

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4598774号  
(P4598774)

(45) 発行日 平成22年12月15日(2010.12.15)

(24) 登録日 平成22年10月1日(2010.10.1)

(51) Int.Cl.

F I

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00 6 1 0 Q

請求項の数 24 (全 21 頁)

(21) 出願番号	特願2006-533134 (P2006-533134)	(73) 特許権者	501113353
(86) (22) 出願日	平成16年5月14日(2004.5.14)		シマンテック コーポレイション
(65) 公表番号	特表2007-503660 (P2007-503660A)		Symantec Corporation
(43) 公表日	平成19年2月22日(2007.2.22)		n
(86) 国際出願番号	PCT/US2004/015383		アメリカ合衆国 カリフォルニア州 95
(87) 国際公開番号	W02004/105332		014 クーパーティノ、スティーブンス
(87) 国際公開日	平成16年12月2日(2004.12.2)		クリーク ブルバード、20330
審査請求日	平成19年5月2日(2007.5.2)		20330 Stevens Creek
(31) 優先権主張番号	60/471,242		Boulevard, Cupertino
(32) 優先日	平成15年5月15日(2003.5.15)		no, California 9501
(33) 優先権主張国	米国 (US)		4, USA
		(74) 代理人	100064621
			弁理士 山川 政樹
		(74) 代理人	100098394
			弁理士 山川 茂樹

最終頁に続く

(54) 【発明の名称】 類似性測度に基づいて電子メール・スパムをフィルタ処理するための方法および装置

(57) 【特許請求の範囲】

【請求項 1】

電子メール・メッセージを受信するための受信手段と、

前記電子メール・メッセージに、第1の文字セット内の文字の位置を示す文字参照が一つ以上含まれているか否かを検出する検出手段と、

前記検出した一つ以上の文字参照の内の少なくとも一つを前記第1の文字セット内の位置に対応する文字に変更することによって、前記電子メール・メッセージのコンテンツを修正する修正手段と、

前記電子メール・メッセージの修正されたコンテンツに基づいて、前記電子メール・メッセージを特徴付けるデータを生成するデータ生成手段と、

前記電子メール・メッセージを特徴付ける前記データを、複数のスパム・メッセージを特徴付けるデータのセットと比較する比較手段と、

前記電子メール・メッセージを特徴付ける前記データと、前記複数のスパム・メッセージを特徴付けるデータのセット内のいずれかのデータ項目の間の類似が、閾値を超えるかどうかを判定する判定手段と

を含む装置。

【請求項 2】

前記一つ以上の文字参照は、「&amp; # &lt; 数値 &gt;」の構文を持つHTML文字参照(ただし、「&lt; 数値 &gt;」は10進文字番号または16進文字番号)である請求項1に記載の装置。

【請求項 3】

10

20

スパム・メッセージを受信するための受信手段と、

前記スパム・メッセージに、第1の文字セット内の文字の位置を示す文字参照が一つ以上含まれているか否かを検出するための検出手段と、

前記検出した一つ以上の文字参照の内の少なくとも一つを前記第1の文字セット内の位置に対応する文字に変更することによって、前記スパム・メッセージのコンテンツを修正する修正手段と、

前記スパム・メッセージの修正されたコンテンツに基づいて、前記スパム・メッセージを特徴付けるデータを生成する生成手段と、

前記スパム・メッセージに類似する着信メッセージを探し出すのに後に使用される、前記スパム・メッセージを特徴付ける前記データを、サーバに転送する転送手段とを含む装置。

【請求項4】

前記第1の文字セットがASCII文字セットである請求項3に記載の装置。

【請求項5】

前記一つ以上の文字参照の内の少なくとも一つは、スパム・メッセージ内のURL内にある請求項3に記載の装置。

【請求項6】

処理システム上で実行されると、

前記処理システムの受信手段が、電子メール・メッセージを受信すること、

前記電子メール・メッセージに、第1の文字セット内の文字の位置を示す文字参照が一つ以上含まれているか否かを、前記処理システムの検出手段が検出すること、

前記検出した一つ以上の文字参照の内の少なくとも一つを前記第1の文字セット内の位置に対応する文字に変更することによって、前記電子メール・メッセージのコンテンツを、前記処理システムの修正手段が修正すること、

前記電子メール・メッセージの修正されたコンテンツに基づいて、前記電子メール・メッセージを特徴付けるデータを、前記処理システムの生成手段が生成すること、

前記処理システムの比較手段が、前記電子メール・メッセージを特徴付ける前記データを、複数のスパム・メッセージを特徴付けるデータのセットと比較すること、

前記処理システムの判定手段が、前記電子メール・メッセージを特徴付ける前記データと、前記複数のスパム・メッセージを特徴付けるデータの前のセット内のいずれかのデータ項目の間の類似が、閾値を超えるかどうかを判定すること

を含む方法を実行するようにさせる実行可能命令を含むコンピュータ可読媒体。

【請求項7】

前記第1の文字セットがASCII文字セットである請求項6に記載のコンピュータ可読媒体。

【請求項8】

前記一つ以上の文字参照は数値文字参照または文字エンティティ参照を含む請求項6に記載のコンピュータ可読媒体。

【請求項9】

処理システム上で実行されると、

前記処理システムの受信手段が、スパム・メッセージを受信すること、

前記スパム・メッセージに、第1の文字セット内の文字の位置を示す文字参照が一つ以上含まれているか否かを、前記処理システムの検出手段が検出すること、

前記処理システムの修正手段が、前記検出した一つ以上の文字参照の内の少なくとも一つを前記第1の文字セット内の位置に対応する文字に変更することによって、前記スパム・メッセージのコンテンツを修正すること、

前記処理システムの生成手段が、前記スパム・メッセージの修正されたコンテンツに基づいて、前記スパム・メッセージを特徴付けるデータを生成すること、

前記処理システムの転送手段が、前記スパム・メッセージに類似する着信メッセージを探し出すのに後に使用される、前記スパム・メッセージを特徴付ける前記データを、サー

10

20

30

40

50

バに転送すること

を含む方法を実行するようにさせる実行可能命令を含むコンピュータ可読媒体。

【請求項 1 0】

前記第 1 の文字セットが A S C I I 文字セットである請求項 9 に記載のコンピュータ可読媒体。

【請求項 1 1】

前記スパム・メッセージのコンテンツを修正することが、

前記スパム・メッセージ内の第 1 のグループの文字参照を前記第 1 の文字セット内の対応する文字に、前記修正手段が変更すること、

前記変更された第 1 のグループの文字参照が新たなグループの文字参照を形成すると  
の決定を、前記修正手段がなすこと、

前記決定に応答して、前記新たなグループの文字参照を前記第 1 の文字セット内の対応する文字に、前記修正手段が変更すること

を含む請求項 9 に記載のコンピュータ可読媒体。

【請求項 1 2】

電子メール・メッセージの中で、スパム・フィルタ処理を回避するように前記電子メール・メッセージに追加された雑音を示すデータを検出し、前記電子メール・メッセージのコンテンツを変更して、前記雑音を低減させるメッセージ・クリーニング手段であって、

前記電子メール・メッセージ内の一つ以上の文字参照を検出するための文字参照検出手段と、

前記一つ以上の文字参照の少なくとも一つを、第 1 の文字セット内の対応する文字に変更する文字参照変更手段と

を有したメッセージ・クリーニング手段と、

前記電子メール・メッセージの前記変更されたコンテンツをスパム・メッセージのコンテンツと比較する類似アイデンティファイア手段と

を含むシステム。

【請求項 1 3】

前記一つ以上の文字参照は数値文字参照を含む請求項 1 2 に記載のシステム。

【請求項 1 4】

前記メッセージ・クリーニング手段は、例外として適格でないフォーマット・データを前記電子メール・メッセージから抜き取ることによって、前記電子メール・メッセージの前記コンテンツを変更することを含む請求項 1 2 に記載のシステム。

【請求項 1 5】

前記一つ以上の文字参照は文字エンティティ参照を含む請求項 1 2 に記載のシステム。

【請求項 1 6】

前記文字参照検出手段が、変更された第 1 のグループの文字参照が新たなグループの文字参照を形成すると決定をなす手段をさらに含み、

前記文字参照変更手段が、前記決定に応答して、前記新たなグループの文字参照を前記第 1 の文字セット内の対応する文字に変更する手段をさらに含む請求項 1 2 に記載のシステム。

【請求項 1 7】

前記第 1 の文字セットが A S C I I 文字セットである請求項 1 2 に記載のシステム。

【請求項 1 8】

前記メッセージ・クリーニング手段は、URL から、事前定義されたカテゴリの一意識別子データを除去することによって、前記電子メール・メッセージの前記コンテンツを変更する手段を含む請求項 1 2 に記載のシステム。

【請求項 1 9】

前記メッセージ・クリーニング手段は、URL から、事前定義されたカテゴリのクエリ・データを除去することによって、前記電子メール・メッセージの前記コンテンツを変更する手段を含む請求項 1 2 に記載のシステム。

## 【請求項 2 0】

前記類似アイデンティファイア手段は、前記電子メール・メッセージの前記変更されたコンテンツが、前記スパム・メッセージの前記コンテンツに類似しているかどうかを判定することによって、前記電子メール・メッセージの前記変更されたコンテンツを、前記スパム・メッセージの前記コンテンツと比較する手段を含む請求項 1 2 に記載のシステム。

## 【請求項 2 1】

電子メール・メッセージの中で、スパム・フィルタ処理を回避するために前記電子メール・メッセージに追加された一つ以上の文字参照であって数値文字参照または文字エンティティ参照を含む文字参照を検出する検出手段と、

前記一つ以上の文字参照の少なくとも一つを、第 1 の文字セット内の対応する文字に変更することによって、前記電子メール・メッセージのコンテンツを変更する変更手段と、

前記電子メール・メッセージの前記変更されたコンテンツを、スパム・メッセージのコンテンツと比較する比較手段と

を含む装置。

## 【請求項 2 2】

前記第 1 の文字セットが A S C I I 文字セットである請求項 2 1 に記載の装置。

## 【請求項 2 3】

処理システム上で実行されると、

前記処理システムの検出手段が、電子メール・メッセージの中で、スパム・フィルタ処理を回避するために前記電子メール・メッセージに追加された一つ以上の文字参照であって数値文字参照または文字エンティティ参照を含む文字参照を検出すること、

前記処理システムの変更手段が、前記一つ以上の文字参照の少なくとも一つを、第 1 の文字セット内の対応する文字に変更することによって、前記電子メール・メッセージのコンテンツを変更すること、

前記処理システムの比較手段が、前記電子メール・メッセージの前記変更されたコンテンツを、スパム・メッセージのコンテンツと比較すること

を含む方法を実行するようにさせる実行可能命令を含むコンピュータ可読媒体。

## 【請求項 2 4】

前記一つ以上の文字参照の内の少なくとも一つは、前記電子メール・メッセージ内の U R L 内にある請求項 2 3 に記載のコンピュータ可読媒体。

## 【発明の詳細な説明】

## 【関連出願】

## 【0 0 0 1】

本出願は、参照により全体が本明細書に組み込まれている、2 0 0 3 年 5 月 1 5 日に出願した米国仮出願第 6 0 / 4 7 1 2 4 2 号の優先権を主張する。

## 【技術分野】

## 【0 0 0 2】

本発明は、電子メール ( e m a i l ) をフィルタ処理することに関し、より詳細には、本発明は、類似性測度に基づいて電子メール・スパムをフィルタ処理することに関する。

## 【背景技術】

## 【0 0 0 3】

インターネットの普及が進んでおり、ますます多くの人々が、電子的大量メールを生成し、送信することによって製品やサービスを宣伝して、インターネットを介してビジネスを行っている。それらの電子メッセージ ( 電子メール ) は、普通、求められておらず、受信者によって迷惑と見られている。というのは、それらのメッセージが、必要で重要なデータ処理のために必要とされる記憶スペースの多くを占有するからである。例えば、メール・サーバは、そのサーバの記憶容量が、宣伝を含む不要な電子メールで最大限度までいっぱいになった場合、重要で、かつ / または所望される電子メールを受け入れることを拒否しなければならなくなる。さらに、セットトップ・ボックス、P D A、ネットワーク・コンピュータ、ポケットベルなどのシン・クライアント・システムはすべて、限られた記

10

20

30

40

50

憶容量を有する。そのようなシステムのいずれにおける不要な電子メールも、ユーザのための有限なリソースを縛りつけることになる。さらに、通常のユーザは、大量であるが、無用な宣伝情報をダウンロードすることによって、時間を浪費する。それらの不要な電子メールは、一般にスパムと呼ばれる。

【 0 0 0 4 】

現在、不要なメッセージをフィルタにかけて除くことができる製品が存在する。例えば、すべてのスパム・エージェント（すなわち、大量の求められていない電子メールを生成する企業）の索引リストを保持し、そのリスト上の企業から送信された電子メールをブロックする手段を提供するスパム・ブロック方法が存在する。

【 0 0 0 5 】

現在、入手可能な別の「ジャンク・メール」フィルタは、事前定義された語や、前述したパターンに基づくフィルタ群を使用する。着信メールは、件名が既知のスパム・パターンを含む場合、不要なメールと指定される。

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

しかし、スパム・フィルタ処理が高度になるにつれ、フィルタを回避するスパム発信者の技術も高度になる。最近の世代のスパム発信者によって採り入れられている戦術の例には、ランダム化、発信元隠蔽、HTMLを使用したフィルタ回避が含まれる。

【 課題を解決するための手段 】

【 0 0 0 7 】

類似性測度に基づいて電子メール・スパムをフィルタ処理するための方法とシステムを説明する。一態様によれば、方法は、着信電子メール・メッセージを受信し、着信電子メール・メッセージのコンテンツに基づいて、着信電子メール・メッセージを特徴付けるデータを生成し、その生成されたデータを、スパム・メッセージを特徴付けるデータのセットと比較することを含む。方法は、着信電子メール・メッセージを特徴付けるデータと、スパム・メッセージを特徴付けるデータのセットからのいずれかのデータ項目の間の類似が、閾値を超えるかどうかを判定することをさらに含む。

【 発明を実施するための最良の形態 】

【 0 0 0 8 】

本発明のその他の特徴は、添付の図面、および以下の詳細な説明から明白となろう。

【 0 0 0 9 】

本発明は、以下に提供する詳細な説明や、本発明の様々な実施形態の添付の図面からより完全に理解されるが、説明と図面は、本発明を特定の実施形態に限定するものと解釈されるべきではなく、単に説明および理解を目的とする。

【 0 0 1 0 】

類似性測度に基づいて電子メール・スパムをフィルタ処理するための方法とシステムを説明する。以下の説明では、多数の詳細を提示する。しかし、本発明は、それらの特定の詳細なしに実施することもできることが、当業者には理解されよう。その他、周知の構造やデバイスは、本発明を不明瞭にするのを避けるため、詳細にではなく、ブロック図の形態で示す。

【 0 0 1 1 】

以下の詳細な説明のいくつかの部分は、コンピュータ・メモリ内のデータ・ビットに対するオペレーションのアルゴリズムと記号表現として提示する。それらのアルゴリズム記述やアルゴリズム表現は、データ処理分野の業者が、自らの作業の内容を他の同業者に最も効果的に伝えるのに使用する手段である。アルゴリズムは、本明細書では、また、一般に、所望の結果をもたらす、自己矛盾のない一連のステップであると考えられる。ステップは、物理的量の物理的な操作を要するステップである。通常、ただし必然的にではなく、それらの量は、格納されること、転送されること、結合されること、比較されること、それ以外の形で操作されることが可能な、電気信号または磁気信号の形態をとる。ときと

10

20

30

40

50

して、主に、一般的な用法の理由で、それらの信号をビット、値、要素、シンボル、文字、項目、数などと呼ぶのが好都合であることが分かっている。

【 0 0 1 2 】

しかし、それらの用語や類似する用語のすべては、適切な物理的量に関連付けられるべきであり、それらの量に適用された便利なラベルに過ぎないことに留意されたい。特に明記しない限り、以下の説明から明白なとおり、説明全体で、「処理する」または「演算する」または「計算する」または「算出する」または「表示する」などの用語を利用する説明は、コンピュータ・システムのレジスタ内やメモリ内の物理的（電子的）量として表されるデータを操作し、変換して、コンピュータ・システムのメモリまたはレジスタ、あるいは他のそのような情報記憶デバイス、情報伝送デバイス、または情報表示デバイスの内部の物理的量として同様に表される他のデータにする、コンピュータ・システム、または類似の電子コンピューティング・デバイスのアクションやプロセスを指すものと理解されるべきである。

10

【 0 0 1 3 】

また、本発明は、本明細書のオペレーションを実行するための装置にも関する。この装置は、要求される目的のために特別に構築されることも、コンピュータの中に格納されたコンピュータ・プログラムによって選択的に起動される、または再構成される汎用コンピュータを含むことも可能である。そのようなコンピュータ・プログラムは、フロッピー・ディスク、光ディスク、CD-ROM、光磁気ディスクを含む任意のタイプのディスク、読み取り専用メモリ（ROM）、ランダム・アクセス・メモリ（RAM）、EPROM、EEPROM、磁気カードまたは光カード、あるいは電子命令を格納するのに適した任意のタイプの媒体などの、ただし、それらには限定されない、コンピュータ可読記憶媒体の中に格納されることが可能であり、各媒体は、コンピュータ・システム・バスに結合される。

20

【 0 0 1 4 】

本明細書で提示するアルゴリズムと表示は、いずれの特定のコンピュータ、またはその他の装置にも本質的に関連していない。様々な汎用システムを、本明細書の教示によるプログラムとともに使用することができ、あるいは、要求される方法ステップを実行する、より特殊化された装置を構築することが好都合であると判明する。様々なそれらのシステムの要求される構造は、以下の説明から明らかとなる。さらに、本発明は、いずれの特定のプログラミング言語に関連しても説明しない。様々なプログラミング言語を使用して、本明細書で説明する本発明の教示を実施することができることが理解されよう。

30

【 0 0 1 5 】

マシン可読媒体には、マシン（例えば、コンピュータ）による読み取りが可能な形態で情報を格納する、または伝送するための任意の機構が含まれる。例えば、マシン可読媒体には、読み取り専用メモリ（「ROM」）、ランダム・アクセス・メモリ（「RAM」）、磁気ディスク記憶媒体、光記憶媒体、フラッシュメモリ・デバイス、電気、光、音響、またはその他の形態の伝搬される信号（例えば、搬送波、赤外線信号、デジタル信号など）、その他が含まれる。

【 0 0 1 6 】

40

類似性測度に基づいて電子メール・スパムをフィルタ処理すること

図1は、スパム電子メール（email）の配信を制御するためのシステムの一実施形態のブロック図である。システムは、公共ネットワーク（例えば、インターネット、無線ネットワークなど）、または私設ネットワーク（例えば、LAN、イントラネットなど）などの通信ネットワーク100に結合された制御センタ102を含む。制御センタ102は、ネットワーク100を介して複数のネットワーク・サーバ104と通信する。各サーバ104は、私設ネットワークまたは公共ネットワークを使用して、ユーザ端末装置106と通信する。

【 0 0 1 7 】

制御センタ102は、スパムと識別されたメッセージを分析し、スパムを検出するため

50

のフィルタ処理規則を開発し、フィルタ処理規則をサーバ群 104 に配信することを担うスパム対策設備である。メッセージは、既知のスパム源（例えば、「スパム・プローブ」、すなわち、可能な限り多くのスパム発信者メーリング・リストに入り込むように特別に選択された電子メール・アドレスを使用して特定された）によって送信されれば、スパムとして識別される。

【0018】

サーバ 104 は、送信された対応するユーザ端末装置のユーザにアドレス指定されたメッセージを受信し、格納するメール・サーバである。代替として、サーバ 104 は、メール・サーバ 104 に結合された、異なるサーバであってもよい。サーバ群 104 は、制御センタ 102 から受け取られたフィルタ処理規則に基づいて、着信メッセージをフィルタ

10

【0019】

ー実施形態では、制御センタ 102 は、スパム攻撃に関連するコンテンツを特徴付けるデータを生成し、そのデータをサーバ群 104 に送信することを担うスパム・コンテンツ準備モジュール 108 を含む。各サーバ 104 は、制御センタ 102 から受信されたスパム・データを格納し、格納されたデータを使用して、スパム・コンテンツに類似する着信電子メール・メッセージを識別することを担う、類似度算出モジュール 110 を含む。

【0020】

代替の実施形態では、各サーバ 104 は、スパム攻撃に関連するコンテンツを特徴付けるデータを生成するスパム・コンテンツ準備モジュール 108 と、生成されたデータを使用して、スパム・コンテンツに類似する電子メール・メッセージを識別する類似度算出モジュール 110 をともにホストする。

20

【0021】

図 2 は、スパム・コンテンツ準備モジュール 200 の一実施形態のブロック図である。スパム・コンテンツ準備モジュール 200 は、スパム・コンテンツ・パーサ 202、スパム・データ・ジェネレータ 206、スパム・データ・トランスミッタ 208 を含む。

【0022】

スパム・コンテンツ・パーサ 202 は、スパム攻撃による電子メール・メッセージの本文（スパム・メッセージと呼ぶ）を解析することを担う。

【0023】

スパム・データ・ジェネレータ 206 は、スパム・メッセージを特徴付けるデータを生成することを担う。ー実施形態では、スパム・メッセージを特徴付けるデータには、スパム・メッセージを構成するトークン（例えば、文字、語、行など）のセットに関して計算されたハッシュ値のリストが含まれる。スパム・メッセージ、または他の任意の電子メール・メッセージを特徴付けるデータを本明細書では、メッセージ・シグネチャと呼ぶ。スパム・メッセージ、または他の任意の電子メール・メッセージのシグネチャは、メッセージ・コンテンツを識別する様々なデータを含むことが可能であり、異なる電子メール・メッセージのシグネチャを比較する際に類似性測度の使用を可能にする、様々なアルゴリズムを使用して作成される。

30

【0024】

ー実施形態では、スパム・コンテンツ準備モジュール 200 は、雑音を示すデータを検出すること、さらに、スパム・メッセージのシグネチャを生成するのに先立って、スパム・メッセージからその雑音を除去することを担う雑音低減アルゴリズム 204 も含む。雑音は、スパムの性質を隠すようにスパム・メッセージに追加されている、受信者には見えないデータを表す。

40

【0025】

ー実施形態では、スパム・コンテンツ準備モジュール 200 は、単一のスパム攻撃を元とするメッセージをグループ化することを担うメッセージ・グループ化アルゴリズム（図示せず）も含む。グループ化は、スパム・メッセージの指定された特性（例えば、含まれる URL、メッセージ部分など）に基づいて実行される。グループ化が使用される場合、

50

スパム・データ・ジェネレータ 206 は、それぞれの個別メッセージに関してではなく、スパム・メッセージのグループに関するシグネチャを生成することができる。

【0026】

スパム・データ・トランスミッタ 208 は、スパム・メッセージのシグネチャを、図 1 のサーバ群 104 のような、参加するサーバ群に配信することを担う。一実施形態では、各サーバ 104 は、コール・センタ 102 に対する接続（例えば、セキュリティで保護された HTTP S 接続）を定期的に（例えば、5 分毎に）開始する。このプル・ベースの接続を使用して、シグネチャは、コール・センタ 102 から妥当なサーバ 106 に伝送される。

【0027】

図 3 は、類似度算出モジュール 300 の一実施形態のブロック図である。類似度算出モジュール 300 は、着信メッセージ・パーサ 302 と、スパム・データ・レシーバ 306 と、メッセージ・データ・ジェネレータ 310 と、類似アイデンティファイア 312 と、スパム・データベース 304 とを含む。

【0028】

着信メッセージ・パーサ 302 は、着信電子メール・メッセージの本文を解析することを担う。

【0029】

スパム・データ・レシーバ 306 は、スパム・メッセージのシグネチャを受信して、そのシグネチャをスパム・データベース 304 の中に格納することを担う。

【0030】

メッセージ・データ・ジェネレータ 310 は、着信電子メール・メッセージのシグネチャを生成することを担う。一部の実施形態では、着信電子メール・メッセージのシグネチャは、その着信電子メール・メッセージを構成するトークン（例えば、文字、語、行など）のセットに関して計算されたハッシュ値のリストを含む。他の諸実施形態では、着信電子メール・メッセージのシグネチャは、電子メール・メッセージのコンテンツ（例えば、着信電子メール・メッセージを構成するトークン・セットのサブセット）を特徴付ける、他の様々なデータを含む。前述したとおり、電子メール・メッセージのシグネチャは、異なる電子メール・メッセージのシグネチャを比較する際に類似性測度の使用を可能にする、様々なアルゴリズムを使用して作成される。

【0031】

一実施形態では、類似度算出モジュール 300 は、以下により詳細に説明するとおり、雑音を示すデータを検出することと、着信電子メール・メッセージから、そのメッセージのシグネチャを生成するのに先立って、その雑音を除去することを担う、着信メッセージ・クリーニング・アルゴリズム 308 も含む。

【0032】

類似アイデンティファイア 312 は、各着信電子メール・メッセージのシグネチャを、スパム・データベース 304 の中に格納されたスパム・メッセージのシグネチャと比較することと、その比較に基づき、着信電子メール・メッセージが、何らかのスパム・メッセージに類似しているかどうかを判定することを担う。

【0033】

一実施形態では、スパム・データベース 304 は、スパム・メッセージが雑音低減プロセスを受ける前のスパム・メッセージ（すなわち、雑音のあるスパム・メッセージ）に対して生成されたシグネチャと、スパム・メッセージが雑音低減プロセスを受けた後のスパム・メッセージ（すなわち、低減された雑音を有するスパム・メッセージ）に対して生成されたシグネチャを格納する。この実施形態では、メッセージ・データ・ジェネレータ 310 がまず、雑音低減に先立って、着信電子メール・メッセージのシグネチャを生成し、類似アイデンティファイア 312 が、そのシグネチャを、雑音のあるスパム・メッセージのシグネチャと比較する。この比較により、着信電子メール・メッセージが、それらのスパム・メッセージの 1 つに類似していることが示された場合、類似アイデンティファイア

10

20

30

40

50



3 1 2 が、その着信電子メール・メッセージにスパムとしてマークを付ける。代替として、類似アイデンティファイア 3 1 2 は、着信メッセージ・クリーニング・アルゴリズム 3 0 8 を呼び出して、着信電子メール・メッセージから雑音を除去する。次に、メッセージ・データ・ジェネレータ 3 1 0 が、その変更された着信メッセージに関するシグネチャを生成し、このシグネチャが、次に、類似アイデンティファイア 3 1 2 によって、低減された雑音を有するスパム・メッセージのシグネチャと比較される。

【 0 0 3 4 】

図 4 は、スパム・メッセージを扱うためのプロセス 4 0 0 の一実施形態を示す流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。一実施形態では、処理論理は、図 1 の制御センタ 1 0 2 に存在する。

10

【 0 0 3 5 】

図 4 を参照すると、プロセス 4 0 0 は、処理論理が、スパム・メッセージを受け取ることから始まる（処理ブロック 4 0 2 ）。

【 0 0 3 6 】

処理ブロック 4 0 4 で、処理論理は、スパム・メッセージを変更して、雑音を低減させる。雑音低減アルゴリズムの一実施形態は、図 9 と図 1 0 に関連して、以下により詳細に説明する。

【 0 0 3 7 】

20

処理ブロック 4 0 6 で、処理論理はスパム・メッセージのシグネチャを生成する。一実施形態では、スパム・メッセージのシグネチャは、図 6 A に関連して以下により詳細に説明するとおり、着信電子メール・メッセージを構成するトークン（例えば、文字、語、行など）のセットに関して計算されたハッシュ値のリストを含む。他の諸実施形態では、着信電子メール・メッセージのシグネチャは、電子メール・メッセージのコンテンツを特徴付ける、他の様々なデータを含む。

【 0 0 3 8 】

処理ブロック 4 0 8 で、処理論理は、スパム・メッセージのシグネチャをサーバ（例えば、図 1 のサーバ 1 0 4 ）に転送し、サーバは、スパム・メッセージのそのシグネチャを使用して、そのスパム・メッセージに類似する着信電子メール・メッセージを探し出す（ブロック 4 1 0 ）。

30

【 0 0 3 9 】

図 5 は、類似性測度に基づいて電子メール・スパムをフィルタ処理するためのプロセス 5 0 0 の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。一実施形態では、処理論理は、図 1 のサーバ 1 0 4 に存在する。

【 0 0 4 0 】

図 5 を参照すると、プロセス 5 0 0 は、処理論理が着信電子メール・メッセージを受け取ることから始まる（処理ブロック 5 0 2 ）。

40

【 0 0 4 1 】

処理ブロック 5 0 4 で、処理論理は、着信メッセージを変更して、雑音を低減させる。雑音低減アルゴリズムの一実施形態は、図 9 と図 1 0 に関連して、以下により詳細に説明する。

【 0 0 4 2 】

処理ブロック 5 0 6 で、処理論理は、着信メッセージのコンテンツに基づき、着信メッセージのシグネチャを生成する。一実施形態では、着信電子メール・メッセージのシグネチャは、図 6 A に関連して以下により詳細に説明するとおり、着信電子メール・メッセージを構成するトークン（例えば、文字、語、行など）のセットに関して計算されたハッシュ値のリストを含む。他の諸実施形態では、着信電子メール・メッセージのシグネチャは

50

、電子メール・メッセージのコンテンツを特徴付ける、他の様々なデータを含む。

【0043】

処理ブロック508で、処理は、着信メッセージのシグネチャをスパム・メッセージのシグネチャと比較する。

【0044】

処理ブロック510で、着信メッセージのシグネチャと何らかのスパム・メッセージのシグネチャの間の類似が、閾値類似性測度を超えている、と処理論理が判定する。2つのメッセージ間の類似を判定するためのプロセスの一実施形態は、図6Bに関連して以下により詳細に説明する。

【0045】

処理ブロック512で、処理論理は、着信電子メール・メッセージにスパムとしてマークを付ける。

【0046】

図6Aは、電子メール・メッセージのシグネチャを作成するためのプロセス600の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。一実施形態では、処理論理は、図1のサーバ104に存在する。

【0047】

図6Aを参照すると、プロセス600は、処理論理が、電子メール・メッセージをトークン・セットに分割することから始まる（処理ブロック602）。各トークン・セットは、電子メール・メッセージからの事前定義された数の順次ユニットを含む。事前定義された数は、1以上である。ユニットは、電子メール・メッセージの中の文字、語、または行を表わす。一実施形態では、各トークン・セットは、電子メール・メッセージの中のそのトークン・セットの出現の回数と組み合わせられる。

【0048】

処理ブロック604で、処理論理は、それらのトークン・セットに関するハッシュ値を計算する。一実施形態では、ハッシュ値は、トークン・セットと対応するトークン出現回数との各組み合わせにハッシュ関数を適用することによって計算される。

【0049】

処理ブロック606で、処理論理は、計算されたハッシュ値を使用して、電子メール・メッセージに関するシグネチャを作成する。一実施形態では、シグネチャは、計算されたハッシュ値のサブセットを選択し、電子メール・メッセージを特徴付けるパラメータを、計算されたハッシュ値の選択されたサブセットに加えることによって作成される。パラメータは、例えば、電子メール・メッセージのサイズ、計算されるハッシュ値の数、電子メール・メッセージに関連するキーワード、添付ファイルの名前などを指定することが可能である。

【0050】

一実施形態では、電子メール・メッセージに関するシグネチャは、図7と図8に関連して以下により詳細に説明する、文字ベースのドキュメント比較機構を使用して作成される。

【0051】

図6Bは、電子メール・メッセージのシグネチャを使用してスパムを検出するためのプロセス650の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。一実施形態では、処理論理は、図1のサーバ104に存在する。

【0052】

図6Bを参照すると、プロセス650は、着信電子メール・メッセージのシグネチャの

10

20

30

40

50

中のデータを、各スパム・メッセージのシグネチャの中のデータと比較する。シグネチャ・データは、電子メール・メッセージのコンテンツを特徴付けるパラメータと、電子メール・メッセージの中に含まれるトークンに関して生成されたハッシュ値のサブセットを含む。パラメータは、例えば、電子メール・メッセージのサイズ、電子メール・メッセージの中のトークンの数、電子メール・メッセージに関連するキーワード、添付ファイルの名前などを指定する。

【 0 0 5 3 】

処理論理は、着信電子メール・メッセージのシグネチャの中のパラメータを、各スパム・メッセージのシグネチャの中の対応するパラメータと比較することから始まる（処理ブロック 6 5 2 ）。

10

【 0 0 5 4 】

判定ボックス 6 5 4 で、処理論理は、いずれかのスパム・メッセージ・シグネチャが、着信メッセージ・シグネチャのパラメータに類似するパラメータを含むかどうかを判定する。類似性は、例えば、2つのパラメータ間の許容される差、または2つのパラメータの許容される比に基づいて判定される。

【 0 0 5 5 】

スパム・メッセージ・シグネチャのいずれも、着信メッセージ・シグネチャのパラメータに類似するパラメータを含まない場合、処理論理は、その着信電子メール・メッセージが正当である（すなわち、スパムではない）と判定する（処理ブロック 6 6 2 ）。

20

【 0 0 5 6 】

代替として、1つまたは複数のスパム・メッセージ・シグネチャが、類似するパラメータを有する場合、処理論理は、最初のスパム・メッセージのシグネチャが、着信電子メールのシグネチャの中のハッシュ値に類似するハッシュ値を有するかどうかを判定する（判定ボックス 6 5 6 ）。類似性閾値に基づき、ハッシュ値は、例えば、ある数のハッシュ値が一致する場合、または一致するハッシュ値と一致しないハッシュ値の比が、指定された閾値を超えた場合、類似していると見る。

【 0 0 5 7 】

最初のスパム・メッセージ・シグネチャが、着信電子メール・シグネチャのハッシュ値に類似するハッシュ値を有する場合、処理論理は、着信電子メールはスパムであると判定する（処理ブロック 6 7 0 ）。それ以外の場合、処理論理は、類似するパラメータを有する、さらなるスパム・メッセージ・シグネチャが存在するかどうかをさらに判定する（判定ブロック 6 5 8 ）。存在する場合、処理論理は、次のスパム・メッセージ・シグネチャが、着信電子メール・シグネチャのハッシュ値に類似するハッシュ値を有するかどうかを判定する（判定ボックス 6 5 6 ）。類似するハッシュ値を有する場合、処理論理は、その着信電子メール・メッセージはスパムであると判定する（処理ブロック 6 7 0 ）。類似するハッシュ値を有さない場合、処理論理は、処理ブロック 6 5 8 に戻る。

30

【 0 0 5 8 】

処理論理は、類似するパラメータを有する他のスパム・メッセージ・シグネチャが全く存在しないと判定した場合、その着信メール・メッセージはスパムではないと判定する（処理ブロック 6 6 2 ）。

40

【 0 0 5 9 】

文字ベースのドキュメント比較機構

図 7 は、ドキュメントの文字ベースの比較のためのプロセス 7 0 0 の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。

【 0 0 6 0 】

図 7 を参照すると、プロセス 7 0 0 は、処理論理がドキュメントを前処理することから始まる（処理ブロック 7 0 2 ）。一実施形態では、ドキュメントは、ドキュメント内のそれぞれの大文字の英字を小文字の英字に変更することによって前処理される。例えば、「

50

I am Sam, Sam I am」というメッセージが前処理されて、「i . a m . s a m . s a m . i . a m」という表現になる。

【 0 0 6 1 】

処理ブロック704で、処理論理は、ドキュメントをトークンに分割し、各トークンが、ドキュメントからの所定の数の順次の文字を含む。一実施形態では、各トークンは、そのトークンの出現回数と組み合わせられる。この組み合わせは、ラベル付きシングル ( s h i n g l e ) と呼ばれる。例えば、トークン内の順次の文字の所定の数が3に等しい場合、上記に規定した表現は、ラベル付きシングルの以下のセットを含む。

【 0 0 6 2 】

【表1】

10

i.al

.aml

am.l

m.sl

.sal

saml

sm.2

m.sl

.sm2

sam2

am.3

m.il

.i.l

i.a2

.am4

20

30

【 0 0 6 3 】

一実施形態では、シングルは、ヒストグラムとして表される。

【 0 0 6 4 】

処理ブロック706で、処理論理は、トークンに関するハッシュ値を計算する。一実施形態では、ハッシュ値は、ラベル付きシングルに関して計算される。例えば、上記の各ラベル付きシングルにハッシュ関数  $H(x)$  が適用された場合、以下の結果がもたらされる。すなわち、

【 0 0 6 5 】

40

## 【表 2】

H(i.a1) -> 458348732	
H(.am1) -> 200404023	
H(am.1) -> 692939349	
H(m.s1) -> 220443033	
H(.sa1) -> 554034022	
H(8am1) -> 542929292	10
H(am.2) -> 629292229	
H(m.s1) -> 702202232	
H(.sa2) -> 322243349	
H(8am2) -> 993923828	
H(am.3) -> 163393269	
H(m.i1) -> 595437753	
H(.i.1) -> 843438583	20
H(i.a2) -> 244485639	
H(.am4) -> 493869359	

## 【 0 0 6 6 】

－実施形態では、処理論理は、ハッシュ値を以下のとおり並べ替える。すなわち、

## 【 0 0 6 7 】

## 【表 3】

163393269	
200604023	30
220643033	
246685639	
322263369	
458368732	
493869359	
542929292	
554034022	40
595637753	
629292229	
692939349	
702202232	
843438583	
993923828	

## 【 0 0 6 8 】

処理ブロック 7 0 8 で、処理論理は、計算されたハッシュ値からハッシュ値のサブセッ 50

トを選択する。一実施形態では、処理論理は、並べ替えられたハッシュ値から、小さい方からX個の値を選択し、それらの値から、ドキュメントの「スケッチ」を作成する。例えば、X = 4 の場合、スケッチは、以下のとおり表現される。すなわち、

[163393269 200404023 220443033 244485639]

【0069】

処理ブロック710で、処理論理は、スケッチに、ドキュメントのトークンに関するパラメータを加えることによって、ドキュメントのシグネチャを作成する。一実施形態では、パラメータは、ドキュメントの中における最初のトークンの数を指定する。前述の例では、最初のトークンの数は、15である。したがって、ドキュメントのシグネチャは、以下のとおり表現される。すなわち、

[15 163393269 200404023 220443033 244485639]

代替として、パラメータは、ドキュメントのコンテンツの他の任意の特性（例えば、ドキュメントのサイズ、ドキュメントに関連するキーワードなど）を指定してもよい。

【0070】

図8は、2つのドキュメントが類似しているかどうかを判定するためのプロセス800の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。

【0071】

図8を参照すると、プロセス800は、処理論理が、ドキュメント1のシグネチャの中で指定されたトークン数と、ドキュメント2のシグネチャの中で指定されたトークン数を比較して、第1のシグネチャの中のトークン数が、第2のシグネチャからのトークン数に対して、許容される範囲内にあるかどうかを判定することから始まる（判定ブロック802）。例えば、許容される範囲は、1つ以内の違い、または90パーセント以上の比である。

【0072】

第1のシグネチャの中のトークン数が、第2のシグネチャからのトークン数に対して、許容される範囲外である場合、処理論理は、ドキュメント1とドキュメント2は異なると判定する（処理ブロック808）。そうではなく、第1のシグネチャの中のトークン数が、第2のシグネチャからのトークン数に対して、許容される範囲内にある場合、処理論理は、シグネチャ1の中のハッシュ値と、シグネチャ2の中のハッシュ値の間の類似が閾値を超えている（例えば、95パーセントを超えるハッシュ値が同一である）かどうかを判定する（判定ボックス804）。超えている場合、処理論理は、2つのドキュメントは類似していると判定する（処理ブロック806）。超えていない場合、処理論理は、ドキュメント1とドキュメント2は異なると判定する（処理ブロック808）。

【0073】

雑音低減を使用する電子メールスパムフィルタ処理

図9は、電子メール・メッセージの中の雑音を低減するためのプロセス900の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。

【0074】

図9を参照すると、プロセス900は、処理論理が、電子メール・メッセージの中で、雑音を示すデータを検出することから始まる（処理ブロック902）。前述したとおり、雑音は、メール・メッセージの受信者に見えず、スパム・フィルタ処理を回避するために電子メール・メッセージに追加されているデータを表す。そのようなデータには、例えば、フォーマット・データ（例えば、HTMLタグ）、数値文字参照、文字エンティティ参照、事前定義されたカテゴリのURLデータなどが含まれる。数値文字参照は、ドキュメ

10

20

30

40

50

ント文字セットのなかの文字のコード位置を指定する。文字エンティティ参照は、シンボリック名を使用して、作成者が、コード位置を憶えている必要がないようにする。例えば、&a r i n g という文字エンティティ参照は、r i n g リングの上に置かれた小文字「a」を指す。

【0075】

処理ブロック904で、処理論理は、電子メール・メッセージのコンテンツを変更して、雑音を低減する。一実施形態では、コンテンツ変更には、フォーマット・データを除去することや、数値文字参照や文字エンティティ参照をASCII等価物に変換すること、URLデータを変更することが含まれる。

【0076】

処理ブロック906で、処理論理は、電子メール・メッセージの変更されたコンテンツをスパム・メッセージのコンテンツと比較する。一実施形態では、比較は、厳密な一致を識別するように実行される。代替として、比較は、2つのドキュメントが類似しているかどうかを判定するように実行される。

【0077】

図10は、電子メール・メッセージを変更して、雑音を低減させるためのプロセス1000の一実施形態の流れ図である。プロセスは、ハードウェア（例えば、専用論理、プログラマブル論理、マイクロコードなど）、ソフトウェア（汎用コンピュータ・システム上、または専用マシン上で実行されるような）、またはそれらの組み合わせを含む処理論理によって実行される。

【0078】

図10を参照すると、プロセス1000は、処理論理が、フォーマット・データ（例えば、HTMLタグ）を求めて、電子メール・メッセージを調べることから始まる（処理ブロック1002）。

【0079】

判定ボックス1004で、処理論理は、見つかったフォーマット・データが、例外として適格であるかどうかを判定する。通常、HTMLフォーマットは、メッセージの情報コンテンツに何も追加しない。しかし、いくつかの例外が存在する。それらの例外は、メッセージのさらなる処理のための有用な情報を含むタグ（例えば、<B O D Y>、<A>、<I M G>、<F O N T>などのタグ）である。例えば、<F O N T>タグと<B O D Y>タグは、「白地に白(white on white)」のテキストを無くすために必要とされ、<A>タグと<I M G>タグは、通常、データをシステムの他のコンポーネントに渡すために使用されるリンク情報を含む。

【0080】

フォーマット・データが、例外として適格でない場合、そのフォーマット・データは、電子メール・メッセージから抜き取られる（処理ブロック1006）。

【0081】

次に、処理論理は、各数値文字参照と各文字エンティティ参照を、対応するASCII文字に変換する（処理ブロック1008）。

【0082】

HTMLでは、数値文字参照は、以下の2つの形態をとる。すなわち、  
1. Dが10進数である、「&#D;」という構文が、ISO10646の10進文字番号Dを指し、  
2. Hが16進数である「&#xH;」または「&#XH;」という構文が、ISO10646の16進文字番号Hを指す。数値文字参照における16進数字は、大文字と小文字の区別がない。

【0083】

例えば、本文中のランダム化された文字は、以下の表現のとおり現れる。すなわち、

【0084】

10

20

30

40

## 【表 4】

Th&#101&#32&#83a&#118&#105n&#103&#115R&#101&#103is  
&#116e&#114&#119&#97&#110&#116&#115&#32yo&#117.

【 0 0 8 5 】

この表現は、「The Savings Register wants you」という句の意味を有する。

【 0 0 8 6 】

ときとして、処理ブロック 1008 において実行される変換は、繰り返される必要がある。例えば、「& # 3 8 ; 」という文字列は、A S C I I において「& 」という文字列に対応し、「& # 3 5 ; 」という文字列は、A S C I I において「# 」という文字列に対応し、「& # 5 1 ; 」という文字列は、A S C I I において 3 に対応し、「# 5 6 ; 」という文字列は、A S C I I において 8 に対応し、「# 5 9 ; 」は、A S C I I において「 ; 」という文字列に対応する。このため、結合された文字列、「& # 3 8 ; & # 3 5 ; & # 5 1 ; # 5 6 ; # 5 9 ; 」は、変換されると、変換される必要がある「& # 3 8 ; 」という文字列をもたらす。

【 0 0 8 7 】

したがって、処理ブロック 1008 における最初の変換オペレーションの後、処理論理は、変換済みのデータが、数値文字参照または文字エンティティ参照を依然として含むかどうかを調べる（判定ボックス 1010）。検査が肯定的であった場合、処理論理は、処理ブロック 1008 において変換オペレーションを繰り返す。肯定的ではなかった場合、処理論理は、処理ブロック 1012 に進む。

【 0 0 8 8 】

処理ブロック 1 0 1 2 で、処理論理は、事前定義されたカテゴリの URL データを変更する。それらのカテゴリには、例えば、処理論理によって、対応する ASCII 文字に変換される、URL の中に含まれる数値文字参照が含まれる。さらに、URL 「パスワード」構文を使用して、URL ホスト名の中で「@」の前に文字を追加することができる。それらの文字は、目標の Web サーバによって無視されるが、相当な量の雑音情報を各 URL に追加する。処理論理は、それらの追加の文字を除去することによって URL データを変更する。最後に、処理論理は、URL の終端における「?」という文字列の後に続く、URL の「クエリ」部分を除去する。

【 0 0 8 9 】

URLの例は、以下のとおりである。すなわち、  
`http%3a%2f%2f1otsofjunk@www.foo.com%2fbar.html?muchmorejunk`  
 処理論理は、上記のURLデータを`http://www.foo.com/bar.html`に変更する。

【 0 0 9 0 】

## 例示的なコンピュータ・システム

図 1 1 は、本明細書で説明するオペレーションの 1 つまたは複数を実行するのに使用することができる例示的なコンピュータ・システム 1 1 0 0 のブロック図である。代替の実施形態では、マシンは、ネットワーク・ルータ、ネットワーク・スイッチ、ネットワーク・ブリッジ、パーソナル・デジタル・アシスタント（PDA）、セルラー電話、Web 機器、あるいはそのマシンによって行われるアクションを指定する一連の命令を実行することができる任意のマシンを含む。

【 0 0 9 1 】

コンピュータ・システム 1100 は、バス 1108 を介して互いに通信する、プロセッサ 1102、メイン・メモリ 1104、スタティック・メモリ 1106 を含む。コンピュータ・システム 1100 は、ビデオ・ディスプレイ・ユニット 1110（例えば、液晶デ



イスプレイ（LCD）または陰極線管（CRT））をさらに含むことが可能である。また、コンピュータ・システム 1100 は、英数字入力デバイス 1112（例えば、キーボード）、カーソル制御デバイス 1114（例えば、マウス）、ディスク・ドライブ・ユニット 1116、信号生成デバイス 1120（例えば、スピーカ）、ネットワーク・インタフェース・デバイス 1122 も含む。

【0092】

ディスク・ドライブ・ユニット 1116 は、前述した方法のいずれか 1 つ、またはすべてを実施する命令セット（すなわち、ソフトウェア）1126 が格納されているコンピュータ可読媒体 1124 を含む。ソフトウェア 1126 は、メイン・メモリ 1104 内、および/またはプロセッサ 1102 内にも、完全に、または部分的に存在するように図示されている。ソフトウェア 1126 は、ネットワーク・インタフェース・デバイス 1122 を介して、さらに送信、または受信される。本明細書では、「コンピュータ可読媒体」という用語は、コンピュータによる実行のための、コンピュータが、本発明の方法のいずれか 1 つを実行するようにさせる一連の命令を格納する、または符号化することができる任意の媒体を含むと解釈される。したがって、「コンピュータ可読媒体」には、ソリッドステート・メモリ、光ディスクや磁気ディスク、搬送波信号が含まれるが、これらに限定されない。

【0093】

以上の説明を読んだ後、本発明の多くの代替形態および変更形態が当業者には明白となるに違いないが、例示として図示し、説明したいずれの特定の実施形態も、限定するものと見なされることは全く意図していないことを理解されたい。したがって、様々な実施形態の詳細についての言及は、本発明に不可欠であると見なされる特徴だけを記載する特許請求の範囲を限定することを意図していない。

【図面の簡単な説明】

【0094】

【図 1】スパム電子メールの配信を制御するためのシステムの一実施形態を示すブロック図である。

【図 2】スパム・コンテンツ準備モジュールの一実施形態を示すブロック図である。

【図 3】類似度算出モジュールの一実施形態を示すブロック図である。

【図 4】スパム・メッセージを扱うためのプロセスの一実施形態を示す流れ図である。

【図 5】類似点測度に基づいて電子メール・スパムをフィルタ処理するためのプロセスの一実施形態を示す流れ図である。

【図 6 A】電子メール・メッセージのシグネチャを作成するためのプロセスの一実施形態を示す流れ図である。

【図 6 B】電子メール・メッセージのシグネチャを使用してスパムを検出するためのプロセスの一実施形態を示す流れ図である。

【図 7】ドキュメントの文字ベースの比較のためのプロセスの一実施形態を示す流れ図である。

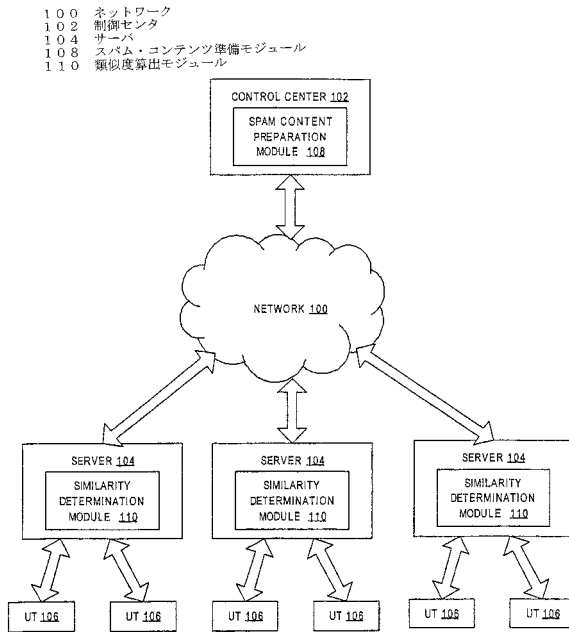
【図 8】2 つのドキュメントが類似しているかどうかを判定するためのプロセスの一実施形態を示す流れ図である。

【図 9】電子メール・メッセージの中の雑音を低減するためのプロセスの一実施形態を示す流れ図である。

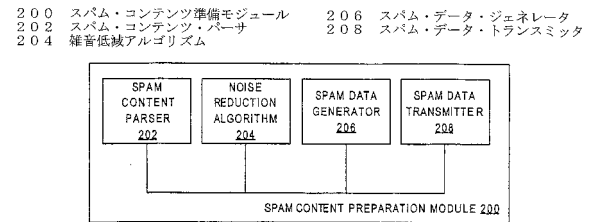
【図 10】電子メール・メッセージを変更して雑音を低減するためのプロセスの一実施形態を示す流れ図である。

【図 11】例示的なコンピュータ・システムを示すブロック図である。

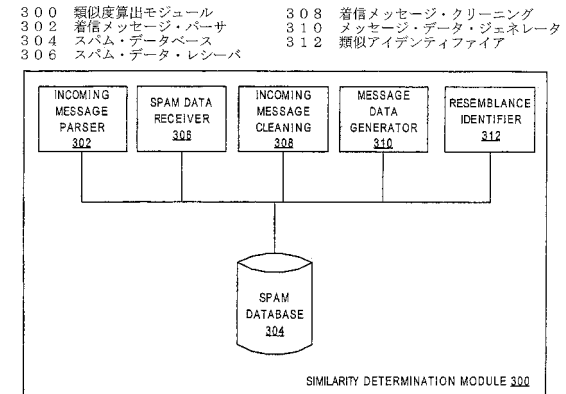
【図 1】



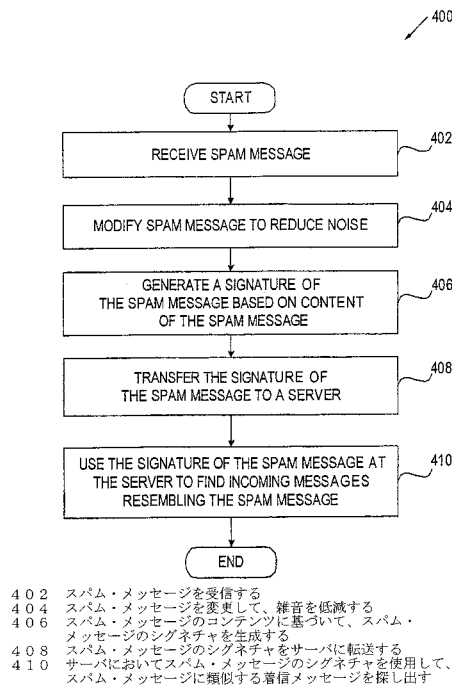
【図 2】



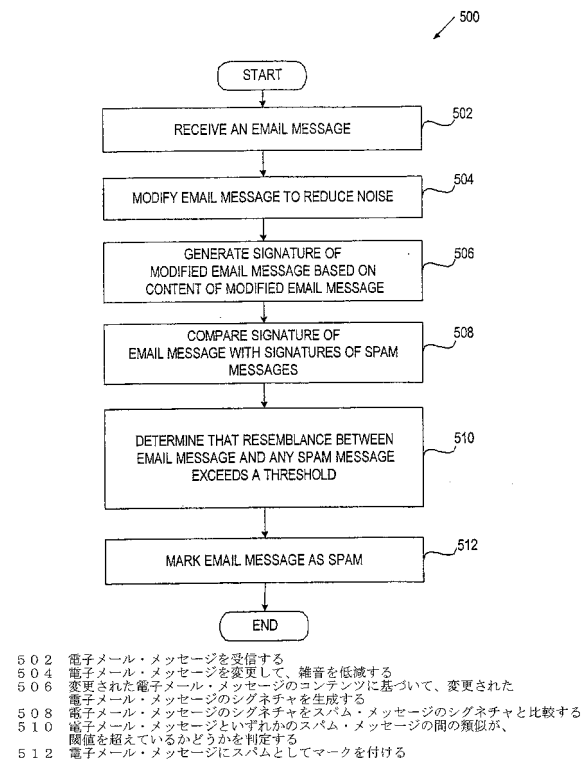
【図 3】



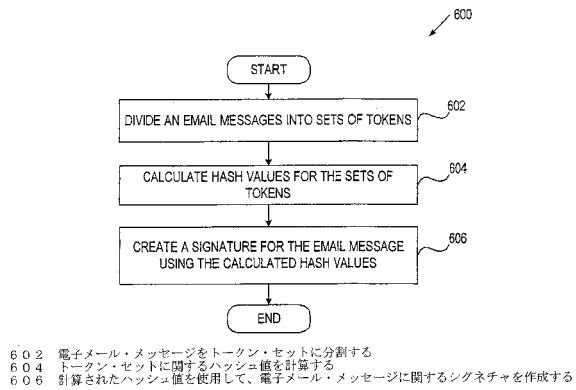
【図 4】



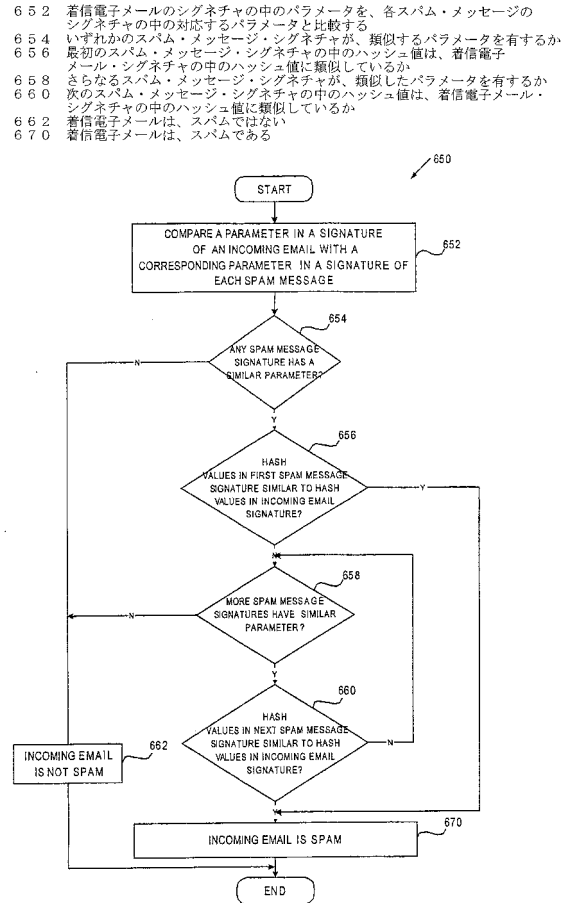
【図 5】



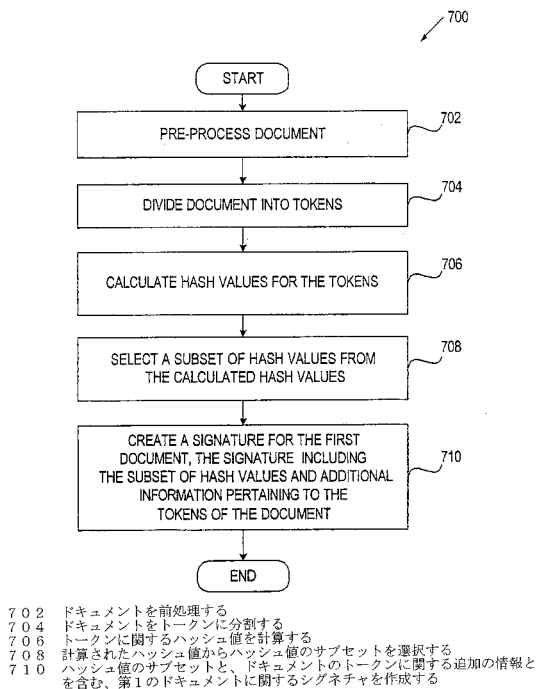
【図 6 A】



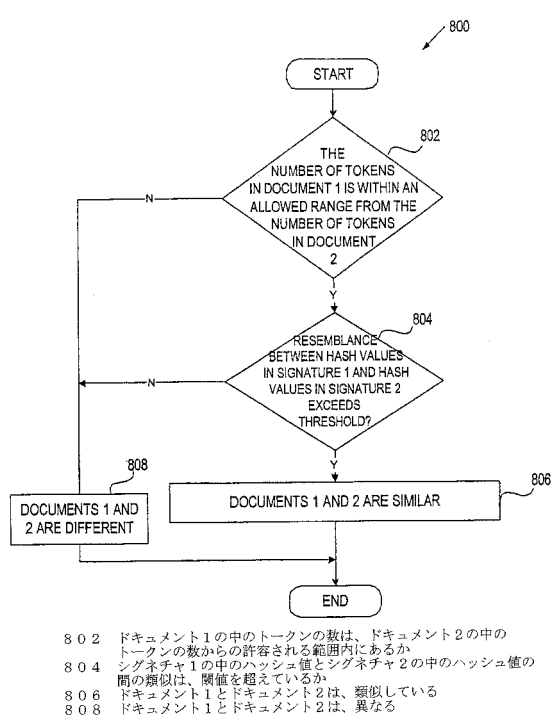
【図 6 B】



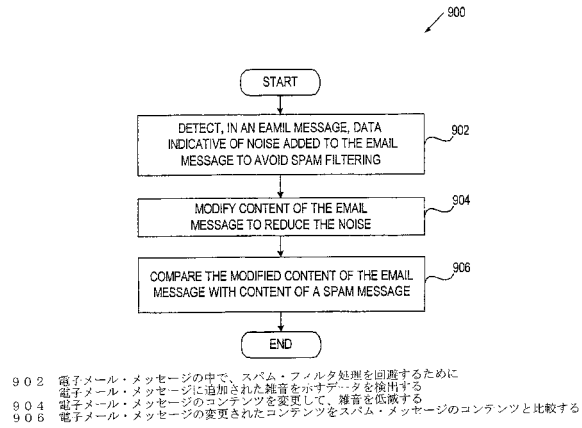
【図 7】



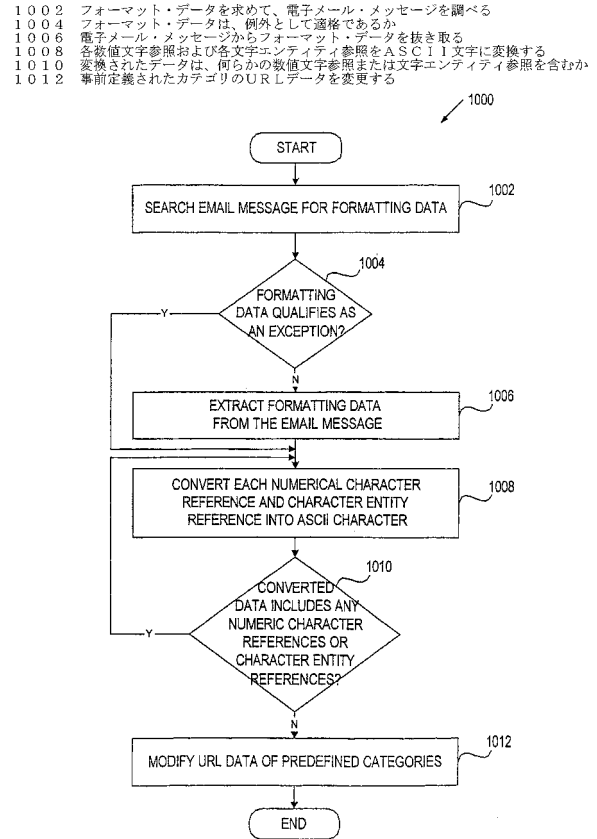
【図 8】



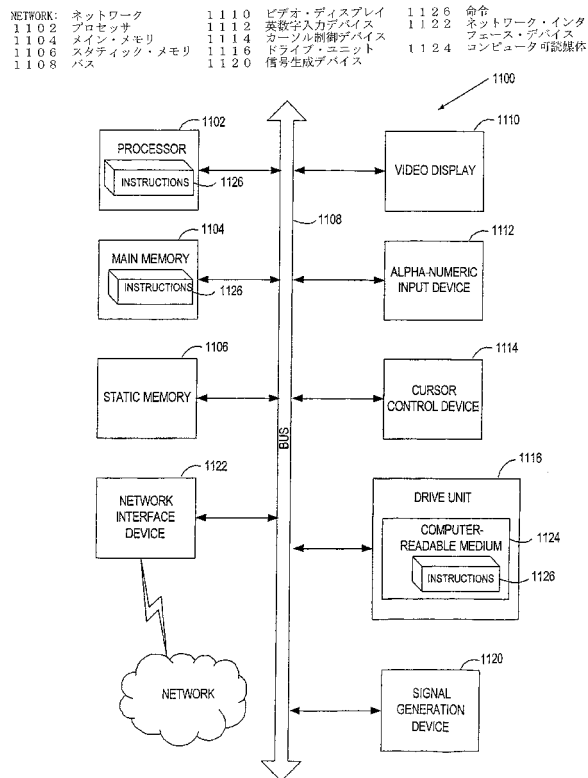
【図 9】



【図 10】



【図 11】



---

フロントページの続き

- (72)発明者 グレッソン, マット  
アメリカ合衆国
- (72)発明者 ホーグストレイト, デイビッド  
アメリカ合衆国
- (72)発明者 ジェンセン, サンディ  
アメリカ合衆国
- (72)発明者 マンテル, エリ  
アメリカ合衆国
- (72)発明者 メドラー, アート  
アメリカ合衆国
- (72)発明者 シュナイダー, ケン  
アメリカ合衆国

審査官 千本 潤介

- (56)参考文献 特開平 1 1 - 0 1 5 7 5 6 ( J P , A )  
特開 2 0 0 0 - 3 5 3 1 3 3 ( J P , A )  
特開 2 0 0 6 - 2 9 3 5 7 3 ( J P , A )  
国際公開第 2 0 0 4 / 1 1 4 6 1 4 ( W O , A 1 )

- (58)調査した分野(Int.Cl. , D B 名)  
G06F 13/00