

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4297345号
(P4297345)

(45) 発行日 平成21年7月15日(2009.7.15)

(24) 登録日 平成21年4月24日(2009.4.24)

(51) Int.Cl. F I
HO4L 12/58 (2006.01) HO4L 12/58 100Z
GO6F 13/00 (2006.01) GO6F 13/00 610Q

請求項の数 4 (全 12 頁)

(21) 出願番号	特願2004-6918 (P2004-6918)	(73) 特許権者	000208891
(22) 出願日	平成16年1月14日 (2004.1.14)		KDDI株式会社
(65) 公開番号	特開2005-202590 (P2005-202590A)		東京都新宿区西新宿二丁目3番2号
(43) 公開日	平成17年7月28日 (2005.7.28)	(73) 特許権者	504016525
審査請求日	平成18年4月27日 (2006.4.27)		吉田 健一
			埼玉県北本市高尾2-232
		(74) 代理人	100084870
			弁理士 田中 香樹
		(74) 代理人	100079289
			弁理士 平木 道人
		(74) 代理人	100119688
			弁理士 田邊 壽二
		(72) 発明者	山崎 克之
			埼玉県上福岡市大原二丁目1番15号 株式会社 KDDI 研究所内
			最終頁に続く

(54) 【発明の名称】 マスメール検出方式およびメールサーバ

(57) 【特許請求の範囲】

【請求項1】

配送対象の電子メールを収集する電子メール収集手段と、

該収集した電子メールを特徴量に変換する特徴量変換手段と、

該変換した特徴量を使ってマスメールを検出するマスメール検出手段とを具備し、

前記特徴量変換手段は電子メールの本文から部分文字列を抽出し、その部分文字列から計算したハッシュ値の集まりを特徴量として用い、

前記マスメール検出手段は、特徴量データベースと特徴量データベースへのポインタを具備し、新規電子メールの特徴量のハッシュ値が該ポインタにエントリされているか否かを判断し、エントリされている場合には、該ポインタを用いて特徴量データベースをアクセスし、該特徴量データベースに既に登録されている電子メールのハッシュ値と比較することにより、新規電子メールと既登録の電子メールとの類似度を判定し、一定数以上のハッシュ値が一致した電子メールを類似電子メールと判定し、該類似電子メールが所定数検出された時に該類似電子メールをマスメールと判定することを特徴とするマスメール検出方式。

【請求項2】

前記特徴量データベースには、類似メール数と、前記ポインタにより参照される被参照数と、電子メールのハッシュ値とが登録され、

類似メールが受信された時には、既登録の類似メールにおける類似メール数が1増加し、被参照数が前記類似メールと既登録の類似メールとのハッシュ値の一致した数に応じて

増加することを特徴とする請求項1に記載のマスメール検出方式。

【請求項3】

前記マスメール検出手段が、directed map cache方式（管理マップキャッシュ方式）またはLRU方式を利用する事の特徴とする請求項2に記載のマスメール検出方式。

【請求項4】

前記請求項1ないし3のいずれかのマスメール検出方式を備えたメールサーバ。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は電子メールのスパム処理に係わり、特に携帯電話やISPなど、大規模な電子メールサーバを運用する事業者が電子メールサーバを経由して配送される電子メールに含まれている未承諾広告などの迷惑メールを検出するのに好適なマスメール検出方式および該マスメール検出方式を備えたメールサーバに関する。

10

【背景技術】

【0002】

電子メールの普及に従い電子メールを搬送手段とした迷惑メールが増加し、社会問題となっている。従来、このような迷惑メールを防止する手段としては、電子メールの受信者が受信に用いる端末に迷惑メールを検知する仕組みを用意し、その仕組みによって迷惑メールを自動的に削除するなどの方法が一般的であった。

【0003】

例えば、SpamAssassinはルールベース方式を用いたソフトであり、bogofilterは機械学習方式を用いたソフトであり、どちらも主としてPCユーザの間で有効な仕組みとして使われている。なお、これらのソフトは、それぞれ、非特許文献1、2に示されている。

20

【非特許文献1】http://www.au.spamassassin.org/presentations/SAGE_IE_2002/

【非特許文献2】http://bogofilter.sourceforge.net/bogofilter_man.html

【発明の開示】

【発明が解決しようとする課題】

【0004】

上記した従来技術は、電子メールの受信者がPCのように一定の水準以上の情報処理能力を持つ受信端末を使用することが前提であり、携帯電話などのような比較的低い能力の受信端末には不向きである。携帯電話などのような比較的低い能力の受信端末を支援するには、事業者側のメールサーバにマスメールを検出する手段を備える事が望ましかった。

30

【0005】

しかし、上記の従来技術は、サーバで用いるには処理速度が遅く大規模な設備を必要とするという問題点があった。また、大人数のユーザに共通したマスメールの検知ルールや機械学習結果を作成するのが困難であり、かつ、新種のスパムに対応するための維持管理のコストが膨大であるといった問題点もあった。

【0006】

本発明の目的は、前記した問題点を解決するために、事前のルール作成や学習が不要で、かつ、高速に動作するマスメール検出方式および該マスメール検出方式を備えたメールサーバを提供することにある。

40

【課題を解決するための手段】

【0007】

上記した目的を達成するために、本発明は、配送対象の電子メールを収集する電子メール収集手段と、該収集した電子メールを特徴量に変換する特徴量変換手段と、該変換した特徴量を使ってマスメールを検出するマスメール検出手段とを具備し、前記特徴量変換手段は電子メールの本文から部分文字列を抽出し、その部分文字列から計算したハッシュ値の集まりを特徴量として用い、前記マスメール検出手段は、特徴量データベースと特徴量データベースへのポインタを具備し、新規電子メールの特徴量のハッシュ値が該ポインタにエントリされているか否かを判断し、エントリされている場合には、該ポインタを用い

50

て特徴量データベースをアクセスし、該特徴量データベースに既に登録されている電子メールのハッシュ値と比較することにより、新規電子メールと既登録の電子メールとの類似度を判定し、一定数以上のハッシュ値が一致した電子メールを類似電子メールと判定し、該類似電子メールが所定数検出された時に該類似電子メールをマスメイルと判定するようにした点に第1の特徴がある。

【0008】

また、本発明は、前記マスメイル検出手段が、記憶領域に記憶しておく電子メールとして頻りに配送される電子メールを優先的に記憶するための手段をもち、該手段としてdirected map cache方式またはLRU方式を用いるようにした点に第2の特徴がある。

【発明の効果】

10

【0009】

本発明によれば、携帯電話やISPなど、大規模な電子メールサーバを運用する事業者が電子メールサーバを経由して配送される電子メールに含まれている未承諾広告などの迷惑メールを検出するのに好適なマスメイル検出方式を提供できる。

【発明を実施するための最良の形態】

【0010】

以下に、図面を参照して本発明を詳細に説明する。図1は、本発明が適用されるシステム構成の一例を示すブロック図である。

【0011】

図において、1は例えば通信事業者（プロバイダ）のメールサーバ群、2はインターネット、3はこれらの間のデータ配送に用いられるスイッチングハブであり、メールサーバ群1はユーザ端末4と接続されている。また、スイッチングハブ3には、本発明のマスメイル検出装置5が接続されている。

20

【0012】

本実施形態は、メールサーバ群1とインターネット2の間でSMTPプロトコルを用いて配送される電子メールの中から、前記マスメイル検出装置5を用いてマスメイルを検出するものである。

【0013】

該マスメイル検出装置5は、電子メール収集手段51、特徴量変換手段52およびマスメイル検出手段53から構成されている。電子メール収集手段51は配送対象の電子メールを収集するものであり適切な計算機上のプログラムで良い。特徴量変換手段52は電子メール収集手段51で収集した電子メールを特徴量に変換するものであり適切な計算機上のプログラムで良い。マスメイル検出手段53は変換した特徴量を使ってマスメイルを検出するものであり適切な計算機上のプログラムで良い。55は検出結果のマスメイルである。

30

【0014】

次に、本実施形態の動作を説明する。電子メール収集手段51は、ネットワーク上に流れる電子メール配送プロトコルを解析し、ネットワーク上に流れる電子メールトラフィックから電子メール本文を抽出する。次に、特徴量変換手段52が電子メール本文から、例えば幾つかのハッシュ値を計算し該メールの特徴量とする。最後に、マスメイル検出手段53が記憶しておいた過去の電子メールと新たに受信した電子メールとを、前記特徴量を用いて比較し、特定の基準に従って類似度を判定し、類似している場合はマスメイルの候補（類似メール）として判定し、一定数以上の類似メールが検出されるとこれをマスメイルと判定する。

40

【0015】

図2は、前記電子メール収集手段51の処理手続きの一例を示すフローチャートである。メールサーバ群1とインターネット2の間では複数のメールが並行して配送されている。そこで、ステップS10では、電子メール収集手段51は、スイッチングハブ3でタッピングすることでTCPパケットを受信する。該受信したTCPパケットは、複数の電子メールの情報が混ったものである。ステップS15では、電子メール収集手段51は、パ

50

ケットの種類を判断する。すなわち、パケットがメールであるか否か、メールであれば、新規メールのパケットであるか、処理中メールのパケットであるか、あるいは処理中メールの終了パケットであるかの判断をする。

【 0 0 1 6 】

そして、受信したパケットが新規メールのものであれば、ステップ S 1 1 に進んで、新規メール用記憶領域を初期設定する。一方、受信したパケットが処理中のメールの終了を意味するパケットであれば、ステップ S 1 3 に進み、処理中のメール本文を特徴量変換手段 5 2 に送信し、次にステップ S 1 4 に進んで、処理中のメール用記憶領域を廃棄 / 解放する。また、受信したパケットが終了以外の処理中のメールパケットであれば、ステップ S 1 2 に進んで、処理中のメール用記憶領域に TCP パケットに含まれるメールの内容を記録する。前記ステップ S 1 5 で、メール以外のパケットであると判断された場合には、何も処理を行わず終了する。図 2 では便宜的にエンド（終了）と記したが、実際の処理は終了することなく、ステップ S 1 0 ~ S 1 4 の処理が継続的に繰り返されることは明らかである。

10

【 0 0 1 7 】

図 3 は前記特徴量変換手段 5 2 の処理の一例を示す説明図であり、図 4 は該特徴量変換手段 5 2 の処理手続きの一例を示すフローチャートである。

【 0 0 1 8 】

本実施形態においては、メール本文の特徴量として、事前に定めた長さ L の文字列（例えば、4 文字）のハッシュ値の集合を用いる。具体的には、図 4 の手順に従い、先にステップ S 2 1 にて、メール本文 1 0 0 の先頭から順番に L 文字ずつ取り出し、そのハッシュ値を計算する。次に、ステップ S 2 2 に進み、計算したハッシュ値をソートし、始めの N 個（例えば 1 0 0 個）を特徴量としてマスメール検出手段 5 3 に送信する。

20

【 0 0 1 9 】

例えば、図 3 に示されているように、長さ L が 4 で、メール本文 1 0 0 が「メール本文の文章」であったとすれば、「メール本」（図 3 の 1 0 1 ）、「イル本文」、「ル本文の」、「本文の文」（図 3 の 1 0 2 ）などのハッシュ値 2 0 1 ~ 2 0 2 を計算する（図 4 のステップ S 2 1 ）。次いで、該ハッシュ値 2 0 1 ~ 2 0 2 をソートした後、始めの N 個を特徴量 2 0 0 としてマスメール検出手段 5 3 に送信する（図 4 のステップ S 2 2 ）。該ハッシュ値としては、例えば 6 4 ビットの整数で表すことができる。

30

【 0 0 2 0 】

図 5 は前記マスメール検出手段 5 3 が利用するデータ構造の例である。3 0 0 は、頻繁に配送される電子メールを特徴量データベース 3 1 0 の中に優先的に記憶するためのデータ構造 Directed Map Cache（以下、DMC）、すなわち管理マップキャッシュ方式を示す。該 DMC 3 0 0 は、特徴量データベース 3 1 0 と、該特徴量データベース 3 1 0 へのポインタ 3 1 1 を有する。特徴量データベース 3 1 0 は、電子メール毎に、特徴量（ハッシュ値 1 ~ N ）、該当電子メールの類似メール数、およびポインタ 3 1 1 のエントリで該当メールを参照しているポインタの数（DMC 被参照数）を記憶した計算機上のデータ構造である。該ポインタ 3 1 1 の各々は、例えば 6 4 ビットで表現することができる。

【 0 0 2 1 】

電子メール収集手段 5 1 が電子メールを抽出すると、特徴量変換手段 5 2 がその電子メール本文から特徴量 2 0 0（図 3 参照）を計算し、最後にマスメール検出手段 5 3 が図 6 に例示した手順に従い、多量に配送されている類似したメールをマスメールとして検知する。具体的には、1 通の電子メールに対して特徴量変換手段 5 2 が計算する特徴量 2 0 0 は N 個（N は、正の整数）のハッシュ値を持つが、マスメール検出手段 5 3 は各メールに対して図 6 の手順に従いステップ S 3 1 からステップ S 4 1 までの処理を最大 N 回繰り返す。

40

【 0 0 2 2 】

ステップ S 3 0 では、前記特徴量 2 0 0 を基に、前記電子メール収集手段 5 1 で収集された電子メールに類似するメールが既にあるか否かの判定が行われる。この処理の一具体

50

例を、図7のフローチャートを参照して説明する。

【0023】

ステップS301では、前記特徴量200の番号を示す数mを1と置き、ステップS302では新規のメールの特徴量200の中のm番目のハッシュ値を抽出する。ステップS303では、該ハッシュ値がポインタ311にエントリされているか否かの判断がなされる。この判断が肯定の場合にはステップS304に進んで、現在のポインタ311から参照されている特徴量データベース310中のエントリと類似度の判定を行う。そして、例えば80%の類似があれば類似メール、80%より小さければ非類似メールと判定する。ステップS305では、 $m = N$ が成立したか否かの判断がなされ、否定の場合には、ステップS306に進んでmに1が加算される。次に、ステップS302に戻って、2番目のハッシュ値が抽出される。以下、同様にして、前記した処理が繰り返し行われ、ステップS305の判断が肯定になると、前記ステップS30の処理は終了する。

10

【0024】

前記ステップS304の類似度判定は、例えば新規の電子メールのハッシュ値200(図3参照)と特徴量データベース310内のハッシュ値が一致した数を利用する。例えば、特徴量の数Nが100個、類似度の閾値が80%の時、80個のハッシュ値が一致すると類似のメールと判断する。一致した数の計測処理を速めるために、予め前記ハッシュ値はソートしておく为好適である。なお、図7ではN個のハッシュ値について類似度を判定したが、必ずしもN個のハッシュ値について類似度を判定する必要はなく、N個より少ないハッシュ値で類似度を判定しても良い。

20

【0025】

図6に戻って説明を続けると、ステップS31において前記mを再度 $m = 1$ とし、ステップS32において、特徴量200の中のm番目のハッシュ値を抽出する。次いで、ステップS33に進み、該m番目のハッシュ値が類似メールのハッシュ値であるか否かの判断がなされる。この判断が否定、すなわち新規の電子メール(非類似メール)の場合には、類似メールは特徴量データベース310に記憶されていないので、ステップS34に進んで、該当メールの特徴量が特徴量データベース310の新規エントリとして登録される。具体的には、新規電子メールの特徴量200(図3参照)を特徴量データベース310のハッシュ値1~ハッシュ値Nとして記憶する。次に、ステップS35に進み、特徴量のベクトル値でDMC300の内容を更新する。

30

【0026】

前記ステップS33の判断が肯定の場合、すなわち類似メールがある場合には、ステップS37に進む。該ステップS37では、特徴量データベース310に既に記憶されているメールの類似メール数(図5参照)を1加算する。次に、ステップS38に進んで、該特徴量データベース310に記憶されたベクトル値でDMC300の内容を更新する。なお、該ステップS38は前記ステップS35と同一の処理であり、その具体例を図8を参照して後述する。

【0027】

ステップS39では、前記類似メール数が予め定められた値S以上になったか否かの判断がなされ、S以上になった場合にはステップS40に進んで該当メールをスパムと判定する。一方、ステップS39の判断が否定の時にはステップS36に進む。ステップS36では、 $m = N$ になったか否かの判断がなされ、この判断が否定の時には、ステップS41に進んで、mに1が加算される。そして、再度ステップS32からの動作が繰り返される。

40

【0028】

図8は図6のステップS35、S38のDMC更新処理の手順を例示したものである。マスメール検出手段53は各電子メールの処理にあたり、まずステップS351の判断をする。すなわち、当該ハッシュ値は、現在のポインタ311から特徴量データベース310の古いエントリを参照しているか否かの判断をする。この判断が否定の時には、ステップS352に進んで、ポインタ311の対応エントリが新しい特徴量データベース310の

50

エントリを指すように設定し、該特徴量データベースの被参照数に1を加算する。

【0029】

一方、前記ステップS351の判断が肯定の時、すなわち当該ハッシュ値が現在のポインタ311から特徴量データベース310の古いエントリを参照している時には、ステップS353に進んで、該ハッシュ値が自分自身のエントリを参照しているか否かの判断をする。すなわち、当該ハッシュ値が前記類似メールの中に含まれているか否かの判断をする。この判断が肯定の時には、何の処理も行わずに図6の処理に抜ける。

【0030】

ステップS353の判断が否定の時、すなわち当該ハッシュ値が前記類似メールの中に含まれていない時には、ステップS354に進み、現在のポインタ311から参照されている特徴量データベース310中の古いエントリのDMC被参照数を1減算する。次いで、ステップS355に進み、DMC被参照数が0であるか否かの判断がなされる。この判断が肯定の時には、ステップS356に進んで、DMC被参照数が0になった過去のメールのエントリを、特徴量データベース310から削除する。前記ステップS355の判断が否定の時には前記ステップS352に進み、ポインタ311の対応エントリが新しい特徴量データベース310のエントリを指すように設定すると共に、該特徴量データベースの被参照数を1加算する。

【0031】

以上の処理によると、類似したメールが多いメールは頻繁に図6のステップS38から起動されて図8の更新処理（より具体的には、ステップS352）が動くのでDMC被参照数は0になりにくい。類似メールがないものはハッシュ値がぶつかったデータを上書きする事で時間の経過とともにDMC被参照数が減少し、前記ステップS356で最終的には削除される。

【0032】

次に、前記した図6～図8の動作を具体例を、図9～図13を参照して参照して説明する。今、インターネットを介して新規のメールが図9に示されているように、メール1、2、3、4の順に収集されたものとし、該メールの特徴量（前記図3の特徴量200）が、メール1に関してはハッシュ値h1, h2, h3, h4、メール2に関してはハッシュ値h2, h3, h6, h7、メール3に関してはハッシュ値h4, h8, h9, h0、メール4に関してはハッシュ値h1, h2, h3, h0であるとする。ここで、類似メールと判定する基準を75%以上の一致とすると、メール4はメール1と類似になる。この判定は図7の処理により行われる。なお、ここでは、説明を簡単にするために、各メールの特徴量が4個であるとした。

【0033】

さて、まずインターネットを介してメール1が抽出されると、図6のステップS33の判断は否定になるので、ステップS34, S35の処理が行われる。ステップS34の処理により特徴量データベース310は図10(b)のハッシュ値1～4にh1～h4が登録され、ステップS35の処理によりDMC300のポインタ311は同図(a)のようになると共に、DMC被参照数が4となる。

【0034】

次に、メール2が抽出されると、前記ステップS33の判断は否定になるので、ステップS34, S35の処理に進む。ステップS34の処理により特徴量データベース310は図11(b)のメール2のハッシュ値1～4にh2, h3, h6, h7が登録され、ステップS35の処理によりDMC300のポインタ311は同図(a)のようになると共に、メール1のDMC被参照数が2となり、メール2のDMC被参照数が4となる。

【0035】

続いて、メール3が抽出されると、前記ステップS33の判断は否定になるので、ステップS34, S35の処理が行われる。ステップS34の処理により特徴量データベース310は図12(b)のように、メール3のハッシュ値1～4にh4, h8, h9, h0が登録され、ステップS35の処理によりDMC300のポインタは同図(a)のようにな

10

20

30

40

50

ると共に、メール 1, 2, 3 のDMC被参照数がそれぞれ、1, 4, 4となる。

【0036】

さらに、メール 4 が抽出されると、このメール 4 は既登録のメール 1 と類似するものであるので、前記ステップ S 3 3 の判断は肯定になり、ステップ S 3 7, S 3 8 の処理が行われる。ステップ S 3 7 の処理により特徴量データベース 3 1 0 のメール 1 の類似メール数に 1 が加算され、図 1 3 (b) のようになる。また、ステップ S 3 8 の処理により、DMC 3 0 0 のポインタは同図 (a) のようになると共に、メール 1, 2, 3 のDMC被参照数はそれぞれ、4, 2, 3 となる。

【0037】

つまり、類似メールが到着すると、前記ステップ S 3 7 でメール 1 の類似メール数に 1 が加算される。次に、ステップ S 3 8 の処理、つまり図 8 の処理において、ハッシュ値 h 1 はポインタ 3 1 が自分自身のメール 1 を指しているのでステップ S 3 5 3 の判断は肯定になり図 8 の処理を抜ける。次のハッシュ値 h 2, h 3 は、ポインタが共にメール 2 を指しているため、ステップ S 3 5 3 の判断は否定となり、ステップ S 3 5 4 以下の処理に移る。そして、ステップ S 3 5 4 でメール 2 のDMC被参照数が 1 減算され、ステップ S 3 5 2 でメール 1 のDMC被参照数が 1 加算される。次のハッシュ値 h 0 についても、同様に処理される。

【0038】

以上のようにして、類似したメールが多いメールは頻繁に図 6 のステップ S 3 8 から起動されて図 8 の更新処理が起動され、DMC被参照数が増加する。一方では、メール 2 を見れば明らかのように、類似メールがないものはハッシュ値がぶつかったデータを上書きする事で時間の経過とともにDMC被参照数が減少する。

【0039】

なお、本発明は前記実施形態に限定されることなく、次のように変形することも可能である。上記の実施形態においては、電子メール収集手段 5 1 はネットワーク上に流れる電子メールをスイッチングハブ 3 でタッピングする事で収集していたが、メールサーバのソフトを変更し、メールサーバが配送対象のメールを直接特徴量変換手段 5 2 に送信するようにしてもよい。またメールの配送プロトコルはSMTPを想定していたが、HTTPを用いたWWWメールのような別の配送形態であってもかまわない。

【0040】

また、メールサーバが配送対象のメールを特徴量変換手段 5 2 に送信する時に既にスパムと判定したメールについてはスパムであるとのマークをつけて送り、その情報を使ってマスメール検出手段 5 3 が、マークのついたメールと類似したメールは即座にスパムと判定してもかまわない。また、メールサーバは前記特徴量変換手段 5 2 までを含むように構成し、該特徴量変換手段 5 2 にて変換された特徴量がネットワーク経由でマスメール検出手段 5 3 に送信される構成にしてもよい。

【0041】

また、上記の実施形態においては、マスメール検出手段 5 3 が、記憶領域に頻繁に配送される電子メールを優先的に記憶するための仕組みとしてDMC 3 0 0 (図 5 参照) を利用していたが、LRU方式のような別の仕組みを利用するのでもかまわない。LRUを使う場合、具体的には特徴量データベース 3 1 0 のエンTRIESを管理するLRUリストを作成し、前記ステップ S 3 7 (図 6 参照) で処理対象とした特徴量データベース 3 1 0 のエンTRIESをLRUの先頭に移動する処理までステップ S 3 7 に含める。また、ステップ S 3 4 で特徴量データベース 3 1 0 に新規のエンTRIESを作る時に必要な記憶領域はLRUの最後のエンTRIESを廃棄して確保し、新規のエンTRIESをLRUの先頭に加える。

【0042】

また、上記の実施形態においては、特徴量変換手段 5 2 の前処理については説明しなかったが、図 1 の電子メール収集手段 5 1 と特徴量変換手段 5 2 との間に前処理手段を設けても良い。この前処理手段は、文字列を抽出する手段であってもよく、この前処理手段により、メールアドレス、電話番号などを抽出するようにしてもよい。また、その他の何ら

10

20

30

40

50

かの前処理を行う手段でもかまわない。この前処理は、電子メールが受信者に表示される際の仕様に従って行うものであってかまわない。この前処理は、例えば受信者の端末の初めの部分（例えば、始めの2頁分）に表示される文字を処理の対象として選択するのでもかまわない。この時に想定される仕様としては、表示に影響を与えるHTMLやMIMEの処理などが上げられるが、これ以外でもかまわない。また、表示の時に大文字と小文字、全角文字と半角文字などを似た文字として扱い、同じ特徴量が計算される仕組みを持つ（例えば、事前に全ての全角大文字を半角小文字に変換する仕組みを持つ）ものでもかまわない。

【0043】

また、上記の実施形態においては、特徴量として電子メール本文に含まれる文字列のハッシュ値を用いたが、バイグラムや単語の出現頻度など、その他の特徴量を用いてもかまわない。

10

【0044】

図14は、本発明の他の実施形態の構成を示すブロック図であり、本発明をメールサーバに組み込んだ構成例である。なお、図14において図1と同一または同等物には同じ符号が付されており、図1と重複する説明は省略する。

【0045】

図14(a)のメールサーバ群1が複数のメールサーバ1a、1b、1cなどから構成されているとすると、本実施形態は、同図(b)に示されているように、電子メール収集手段51、特徴量変換手段52、マスメール検出手段53およびメール処理手段57を、各メールサーバ1a、1b、1cに組み込んだ点に特徴がある。

20

【0046】

本実施形態では、マスメールであるか否かの検出結果は、メール処理手段57へ送られる。メール処理手段57の行う処理は、マスメール検出結果に基づいて、当該マスメールの削除、メール表題部へのマスメールの表示などを行う。また、メールサーバ運用者へのマスメールの通知であってもかまわない。

【0047】

以上の説明から明らかなように、本発明によれば、事前のルール作成や学習が不要である。また、単に電子メールの特徴量を比較することにより類似メールを検出し、該類似メールが一定数に達するとマスメールと判定するので、高速にマスメール検出動作を行うことができる。

30

【図面の簡単な説明】

【0048】

【図1】本発明を含むシステム構成を示すブロック図である。

【図2】電子メール収集手段の処理手続きの一例を示すフローチャートである。

【図3】特徴量変換手段の処理の一例を示す説明図である。

【図4】特徴量変換手段の処理手続きの一例を示すフローチャートである。

【図5】マスメール検出手段が利用するデータ構造の一例であるDirected Mapped Cacheを示す。

【図6】マスメール検出手段の処理の一例を示すフローチャートである。

40

【図7】図6のステップS30の一例を示すフローチャートである。

【図8】図6のステップS35およびS38の一例を示すフローチャートである。

【図9】順次抽出されるメール1、2、3、4の特徴量の説明図である。

【図10】メール1に対するDirected Mapped Cacheのデータ例を示す説明図である。

【図11】メール1、2に対するDirected Mapped Cacheのデータ例を示す説明図である。

【図12】メール1、2、3に対するDirected Mapped Cacheのデータ例を示す説明図である。

【図13】メール1、2、3、4に対するDirected Mapped Cacheのデータ例を示す説明図である。

50

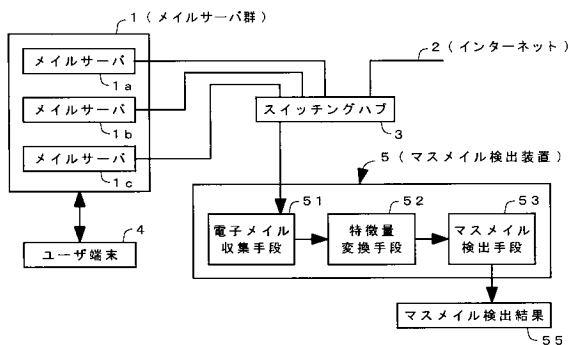
【図14】本発明の他の実施形態の要部を示すブロック図である。

【符号の説明】

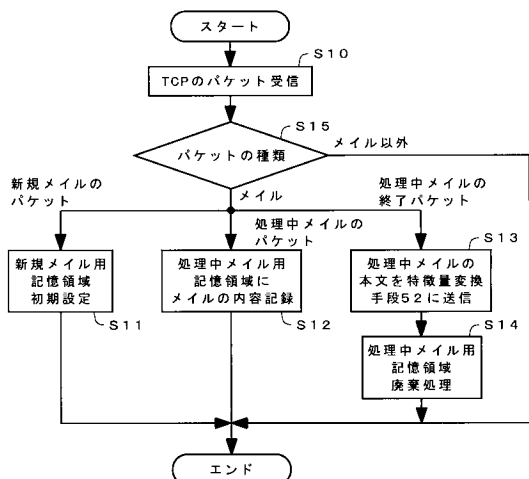
【0049】

1・・・メールサーバ群、2・・・インターネット、3・・・スイッチングハブ、5・・・マスメール検出装置、51・・・電子メール収集手段、52・・・特徴量変換手段、53・・・マスメール検出手段、55・・・マスメール検出結果、300・・・Directed Mapped Cache、310・・・特徴量データベース、311・・・ポインタ。

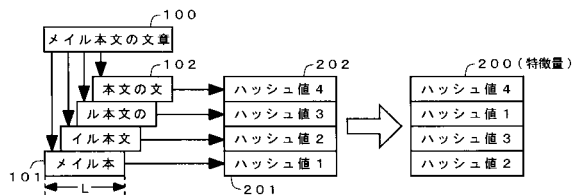
【図1】



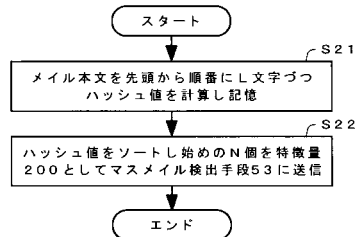
【図2】



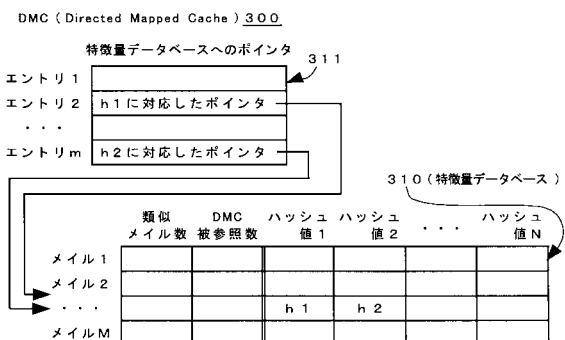
【図3】



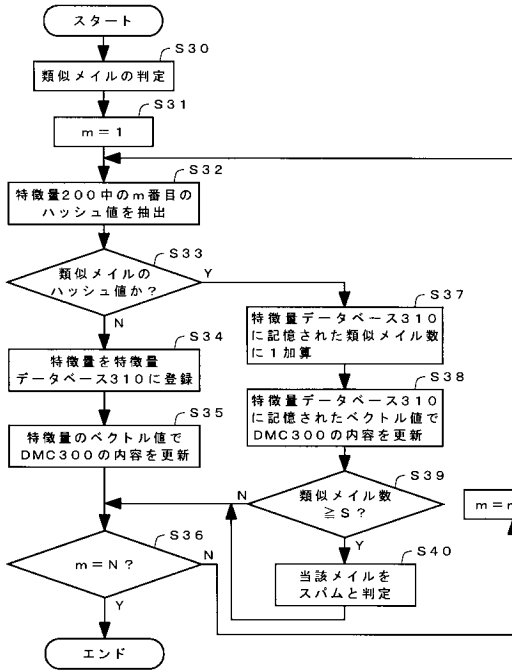
【図4】



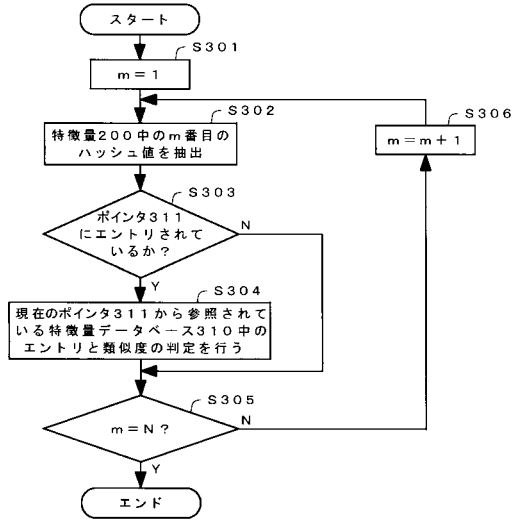
【図5】



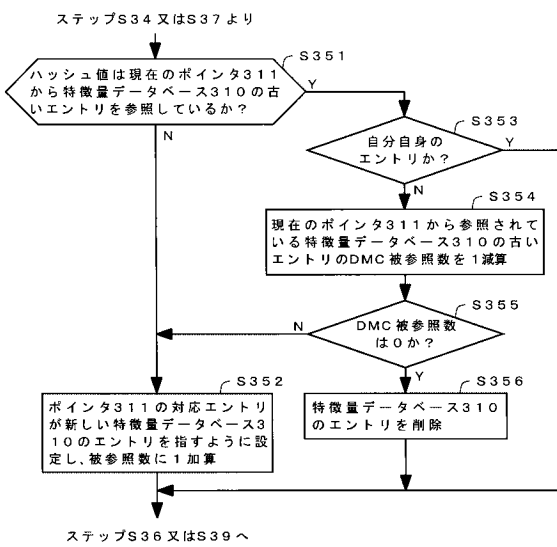
【図6】



【図7】



【図8】



【図10】

(a)

h 0	
h 1	メール 1
h 2	メール 1
h 3	メール 1
h 4	メール 1
h 5	
h 6	
h 7	
h 8	
h 9	

(b)

類似メール数	DMC被参照数	ハッシュ値 1	ハッシュ値 2	ハッシュ値 3	ハッシュ値 4
0	4	h 1	h 2	h 3	h 4

【図9】

メール 1	メール 2	メール 3	メール 4
h 1	h 2	h 4	h 1
h 2	h 3	h 8	h 2
h 3	h 6	h 9	h 3
h 4	h 7	h 0	h 0

(類似は、3/4=75%とする。よって、メール4はメール1と類似)

【図 1 1】

(a)

h 0	
h 1	メール 1
h 2	メール 2
h 3	メール 2
h 4	メール 1
h 5	
h 6	メール 2
h 7	メール 2
h 8	
h 9	

(b)

類似 メール数	DMC 被参照数	ハッシュ 値 1	ハッシュ 値 2	ハッシュ 値 3	ハッシュ 値 4	
メール 1	0	2	h 1	h 2	h 3	h 4
メール 2	0	4	h 2	h 3	h 6	h 7

【図 1 3】

(a)

h 0	メール 3
h 1	メール 1
h 2	メール 1
h 3	メール 1
h 4	メール 1
h 5	
h 6	メール 2
h 7	メール 2
h 8	メール 3
h 9	メール 3

(b)

類似 メール数	DMC 被参照数	ハッシュ 値 1	ハッシュ 値 2	ハッシュ 値 3	ハッシュ 値 4	
メール 1	1	4	h 1	h 2	h 3	h 4
メール 2	0	2	h 2	h 3	h 6	h 7
メール 3	0	3	h 4	h 8	h 9	h 0

【図 1 2】

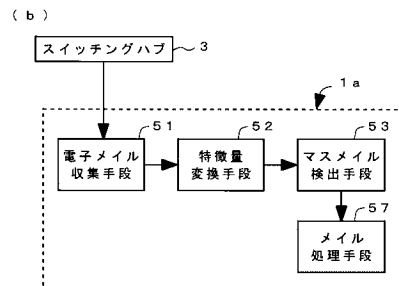
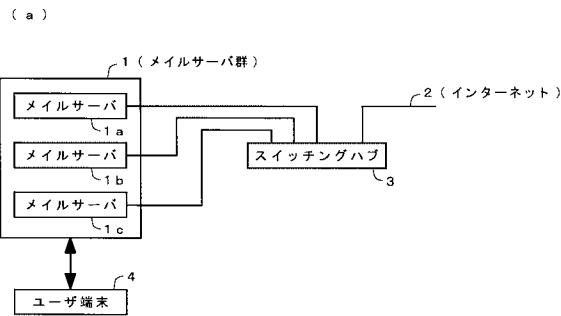
(a)

h 0	メール 3
h 1	メール 1
h 2	メール 2
h 3	メール 2
h 4	メール 3
h 5	
h 6	メール 2
h 7	メール 2
h 8	メール 3
h 9	メール 3

(b)

類似 メール数	DMC 被参照数	ハッシュ 値 1	ハッシュ 値 2	ハッシュ 値 3	ハッシュ 値 4	
メール 1	0	1	h 1	h 2	h 3	h 4
メール 2	0	4	h 2	h 3	h 6	h 7
メール 3	0	4	h 4	h 8	h 9	h 0

【図 1 4】



フロントページの続き

- (72)発明者 藤川 裕充
埼玉県上福岡市大原二丁目1番15号 株式会社 KDDI 研究所内
- (72)発明者 中島 昭浩
東京都新宿区西新宿二丁目3番2号 KDDI 株式会社内
- (72)発明者 本間 輝彰
東京都新宿区西新宿二丁目3番2号 KDDI 株式会社内
- (72)発明者 吉田 健一
埼玉県北本市高尾2-232

審査官 須藤 竜也

(56)参考文献 特表2004-500761(JP,A)

(58)調査した分野(Int.Cl., DB名)
H04L 12/58
G06F 13/00