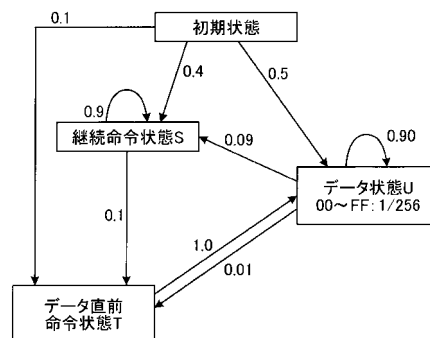
【図6】



【図5】



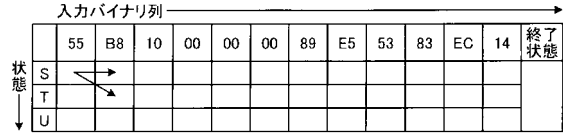
【図7】



【図 8 - 1】

55	B8	10	00	00	00	89	E5	53	83	EC	14
----	----	----	----	----	----	----	----	----	----	----	----

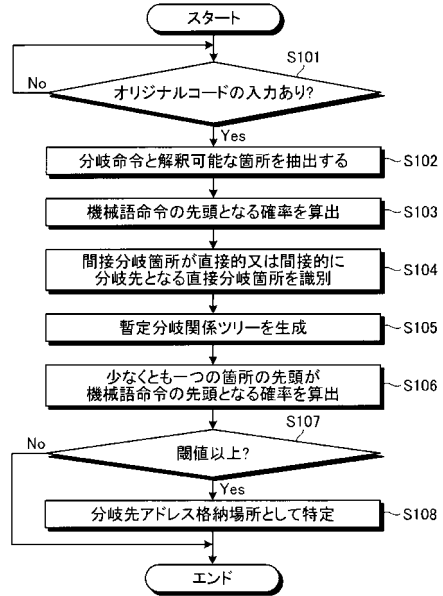
【図 9】



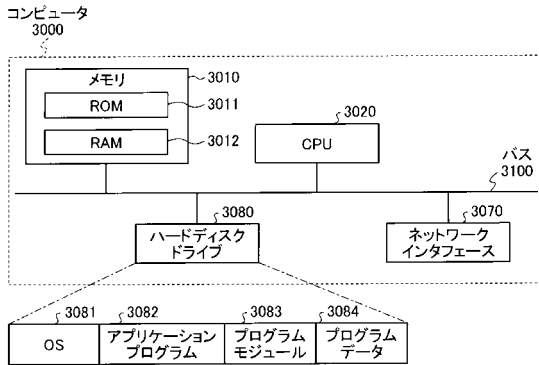
【図 8 - 2】

入力バイナリ列	命令長	(参考) ニーモニック
55	1	PUSH EBP
B8	5	MOV EAX,10
10	2	ADC BYTE PTR DS:[EAX],AL
00	2	ADD BYTE PTR DS:[EAX],AL
00	2	ADD BYTE PTR DS:[EAX],AL
00	6	ADD BYTE PTR DS:[ECX+EC8353E5],CL
89	2	MOV EBP,ESP
E5	2	IN EAX,53
53	1	PUSH EBX
83	3	SUB ESP,14
EC	1	IN AL,DX
14	-	(命令として解釈できない)

【図 10】



【図 11】



---

フロントページの続き

- (56)参考文献 特開2010-092179(JP,A)  
特開2009-193161(JP,A)  
米国特許出願公開第2010/0281540(US,A1)  
岩村 誠, マルウェアのエントリーポイント検出後におけるコード領域識別手法, 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2010年 6月10日, Vol.110, No.78, 第19-24頁  
泉田 大宗, 動的エミュレーションと静的解析を併用したバイナリコードの解析手法, 情報処理学会研究報告 平成21年度 6 [DVD-ROM], 日本, 社団法人情報処理学会, 2010年 4月15日, 第1-8頁

(58)調査した分野(Int.Cl., DB名)

G06F21/00-21/88  
G09C1/00-5/00  
H04K1/00  
H04L9/00  
G06F9/06, 9/44, -9/445, 9/48-9/50