

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4950439号
(P4950439)

(45) 発行日 平成24年6月13日(2012.6.13)

(24) 登録日 平成24年3月16日(2012.3.16)

(51) Int. Cl. F I
G 0 6 F 13/00 (2006.01) G O 6 F 13/00 5 2 O D
G 0 6 F 15/00 (2006.01) G O 6 F 15/00 4 4 O A

請求項の数 30 (全 24 頁)

(21) 出願番号	特願2005-139956 (P2005-139956)	(73) 特許権者	500046438
(22) 出願日	平成17年5月12日 (2005.5.12)		マイクロソフト コーポレーション
(65) 公開番号	特開2005-327291 (P2005-327291A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成17年11月24日 (2005.11.24)		2-6399 レッドモンド ワン マイ
審査請求日	平成20年5月9日 (2008.5.9)		クロソフト ウェイ
(31) 優先権主張番号	10/845,625	(74) 代理人	100077481
(32) 優先日	平成16年5月14日 (2004.5.14)		弁理士 谷 義一
(33) 優先権主張国	米国 (US)	(74) 代理人	100088915
			弁理士 阿部 和夫
		(72) 発明者	チャ ジャン
			アメリカ合衆国 98052 ワシントン
			州 レッドモンド ワン マイクロソフト
			ウェイ マイクロソフト コーポレーシ
			ョン内

最終頁に続く

(54) 【発明の名称】 部分的複製を使用する、Web コンテンツの分散型ホスティング

(57) 【特許請求の範囲】

【請求項1】

ホスト・コンピューティング・デバイスおよび1または複数のピア・コンピューティング・デバイスを含むコンピュータ・クラスタ内の複数のコンピューティング・デバイスに、Webサイトのコンテンツを分配する方法であって、前記コンテンツは、複数のファイルを含み、前記方法は、

前記複数のファイルのそれぞれについて、前記Webサイトのディレクトリ間の相対的重要度を表すサイト重み(S)に、コンテンツのタイプに応じた重要度を表すタイプ重み(T)を乗算したファイル重み(FW)を計算すること、

前記複数のファイルのそれぞれの複製係数に基づいて決定されている前記コンテンツの複数の部分集合を作成することであって、前記複製係数は、前記部分集合の1つを格納するコンピュータ・デバイスの記憶容量に基づく複製量(D)に基づいて、下式により決定されること、

【数1】

$$D(\lambda_i) = \sum_m |Pm| \max\{1, FW_m \times \lambda_i\}$$

ここで、D(i)はピア・コンピューティング・デバイスに複製される総量であり、iはピア・コンピューティング・デバイスについての複製係数であり、FW_mは前記Webサイトの内のm番目のファイルのファイル重みであり、P_mはWebサイトの内のm番目のファイルのファイルサイズであること、および

前記複数の部分集合の1つを前記複数のコンピューティング・デバイスのそれぞれに分配すること

を備えることを特徴とする方法。

【請求項2】

前記コンピュータ・クラスタは、サーバ・アレイを備えることを特徴とする請求項1に記載の方法。

【請求項3】

前記コンピュータ・クラスタは、ピア・ツー・ピア・ネットワークを備えることを特徴とする請求項1に記載の方法。

【請求項4】

前記タイプ重みは、ファイル拡張子に基づくことを特徴とする請求項1に記載の方法。

【請求項5】

前記サイト重みは、前記Webサイトのディレクトリに割り当てられ、前記ディレクトリ内のそれぞれのファイルは、前記ディレクトリの前記サイト重みを持つことを特徴とする請求項1に記載の方法。

【請求項6】

前記ディレクトリのサブ・ディレクトリは、前記ディレクトリに割り当てられている前記サイト重みを継承することを特徴とする請求項5に記載の方法。

【請求項7】

前記複数のファイルは、複数のメッセージに分割され、それぞれの部分集合は、少なくとも1つのメッセージを備えることを特徴とする請求項1に記載の方法。

【請求項8】

前記複数のメッセージの大部分は、固定サイズであることを特徴とする請求項7に記載の方法。

【請求項9】

さらに、前記複数のファイルを複数のオリジナル・メッセージに分割すること、イレージャ・コーディング方式を適用して複数の符号化メッセージを派生させること、および

前記符号化メッセージのうちの1つまたは複数をそれぞれの部分集合に追加することを備えることを特徴とする請求項1に記載の方法。

【請求項10】

前記イレージャ・コーディング方式は、リード・ソロモン・コーデックを備えることを特徴とする請求項9に記載の方法。

【請求項11】

前記複数の部分集合のうちの1つを分配することは、前記イレージャ・コーディング方式に関連付けられているキーに基づき前記部分集合を分配することを備えることを特徴とする請求項9に記載の方法。

【請求項12】

請求項1ないし11のいずれかに記載の方法を実行するためのコンピュータ実行可能命令を備えることを特徴とするコンピュータ可読記録媒体。

【請求項13】

ホスト・コンピューティング・デバイスおよび1または複数のピア・コンピューティング・デバイスを含むコンピュータ・クラスタ内の複数のコンピューティング・デバイスから、Webサイトのコンテンツを取り出す方法であって、

前記コンテンツは、複数のファイルを含み、

前記複数のファイルのそれぞれに、前記Webサイトのディレクトリ間の相対的重要度を表すサイト重み(S)に、コンテンツのタイプに応じた重要度を表すタイプ重み(T)を乗算したファイル重み(FW)が割り当てられ、

前記複数のファイルは、下式により決定される複製比に基づいて前記複数のコンピューティング・デバイスに分配されており、

10

20

30

40

50

$$w_{m,i} = \max \{ 1, F W_m \times i \}$$

ここで、 $w_{m,i}$ はピア*i*に対する*m*番目のファイルの複製比であり、 i はピア*i*について計算された複製係数であり、 $F W_m$ は前記Webサイトの内の*m*番目のファイルのファイル重みであり、前記方法は、

前記WebサイトのコンテンツのURLを受け取ること、

前記コンテンツの部分集合であるファイルを格納しているコンピューティング・デバイスであって、前記コンテンツを取得することになる複数のコンピューティング・デバイスを決定すること、および

前記複数のコンピューティング・デバイスとのネゴシエーションにより、前記コンテンツに含まれる前記複数のファイルを取得して、前記Webサイトの前記コンテンツを再作成すること

10

を備えたことを特徴とする方法。

【請求項14】

前記URLは、前記URLが前記分散型Webサイトに関連付けられていることを示すタグを含むことを特徴とする請求項13に記載の方法。

【請求項15】

前記複数のコンピューティング・デバイスを決定することは、前記URLを大きな整数値に変換し、前記複数のコンピューティング・デバイスを識別することを備えることを特徴とする請求項13に記載の方法。

【請求項16】

20

前記大きな整数値は、GUID (Global Unifier Identifier) であることを特徴とする請求項15に記載の方法。

【請求項17】

前記複数のコンピューティング・デバイスを決定することは、さらに、前記大きな整数値であるキーを使用する分散型ハッシュ・テーブル (DHT) プロトコルを使用することを備えることを特徴とする請求項15に記載の方法。

【請求項18】

さらに、前記URLに関連付けられているサマリ・ファイルを取得することを含み、前記サマリ・ファイルにより、前記Webファイルを備えるメッセージのインデックスおよび前記Webファイルの修正時刻を識別することを特徴とする請求項13に記載の方法。

30

【請求項19】

前記サマリ・ファイルは、ローカルでキャッシュされることを特徴とする請求項18に記載の方法。

【請求項20】

前記サマリ・ファイルは、前記複数のコンピューティング・デバイスのうちの少なくとも1つから取得された複数のアイテムから再作成されることを特徴とする請求項18に記載の方法。

【請求項21】

さらに、前記複数のファイルを複数のオリジナル・メッセージに分割すること、
イレージャ・コーディング方式を使用して複数の符号化メッセージを派生させること、
および

40

前記符号化メッセージのうちの1つまたは複数をそれぞれの部分集合に追加することを備えることを特徴とする請求項13に記載の方法。

【請求項22】

前記イレージャ・コーディング方式は、リード・ソロモン・コーデックを備えることを特徴とする請求項21に記載の方法。

【請求項23】

前記複数のコンピューティング・デバイスとのネゴシエーションを実行することは、前記イレージャ・コーディング方式に関連付けられたキーを送信することを含むことを特徴とする請求項21に記載の方法。

50

【請求項 2 4】

さらに、前記再作成されたファイルをブラウザで表示すること備えることを特徴とする請求項 1 3 に記載の方法。

【請求項 2 5】

請求項 1 3 ないし 2 4 のいずれかに記載の方法を実行するためのコンピュータ実行可能命令を備えることを特徴とするコンピュータ可読記録媒体。

【請求項 2 6】

ホスト・コンピューティング・デバイスおよび 1 または複数のピア・コンピューティング・デバイスを含むコンピュータ・クラスタ内の複数のコンピューティング・デバイスに、Web サイトのコンテンツを分配しているシステムであって、

前記コンテンツは、複数のファイルを含み、

前記複数のファイルのそれぞれに、前記 Web サイトのディレクトリ間の相対的重要度を表すサイト重み (S) に、コンテンツのタイプに応じた重要度を表すタイプ重み (T) を乗算したファイル重み (F W) が割り当てられ、

前記複数のファイルのそれぞれの複製係数に基づいて決定されている前記コンテンツの複数の部分集合が作成され、前記複数のコンピューティング・デバイスのそれぞれに分配されていることであって、前記複製係数は、前記部分集合の 1 つを格納するコンピュータ・デバイスの記憶容量に基づく複製量 (D) に基づいて、下式により決定され、

【数 2】

$$D(\lambda_i) = \sum_m |P_m| \max\{1, F W_m \times \lambda_i\}$$

ここで、 $D(\lambda_i)$ はピア・コンピューティング・デバイスに複製される総量であり、 λ_i はピア・コンピューティング・デバイスについての複製係数であり、 $F W_m$ は前記 Web サイトの内の m 番目のファイルのファイル重みであり、 P_m は Web サイトの内の m 番目のファイルのファイルサイズであることを特徴とするシステム。

【請求項 2 7】

さらに、前記ファイルを複数のメッセージに分割することを含み、前記ファイルを送信することは、前記ファイルに関連付けられている前記複数のメッセージのうち少なくとも 1 つを前記 1 つのコンピューティング・デバイスに送信することを備えることを特徴とする請求項 2 6 に記載のシステム。

【請求項 2 8】

さらに、イレージャ・コーディング方式を適用して複数の符号化メッセージを派生させ、前記符号化メッセージを前記メッセージとすることを特徴とする請求項 2 7 に記載のシステム。

【請求項 2 9】

ホスト・コンピューティング・デバイスおよび 1 または複数のピア・コンピューティング・デバイスを含むコンピュータ・クラスタ内の複数のコンピューティング・デバイスから、Web サイトのコンテンツを取り出すクライアント・コンピューティング・デバイスであって、

前記コンテンツは、複数のファイルを含み、

前記複数のファイルのそれぞれに、前記 Web サイトのディレクトリ間の相対的重要度を表すサイト重み (S) に、コンテンツのタイプに応じた重要度を表すタイプ重み (T) を乗算したファイル重み (F W) が割り当てられ、

前記複数のファイルは、下式により決定される複製比に基づいて前記複数のコンピューティング・デバイスに分配されており、

$$w_{m,i} = \max\{1, F W_m \times \lambda_i\}$$

ここで、 $w_{m,i}$ はピア i に対する m 番目のファイルの複製比であり、 λ_i はピア i について計算された複製係数であり、 $F W_m$ は前記 Web サイトの内の m 番目のファイルのファイル重みであり、前記クライアント・コンピューティング・デバイスは、

前記 Web サイトのコンテンツの URL を受け取り、

10

20

30

40

50

前記コンテンツの部分集合であるファイルを格納しているコンピューティング・デバイスであって、前記コンテンツが取り出されることになる複数のコンピューティング・デバイスを決定して、

前記複数のコンピューティング・デバイスとのネゴシエーションにより、前記コンテンツに含まれる前記複数のファイルを取得して、前記Webサイトの前記コンテンツを再作成することを特徴とするクライアント・コンピューティング・デバイス。

【請求項30】

前記複数のコンピューティング・デバイスを決定することは、前記URLを大きな整数に変換すること、分散型ハッシュ・テーブル(DHT)プロトコルを適用して、前記大きな整数をキーとして保持している前記コンピューティング・デバイスの位置を突きとめることを備えることを特徴とする請求項29に記載のクライアント・コンピューティング・デバイス。

10

【発明の詳細な説明】

【技術分野】

【0001】

本明細書は、概して、Webコンテンツのホスティングに関し、より詳細には、コンピュータ・クラスタ内のWebコンテンツのホスティングに関する。

【背景技術】

【0002】

インターネットで容易に入手可能な情報量は過去数年の間に著しく増大した。最近まで、ほとんどの情報は企業Webサイトにより提供されていた。しかし、今日では、多くの個人ユーザがパーソナルWebページとして情報を公開(publish)している。これらのパーソナルWebページは、日記、Webログ(通例、ブログ(blog)と呼ばれる)、パーソナル写真/ビデオ・コレクション、個人的な意見(パーソナル・アドバイス)、個人的な体験(パーソナル・エクスペリエンス)などの様々なWebコンテンツを含む場合がある。インターネットは、このようなWebコンテンツを公開する優れた手段を提供するが、Webコンテンツをパブリッシュする個人ユーザの能力に影響を与える様々な障害がある。

20

【0003】

一般に、個々のコンテンツ所有者は、自分のWebコンテンツのホスティングを行ううえで、(1)商用データ・センター上でのホスティングすること、または(2)自分のインターネットリンクを使用して自分のパーソナル・コンピュータ上でのホスティングすること、の2つのオプションが与えられる。商用データ・センターのオプションでは、信頼できるサーバおよび帯域幅のリンクが提供される。しかしながら、いくつかの短所もある。例えば、コンテンツ所有者は、ホスト・サービスについて追加料金を支払わなければならない。彼らは、データ・センターでホスティングできるコンテンツの量において制約を受ける。また、彼らは、データ・センターを通じて利用できる毎日および毎月のトラフィック量も制約を受ける。コンテンツ所有者は、お気に入りのアプリケーションまたはツールを、データ・センター側でそれらのアプリケーションまたはツールをサポートしていないという理由で、利用することができないかも知れない。さらに、データ・センターを利用する多数の人々の間で強力なサーバおよび高帯域リンクが共有されるけれども、個々のユーザの要求に応じて割り当てられる計算能力およびネットワーク帯域幅は制限される可能性がある。

30

40

【0004】

一方、上述のオプション2では、コンテンツ所有者は、自分のコンピュータ上でWebコンテンツをホスティングすることができる。所有するコンピュータを使用してWebコンテンツをホスティングすれば、追加料金の発生が回避され、利用できるホスティング・スペースのサイズは事実上無制限であり(使用ハード・ドライブのサイズのみにより制限される)、所有者はアプリケーション/データベースを自由に選択できる。しかし、そこには障害もある。障害の1つは、Webコンテンツを他のユーザに配信する際に、信頼性

50

が低い、不十分であるという特性である。Webコンテンツに連続したアクセスを提供するためには、所有者のホーム・コンピュータおよびインターネット接続が、常時、稼働し、正常に機能していなければならない。ホーム・コンピュータで障害が発生したり、またはコンピュータの電源が不注意に切られたりすると、コンテンツを提供できなくなる。さらに、所有者のインターネット接続がダウンした場合も、コンテンツを提供できなくなる。コンピュータおよびインターネット接続が絶対にダウンしないようにすることが可能であっても、さらに、十分な帯域幅を確保するなどの、克服すべき他の障害が残っている。インターネット・サービス・プロバイダ（ISP）がユーザのホーム・コンピュータからのインターネット接続のアップロード速度を制限するのはまれなことではない。このアップロード速度がWebコンテンツを他のユーザに高速に配信するのに十分であるということ 10
 はめったにない。例えば、ブロードバンド接続の場合であっても、アップロード速度は、通常、128Kbpsに制限されるのはふつうであり、これはWebコンテンツ・アクセス要求に応えるのに十分な帯域幅とはいえない。

【0005】

Webコンテンツを公開する企業は、高価なサーバ・アレイおよびより高速なインターネット接続に投資することでこのような障害を克服することも可能である。しかし、このような選択肢は、費用が高く、大半の個人ユーザにとっては利用できないものである。幸いなことに、いくつかの種類 20
 のWebコンテンツについては、代替えとなる費用効果の高いソリューションが出現している。この代替えソリューションは、ピア・ツー・ピア（P2P）ネットワークを構築するものである。P2Pコンシューマ・アプリケーションの例として、「NAPSTER」、「KAZAA」、および「gnutella」がある。これらのP2Pコンシューマ・アプリケーションはそれぞれ、複数のコンピュータ間でファイルを共有することに重点を置いている。ファイルの共有はWebサイトの共有に似ているように見えるが、Webサイトの共有には固有の問題点がある。

【0006】

【非特許文献1】"Reed-Solomon Codes and their applications", by S. B. Wicker and V. K. Bhargava, IEEE Press, New York, 1994

【非特許文献2】Bayardo, Jr., R. et al., "YouServ: A Web-Hosting and Content Sharing Tool for the Masses", Proceedings of the International World Wide Web Conference, pp. 345-354, May 2002. 30

【非特許文献3】Ratnasamy, S., et al., "Routing Algorithms for DHTs: Some Open Questions", IPTPS '02, pp. 45-52, 2002.

【非特許文献4】Petar Maymounkov et al., "Kademlia: A Peer-to-Peer Information System Based on the XOR Metric", IPTPS 2002, pp. 52-65, July 2002.

【非特許文献5】Ngan, et al., "Enforcing Fair Sharing of Peer-to-Peer Resources", IPTPS '03, pp. 149-159, 2003.

【非特許文献6】Cardellini, V., et al. "The state of the art in locally distributed web-server systems", ACM Computing Surveys, Vol. 34, No. 2, pp. 263-311, September 2002.

【発明の開示】 40

【発明が解決しようとする課題】

【0007】

問題点の1つは、共有されなければならない情報の量である。Webサイトは、複数のWebページを有し、それぞれのWebページはファイルの集合からなる。複数のWebページ 40
 の1つの関連付けられているハイパーリンクが選択された場合、ファイルの集合全体が直ちに使用できなければならない。したがって、Webサイトの共有は、1個のファイルを共有するのと比べて、大量の記憶領域を必要とし、また広い帯域幅を消費する。他の問題点としては、ファイルの集合を提供する際の取り出し速度および応答時間がある。現在のP2Pコンシューマ・アプリケーションでは、1つのファイルの取り出しの実行は、きわめて低速であり、数時間から数日を要することさえある。さらに、取り出しは、ネ 50

ットワークが混雑していない時間に、またはコンテンツを保持しているコンピュータがオンラインになるとすぐに、スケジュールすることができる。要求側クライアントがWebページを表示するまで何時間も待つ、またはコンピュータがオンラインになるまで待つことはありえないため、Webページを取り出すときに、このように取り出し速度が遅いのは受け入れがたいことである。

【0008】

これらの問題点を克服するために、Webコンテンツを複数のコンピュータに複製する試みがいくつかなされた。すると、クライアントによりコンテンツがアクセスされたときに、コンテンツは、所有者のホーム・コンピュータから、またはコンテンツ全体をホスティングする他のコンピュータのうちの1つからアクセス可能である。複数のコンピュータにWebコンテンツを複製すると、すべてのコンピュータおよびその関連するネットワークリンクが同時にダウンする可能性は低いため、Webコンテンツの信頼性が向上する。しかし、コンテンツを提供する帯域幅は、コンテンツ全体がまだ1台のコンピュータおよびその関連するネットワーク接続から取り出されるので、同じままである。このタイプのシステムでは、Webページへのアクセスの信頼性は向上するが、システムは、依然として、コンテンツ全体を格納するため大量の記憶容量を必要とし、またコンテンツ全体を配信するために広い帯域幅を必要とする。したがって、現在まで、一般大衆が使用するのに好適な、Webコンテンツを公開するための満足のいくソリューションは得られていない。

【課題を解決するための手段】**【0009】**

本明細書で説明されている手法およびメカニズムは、部分的複製を使用する複数のコンピューティング・デバイス上で、Webサイトのコンテンツをホスティングすることを対象とする。

【0010】

Webサイトに関連付けられているそれぞれのファイルに対する相対的重要度が計算される。この相対的重要度を使用して、サーバ・アレイ、ピア・ツー・ピア・ネットワークなどの、コンピュータ・クラスタ内の複数のデバイスに分配されるコンテンツの複数の部分集合を計算する。これらの部分集合は、1つまたは複数のファイルの一部を含むパケットに関してイレージャ・コーディング方式(erasure coding scheme)を使用して作成された符号化メッセージを含むことができる。ファイルの取得時、一定数のはっきりと識別可能なパケット(distinct packets)が、イレージャ・コーディング方式に基づいてデバイスから取り出される。ファイルは、これらのはっきりと識別可能なパケットにより再作成される。複数のデバイスがコンテンツを保持するので、Webサイトに対する取り出しをかなり高速化することができ、どのコンピューティング・デバイスも大量の記憶領域または帯域幅を必要とせずに、信頼性が向上する。

【0011】

以下の図を参照しつつ、非限定的非網羅的实施形態を説明するが、類似の参照番号は、特に断りのない限り、様々な図面全体を通して類似の部分を目指す。

【発明を実施するための最良の形態】**【0012】**

要するに、本発明のWebホスティング・メカニズムは、それぞれがWebコンテンツの部分集合を格納するコンピュータ・クラスタ内の複数のコンピューティング・デバイスからWebコンテンツを取り出すことをサポートするということである。本発明のWebホスティング・メカニズムは、さらに、コンピュータ・クラスタ内の複数のコンピューティング・デバイス上に格納される複数の部分集合にWebコンテンツを分散させることもサポートする。以下で詳述するように、この分配方法では、Webコンテンツを複製する際に必要な記憶領域のサイズを極力抑える。さらに、この分配および取り出し方法は、Webコンテンツに対する信頼性およびアクセス時間を向上させる。これらおよびその他の利点は、以下の詳細な説明を読むと明らかになるであろう。

【0013】

図1は、本発明のWebホスティング・メカニズムを実装するためのシステムの一実施例を示している。システムは、コンピューティング・デバイス100などのコンピューティング・デバイスを備える。最も基本的な構成では、コンピューティング・デバイス100は、少なくとも1つの処理ユニット102およびシステム・メモリ104を備えるのがふつうである。コンピューティング・デバイスの正確な構成と種類に応じて、システム・メモリ104は揮発性(RAMなど)、不揮発性(ROM、フラッシュ・メモリなど)、またはこれら2つの何らかの組合せとすることができる。システム・メモリ104は、通常、オペレーティング・システム105、1つまたは複数のプログラム・モジュール106を格納し、プログラム・データ107を含むこともある。プログラム・モジュール106は、コンテンツを複数のコンピューティング・デバイスに分配し、それらからコンテンツを取り出す本発明のWebホスティング・メカニズムを実装するモジュール130を備える。さらに、システム・メモリ104は、Webページを特定し、表示するためのブラウザを備える。基本構成は、図1において点線108内のコンポーネントにより示されている。

【0014】

コンピューティング・デバイス100は、さらに特徴または機能を追加することもできる。例えば、コンピューティング・デバイス100は、磁気ディスク、光ディスク、またはテープなどの追加データ記憶デバイス(取り外し可能および/または取り外し不可能)を備えることもできる。このような追加記憶装置は、図1では、取り外し可能記憶装置109および取り外し不可能記憶装置110により例示されている。コンピュータ記憶媒体は、コンピュータ可読命令、データ構造体、プログラム・モジュール、またはその他のデータなどの情報を格納する方法または技術で実装される揮発性および不揮発性、取り外し可能および取り外し不可能媒体を含むことができる。システム・メモリ104、取り外し可能記憶装置109、および取り外し不可能記憶装置110は、すべてコンピュータ記憶媒体の実施例である。したがって、コンピュータ記憶媒体としては、限定するものではないが、RAM、ROM、EEPROM、フラッシュ・メモリもしくはその他のメモリ技術、CD-ROM、デジタル多目的ディスク(DVD)もしくはその他の光ディスク記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置もしくはその他の磁気記憶デバイス、または所望の情報を格納するために使用することができコンピューティング・デバイス100によりアクセスできるその他の媒体がある。このような任意のコンピュータ記憶媒体をデバイス100の一部とすることができる。さらにコンピューティング・デバイス100は、キーボード、マウス、ペン、音声入力デバイス、タッチ入力デバイスなどの入力デバイス112を備えることもできる。ディスプレイ、スピーカ、プリンタなどの出力デバイス114を備えることもできる。これらのデバイスは、当業ではよく知られているため、本明細書でさらに詳しい説明をする必要はない。

【0015】

また、コンピューティング・デバイス100は、デバイスがネットワークなどを経由して他のコンピューティング・デバイス118と通信するために使用する通信接続116も含むことができる。(複数の)通信接続116は、通信媒体の一実施例である。通信媒体は、通常、搬送波もしくはその他のトランスポート・メカニズムなどの変調データ信号を介して、コンピュータ可読命令、データ構造体、プログラム・モジュール、またはその他のデータによって具現化することができ、情報配信媒体を含む。「変調データ信号」という用語は、信号内の情報を符号化する方法によりその特性のうち1つまたは複数が設定または変更された信号を意味する。例えば、通信媒体としては、限定するものではないが、有線ネットワークまたは直接配線接続などの有線媒体、ならびに、音響、RF、赤外線、およびその他の無線媒体などの無線媒体がある。コンピュータ可読媒体は、コンピュータによりアクセスできる入手可能な媒体であればどのようなものでもよい。例えば、コンピュータ可読媒体は、限定するものではないが、「コンピュータ記憶媒体」および「通信媒体」を含むことができる。

【0016】

10

20

30

40

50

様々なモジュールおよび手法は、1つまたは複数のコンピュータまたはその他のデバイスにより実行される、プログラム・モジュールなどのコンピュータ実行可能命令の一般的なコンテキストで、本明細書で説明することができる。一般に、プログラム・モジュールは、特定のタスクを実行する、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などを含む。これらのプログラム・モジュールなどは、ネイティブ・コードとして実行するか、または仮想マシンまたはその他のジャスト・イン・タイム・コンパイル実行環境などで、ダウンロードし実行することができる。通常、プログラム・モジュールの機能は、様々な実施形態で望まれているように組み合わせるか、または分散させることができる。これらのモジュールおよび手法の実装は、コンピュータ可読媒体に格納するか、または何らかの形のコンピュータ可読媒体で伝送

10

【0017】

図2は、図1に示されている、コンピューティング・デバイス100などの2つ以上のコンピュータデバイスが本発明の分散型Webホスティング・メカニズムの手法およびメカニズムを実装するように構成されているネットワークの図である。ネットワークは、コンピュータ・クラスタ200と呼ぶことができる。一実施形態では、コンピュータ・クラスタ200は、サーバ・アレイとして構成することができる。他の実施形態では、コンピュータ・クラスタ200は、ピア・ツー・ピア(P2P)ネットワークを使用して構成することができる。一般に、コンピュータ・クラスタ200は、複数のコンピューティング・デバイス202~212を備える。それぞれのコンピューティング・デバイス202~212は、他のコンピューティング・デバイス202~212のうちの1つまたは複数と通信するように構成されている。通信路は、図2では、コンピューティング・デバイス2002~212のうちの2つの間の実線により表されている。通信路は、ローカル・エリア・ネットワーク、インターネット、無線などを經由することができる。コンピューティング・デバイスのうちの1つ(例えば、コンピューティング・デバイス202)は、オリジナルのWebコンテンツ(例えば、Webコンテンツ222)を保持する。以下の説明全体を通して、コンピューティング・デバイス202は、ホスト・コンピュータ202と呼ばれる。他のコンピューティング・デバイス(例えば、コンピューティング・デバイス204~212)は、それぞれ、Webコンテンツ222の部分集合(例えば、部分集合224~232)を保持する。他のコンピューティング・デバイス204~212はホスト・コンピュータのWebコンテンツの部分集合を保持するので、コンピューティング・デバイス204~212は、以下の説明全体を通してホスト・コンピュータ202のピア204~212と呼ばれる。ピアという用語を使用するといっても、コンピュータ・クラスタ200がピア・ツー・ピア・ネットワークとして構成されている必要はないことに留意されたい。むしろ、ピアという用語は、ピアが他のコンピューティング・デバイスのためにコンテンツを保持することを反映する。

20

30

【0018】

概して、ホスト・コンピュータ202では、Webコンテンツ222を、クライアントコンピューティングデバイス260などの1つまたは複数のコンピューティング・デバイス(これ以降クライアント260と呼ぶ)に配信することを望む。クライアント260は、コンピュータ・クラスタ200の一部として示されていないことに留意されたい。これは標準的な構成であるが、本発明の分散型Webホスティング・メカニズムは、クライアント260がコンピュータ・クラスタ200の一部である場合と同様に等しく機能する。いずれの構成でも、クライアント260は、以下で説明される本発明の分散型Webホスティング・メカニズムと連携し、よく知られている手法を使用してインターネット経由でピア上の部分集合にアクセスすることができる。パーソナル・コンテンツを公開するための従来の実装では、クライアント260はホスト・コンピュータ202、またはWebコンテンツ222全体を複製した他のコンピューティング・デバイスからWebコンテンツ203全体を取得する。しかし、上述のように、これらの実装では、かなりの記憶領域を消費し、ホスト202または他のコンピューティング・デバイスの帯域幅のみを利用する。

40

50

したがって、本発明の分散型Webホスティング・メカニズムは、ピア204～212のそれぞれにWebコンテンツ222の部分集合224～232を複製することを重視する。それぞれの部分集合224～232は、コンテンツ222の異なる集合を含むことができる。ピア204～212はWebコンテンツ222全体を格納しないため、クライアント260は、本発明の分散型Webホスティング・メカニズムにより複数のピア（例えば、ピア208～212）からWebコンテンツ222全体を取得する。

【0019】

Webコンテンツ222の部分集合を決定し、これらの部分集合224～232をピア204～212に分配するメカニズムは、図3に示されている流れ図300で例示される。概して、分配プロセス300は、様々なピア上にコンテンツのどの部分集合を複製するかを決定する。コンテンツ全体がピアのそれぞれに分配されるわけではないため、ピアに対する格納コストは極力抑えられ、部分集合をピアに分配することに関連するコスト/時間も極力抑えられる。

10

【0020】

ブロック302で、タイプ重みがWebサイト階層内の各ファイルについて割り当てられる。Webコンテンツは、テキスト、画像、ビデオなどのいくつかのタイプのコンテンツを含むことができる。これらのいくつかのタイプのコンテンツは、それぞれ、コンテンツのタイプに関連付けられた異なる属性を持つことができる。本発明の分散型Webホスティング・メカニズムにより、それぞれのタイプのコンテンツに特定の重み（つまり、タイプ重み）が割り当てられる。タイプ重みは、他のタイプのコンテンツに関連した、そのタイプのコンテンツの重要度を反映する。一実施形態では、重みは大きいほど、このタイプのコンテンツを含むWebページのどれかのページが取り出されるときにこのタイプのコンテンツが欠落していないことを保証したいことを示す。例えば、アイコンをWebページに表示させることは、Webページ上にテキストを表示させることに比べればあまり重要でないといえる。したがって、アイコン・コンテンツに割り当てられるタイプ重みは、テキスト・コンテンツに割り当てられるタイプ重みよりも低くできる。一実施形態では、ユーザはコンテンツのそれぞれの個々のタイプに対し、このタイプ重みを割り当てることができる。例えば、ユーザは、各タイプのファイル拡張子に対しタイプ重みを割り当てることができる。その後、分配プロセスは、ファイル拡張子が認識されると、コンテンツに対し指定されたタイプ重みを割り当てる。他の実施形態では、異なるタイプのコンテンツに対し、既定のタイプ重み設定を設定できる。ユーザは、後から、既定の設定を指定変更することもできる。処理は、ブロック304に続く。

20

30

【0021】

ブロック304で、サイト重みがWebサイト階層内の各ファイルについて割り当てられる。一般に、Webサイトは、興味、時間/イベントなどのトピック別にまとめることができる。特定のWebページに関連付けられているファイルは、ディレクトリに入れることにより、記憶媒体上に編成することができる。したがって、Webサイト階層は、図4に示されているように、階層型ディレクトリ・ツリーとして図形で例示することができる。

【0022】

図4では、Webサイト階層例は、workディレクトリ410およびpersonalディレクトリ430の2つのサブ・ディレクトリを含むrootディレクトリ402を含む。これらのサブ・ディレクトリは両方とも、さらに、サブ・ディレクトリを持つ。workディレクトリ410は、project 1サブ・ディレクトリ412とproject 2サブ・ディレクトリ418を持つ。project 1とproject 2は両方とも、それぞれdata 414および420ならびにsummary 416および422の2つのサブ・ディレクトリを持つ。personalディレクトリ420は、trip 2002 432およびtrip 2004 434の2つのサブ・ディレクトリを持つ。Webサイト階層の実際のレイアウトに特に興味があるわけではないが、それぞれのディレクトリはそのディレクトリに関連付けられたサイト重みを持つことには

40

50

留意されたい。ディレクトリのサイト重みは、Webサイト内の様々なディレクトリ間の相対的重要度を反映するように、Webサイト階層の所有者により割り当てられる。ディレクトリ間の相対的重要度は、Webサイトを閲覧している人が、特定のコンテンツが欠落している場合に経験するであろう困惑の様々なレベルを反映する。相対的重要度は、さらに、特定のコンテンツへのアクセスの可能性が他のコンテンツへのアクセスに比べて高いことも反映する。例えば、最近作成されたWebページは、古いWebページよりもアクセスされる可能性が高い。

【0023】

サイト重みの割り当ては、図4に示されているようなグラフィックを用いたWebサイト階層を形成し、ユーザがそれぞれのディレクトリまたはファイルを選択し、サイト重みを割り当てられるようにすることで実現できる。サイト重みを割り当てるプロセスを減らすため、サイト重み割り当てを少数のディレクトリだけに絞るとよい。そこで、未割り当てディレクトリは、親ディレクトリから自サイト重みを継承することができる。したがって、ディレクトリまたはファイルが自サイト重みを割り当てられていない場合、そのディレクトリまたはファイルは、親からサイト重みを継承することができる。例えば、workディレクトリ410には、サイト重み2.0が割り当てられている。project2ディレクトリおよびそのサブ・ディレクトリには、具体的サイト重みは割り当てられていない。したがって、project2 418およびdata 420、およびsummary 422は、workディレクトリ410からサイト重み2.0を継承する。他の実施形態では、サイト重みをファイル内にリストとして記述することができる。サイト重みを割り当てるためのこれらの変更形態およびその他の変更形態により、本発明の分散型Webホスティング・メカニズムがサイト重みを入手することができる。サイト重みおよびタイプ重みが割り当てられた後、処理は図3のブロック306に続く。

【0024】

ブロック306で、ファイル重みFWがWebサイト階層内の各ファイルについて計算される。一般に、以下に示されるように、ファイル重みFWは、関連付けられているファイルmに対する複製比 $w_{m,i}$ の計算に影響を及ぼす。複製比 $w_{m,i}$ は、コンテンツ（またはファイル）がピア上に何回複製されるかを決定する。したがって、複製比 $w_{m,i}$ は、特定のファイルmに対する取り出し信頼度に影響を及ぼす。ファイルmのファイル重みFWは、以下の式に示されているように、ファイルのサイト重みSにタイプ重みWを掛けることで計算される。

$$FW_m = S_m \times T_m \quad (\text{式1})$$

ただし、mはWebサイト内のm番目のファイルを表す。一般に、以下に示されるように、ファイル重みFWを2倍にすると、コンテンツ複製比も2倍になり、その結果得られるファイルは提供帯域幅および取り出し信頼度の2倍で取り出される。個々のファイル重みFWがWebサイト階層内の各ファイルmについて計算された後、処理はブロック308に続く。

【0025】

ブロック308で、相対的ピア複製係数 λ_i がピアについて決定される（例えば、図2に示されているピア204～212）。相対的ピア複製係数 λ_i は、ホスト・コンピューティング・デバイスについてピアが格納することに同意した同意複製量Dに基づく。相対的ピア複製係数 λ_i は、相対的ピア複製係数 λ_i について以下の方程式を解くことにより決定することができる。

【0026】

【数1】

$$D(\lambda_i) = \sum_m |P_m| \max\{1, FW_m \times \lambda_i\} \quad (\text{式2})$$

【0027】

ここで、 $D(\lambda_i)$ は、ピアiに複製されるコンテンツの総量を表し、 λ_i は、ピアi

10

20

30

40

50

について計算されるピア複製係数を表し、 FW_m は、Webサイト内のm番目のファイルに対するファイル重みを表し、 P_m は、Webサイト内のm番目のファイルのファイル・サイズを表す。したがって、ピアiに複製されるコンテンツの総量 $D(i)$ は同意した量であり、ファイル・サイズおよびファイル重みは、Webサイト内のファイル毎に知られているため、相対的ピア複製係数 α_i を決定できる。

【0028】

二分探索法 (bi-sectional search) を実行することによって、相対的ピア複製係数を決定する1つの方法がある。上記の方程式は、相対的ピア複製係数 α_i が増加する場合に単調増加関数となるので、二分探索法は成功する。同意した複製量について、各ピアとネゴシエートすることができる。例えば、同意した複製量は、ホスト・コンピュータがピアのWebサイトの複製をサポートするためホスト・コンピュータ上の同意した量の複製をやり取りするというホスト・コンピュータによる双方の合意を表すことができる。複製係数 α_i がピアとの格納契約、つまり D (式2を参照) に依存することに留意されたい。すべてのピアの同意した量がWebサイト全体に必要な最低限の記憶装置サイズに等しいという制約はない。実際、本発明のホスティング・メカニズムによれば、各ピアにより格納される複製データの量がどうあれ、信頼性および取り出し速度は改善される。さらに、複製係数 α_i が小さくても、ファイル重み FW_m の大きいいくつかのファイルを大規模に複製することができ、それにより、ファイルの信頼性および取り出し速度が向上する。個々のピアに対するピア複製係数が決定された後、処理はブロック310に続く。

【0029】

ブロック310で、複製のためピアに送られるコンテンツを決定するため、Webサイト内のそれぞれのファイルのピア複製比を計算する。それぞれのファイルのピア複製比は、以下の式を使用して計算される。

$$w_{m,i} = \max \{ 1, FW_m \times \alpha_i \} \quad (\text{式3})$$

ここで、 $w_{m,i}$ は、ピアiに対するm番目のファイルのピア複製比を表し、 α_i は、ピアiについて計算された相対的ピア複製係数を表し、 FW_m は、Webサイト内のm番目のファイルに対するファイル重みを表す。ピア複製比 $w_{m,i}$ は、m番目のファイルのサイズ $|P_m|$ に比例する、ピアiに送られるm番目のファイルに関するコンテンツの量として見る事ができる。ピア複製率が決定されると、処理はブロック312に続く。

【0030】

ブロック312で、これらのファイルはメッセージ内に配置される。一時的に、図5を参照すると、これらのメッセージの作成が説明されている。すでに述べたように、WebサイトはWebページの集合である。それぞれのWebページは、複数のファイルを含む (つまり、Webファイル501)。これらのWebファイル501は、多数のパケットに分けられる (例えば、パケット502 ~ 510)。例えば、大きなファイルを複数のパケットに分割したり、または複数の小さなファイルを1つのパケットにまとめたりすることができる。一実施形態では、これらのパケット502 ~ 510は、データ操作がしやすいように固定サイズである。しかし、この実施形態では、パケットの一部が固定サイズでなくてよいことに留意されたい。例えば、大きなファイルを複数のパケットに分割するときこのようなことが生じることがある。最後のパケットは、他のパケットよりも小さくなる可能性がある。一般に、この実施形態については、パケットの大半は固定サイズである。これらのパケットはそれぞれ、k個のメッセージ (例えば、メッセージ512 ~ 520) にさらに分割できる。再び、これらk個のメッセージは、サイズ固定としてよい。1つのメッセージ (例えば、メッセージ512) は、オリジナルのパケット506のコンテンツに応じて、オリジナルの1つのWebサイト・ファイルの一部を含むか、または複数のWebサイト・ファイルの一部を持つことができる。処理は、図3のブロック316に続く。

【0031】

図3を参照すると、ブロック316で、多数のこれらk個のメッセージがピアに分配される。ホスト・コンピュータとピアiとの間の接続上のトラフィックが最小の場合にメッ

10

20

30

40

50

セージを分配すると都合がよい。一実施形態では、ランダムな数 k のメッセージがピア i に送られる。メッセージのランダムな個数は、 m 番目のファイルについてピア i に対し計算されたピア複製比 $w_{m, i}$ に比例することができる。例えば、パケット A が 16 このメッセージに分割され、ピア Z に対するパケット A に関連付けられているファイルに対するピア複製比は、.5 であり、これら 16 個のメッセージのうちの半分（つまり、8 個のメッセージ）がピア Z に送られる。概して、図 6 について以下で詳述するように、取り出し時に、クライアントは、パケット A を再作成するためにコンピュータ・クラスタ内のピアのうちの 1 つまたは複数から 16 個のメッセージのそれぞれを特定する。この実施形態にはいくつかの短所がある。複数のピアに多数のメッセージが分配する場合でも、クライアントがパケット A に対する要求を送信する時点で、そのコンピュータ・クラスタが、パケット A に関連付けられているそれぞれの特定のメッセージを利用可能な状態で持たない場合がある。このようなことになったら、パケット A だけでなく、パケット A を含むファイルも、取り出せない。各パケットを再作成できる確率を高める方法の 1 つは、分配する前に、メッセージに関してイレージャ・コーディングを実行することである。そこで、ブロック 316 の前にブロック 314 を実行できる。

10

【0032】

ブロック 314 で、任意選択で、メッセージに関してイレージャ・コーディングを実行することができる。以下に示されるように、イレージャ・コーディングをファイル重みと連携して適用すると、さらに、Web コンテンツの信頼性および取り出し速度が向上する。図 5 を再び参照すると、メッセージ 512 ~ 520 は、イレージャ・コーディングと呼ばれるよく知られている数学的ツールを使用して処理される。イレージャ・コーディングはデータの符号化に関してよく知られているが、Web サイトのコンテンツへの応用については、これまで考えられたことがなかった。概して、メッセージ 512 ~ 520 は、 (n, k) イレージャ・コーデックを通じて処理され、 n 個の符号化メッセージ（例えば、符号化メッセージ 522 ~ 550）から成るイレージャ・コーディング空間を形成する。リード・ソロモン・イレージャ・コーデック、トーネイド・コーデック、および LPDC コーデックなど、使用可能な複数のイレージャ・コーディング技術がよく知られている。

20

【0033】

メッセージ誤り訂正符号として、 (n, k) イレージャ・レジリエント符号 (erasure resilient code) のオペレーションは、ガロア体 $GF(p)$ 上の行列乗算を介して記述することができる。

30

【0034】

【数 2】

$$\begin{bmatrix} c_0 \\ c_1 \\ M \\ M \\ c_{n-1} \end{bmatrix} = G \begin{bmatrix} x_0 \\ x_1 \\ M \\ x_{k-1} \end{bmatrix} \tag{式4}$$

40

【0035】

ここで、 p はガロア体の次数であり、 $\{x_0, x_1, \dots, x_{k-1}\}$ はオリジナル・メッセージであり、 $\{c_0, c_1, \dots, c_{n-1}\}$ は符号化メッセージであり、 G は生成行列 (generator matrix) である。一実施形態では、符号化されたメッセージは、すべて同時に生成されるわけではない。むしろ、本発明のホスティング・メカニズムでは、生成行列 G を使用して、符号化メッセージ空間を定義する。クライアントが k 個の符号化メッセージ $\{c'_0, c'_1, \dots, c'_{k-1}\}$ を受信した場合、これら k 個の符号化メッセージは以下のように表すことができる。

【0036】

50

【数3】

$$\begin{bmatrix} c'_0 \\ c'_1 \\ \vdots \\ M \\ \vdots \\ c'_{k-1} \end{bmatrix} = G_k \begin{bmatrix} X_0 \\ X_1 \\ \vdots \\ M \\ \vdots \\ X_{k-1} \end{bmatrix} \quad (\text{式5})$$

【0037】

ここで、 G_k は、符号化メッセージに対応する生成行列 G の k 個の行 (row) により形成される部分生成行列である。この部分生成行列 G_k が最大階数 k を持つ場合、行列 G_k は反転させる (inverse) ことができ、したがって、オリジナル・メッセージは復号化できる。

10

【0038】

それぞれのイレージャ・コーディング技術は、それ独自の利点を有している。例えば、リードソロモン符号は、 k 個のはっきりと識別可能な符号化メッセージが取り出される限り復号化を保證する最大距離分離可能 (MDS ; maximum distance separable) 特性を持つ。Webホスティング・アプリケーションにおける誤りの主要な形態は、ネットワーク伝送における接続の途絶またはパケットの喪失により引き起こされる符号化メッセージの喪失であるため、リードソロモン符号は、本発明の分散型Webホスティング・メカニズムに特に好適である。リードソロモン符号の詳細については、書籍を参照されたい (例えば、非特許文献1を参照)。

20

【0039】

イレージャ符号のパラメータ k は、パケットの粒度とともにイレージャ・コーディング空間のサイズを決定する。オリジナルのパケット 502 ~ 510 は、それぞれ、 k 個のサイズの等しいメッセージに分解されるため、パラメータ k の値が大きいくほど、それぞれのパケットから得られるメッセージ 512 ~ 520 の個数は多くなる。このため、アクセスの粒度とイレージャ・エンコーディングのオーバーヘッドの両方が増大することになる。他方、 k はユーザがコンテンツを同時に取り出すことができるピアの最大数を決める。したがって、ユーザが、多数のピアからコンテンツを取り出して、可能な最大の速度増大が得られるようにするために、適度なサイズ k を選択することが有益である。パラメータ n は、イレージャ符号から生成可能な符号化メッセージの個数を決定する。 n に対し十分に大きな値を取ることで、異なるピアが異なる符号化メッセージを保持するようである。例示されているパラメータ群は、 $k = 16$ および $n = 2^k = 65536$ である。このパラメータ群を使って、4096 (65536 / 16) 個のピアに対応できる。

30

【0040】

図3をもう一度参照して、ブロック316で、これらの符号化メッセージ 522 ~ 550 の部分集合が、オリジナルのメッセージの代わりに、ピアに分配される。イレージャ・コーディングが適用されると、ブロック314で、それぞれのピアは、イレージャ・コーディング空間内の n 個の符号化メッセージから Z 個のはっきりと識別可能な符号化メッセージを受け取る。個数 Z は、 m 番目のファイルおよびピア i について以下のように計算されるピア複製比 $W_{m,i}$ に基づく。

40

$$Z_i = W_{m,i} \times k \quad (\text{式6})$$

ここで、 Z_i は、符号化メッセージのはっきりと識別可能な個数を表し、 $W_{m,i}$ は、 m 番目のファイルとピア i に対するピア複製比を表す。例えば、ピアが複製比 0.5 を持つ場合、そのピアはパケットに対するオリジナル・メッセージの個数の半分に等しい個数の符号化メッセージを受け取らなければならないことを意味する (つまり、50%の重複率)。はっきりと識別可能な符号化メッセージの個数 Z_i は、分数値でもよい。このような場合、 Z_i は、

【0041】

50

【数4】

$$\lfloor x \rfloor$$

【0042】

をフロア関数とすると、確率

【0043】

【数5】

$$(1 + \lfloor Z_i \rfloor - Z_i)$$

【0044】

であれば

【0045】

【数6】

$$\lfloor Z_i \rfloor$$

【0046】

であり、確率

【0047】

【数7】

$$(Z_i - \lfloor Z_i \rfloor)$$

【0048】

であれば

【0049】

【数8】

$$\lfloor Z_i \rfloor + 1$$

【0050】

である、と単に解釈されるだけである。したがって、この確率の調整により、いくつかのピアは1つ多い符号化メッセージを持ち、それ以外のピアは余分な符号化メッセージを持たないことになる。ピアに分配される符号化メッセージが一意であることを保証するために、異なるイレージャ・コーディング・キーをそれぞれのピアに割り当てることができる。イレージャ・コーディング・キーは、イレージャ符号に関連付けられている行列の行 (row) インデックスから導くことができる。ファイル P_m の集計されたコンテンツ複製比は、 C_m 、即ち

【0051】

【数9】

$$C_m = \sum_i W_{m,i} \tag{式7}$$

【0052】

と表され、コンピュータ・クラスタ内で複製されたファイル P_m のコピーの総数である。それらのピアに分配されるWebサイト・コンテンツのこれらの部分集合に関して、クライアントは要求を開始することができ、 k 個のはっきりと識別可能な符号化メッセージが見つかることになる。

【0053】

図6は、ピアの1つまたは複数から欠落していたアイテム (missing item) を取り出すための取り出しプロセス600を例示する流れ図である。欠落しているアイテムは、実際のWebサイト・ファイル、パケット、メッセージ、または符号化メッセージである可能性がある。概して、それぞれのWebページは複数のピアから取り出されるため、Web

10

20

30

40

50

ページを取り出すスループットは増大する。さらに、要求されたWebページを取得する信頼性も、所望のコンテンツを複数のコンピュータ（つまり、ピア）から取得できるため向上する。取り出し時の一般的原理は、オリジナルの個数（つまり、 k 個）のアイテム（例えば、メッセージ、符号化メッセージ）をそれぞれのパケットに対する任意の数のピアから取り出すというものである。その後、取り出されたアイテムを組み合わせ、オリジナル・パケットを再作成し、これらを組み合わせ、オリジナルのWebファイルを作成する。アイテムが符号化メッセージであれば、符号化メッセージをイレージャ復号化して、オリジナル・メッセージを取得し、続いて、これらを上述のように組み合わせる。個々のパケットについては、オンラインのすべてのピアに格納されているはつきりと識別可能なメッセージの個数がそのパケットに対するオリジナル・メッセージの個数（つまり、 k ）よりも大きい限り、パケットを取り出すことができる。それぞれのオリジナルWebファイルは、それを構成する複数の要素パケットが取り出し可能である限り、復元できる。この場合、イレージャ・コーディングを使用せずに、各オリジナル・メッセージを取り出して、パケットを再作成することができる。しかし、イレージャ・コーディングを使用すると、 k をオリジナル・メッセージの個数とした場合、 k 個のはつきりと識別可能な符号化メッセージのみを取り出して、パケットを再作成する。上述のように、リード・ソロモン・コーデックを使用するイレージャ・コーディングが適用され、パラメータ k が16に設定された場合、16個のはつきりと識別可能な符号化メッセージが、オリジナル・パケットの復元を可能にする。

10

【0054】

20

取り出しプロセス600は、ブラウザが要求を送出したときに開始する。その要求の中で、ホストが、そのホストの下部の相対アドレスを加えて、識別される。例えば、典型的な要求は、`<web>www.xyz.com`のように表示される。ここに示すように、`<web>`などのタグは、ホスト・アドレスと一緒にされる。このタグは、その後続くホスト・アドレスが分散型Webサイトであることを示す。このようなタグがないと、ブラウザ要求を受け取ったプログラムは、分散型Webサイトと単一サーバによりホスティングされている通常のWebサイトとを区別できない場合がある。一実装形態では、このプログラムはプロキシとして実装される。この実装形態では、このプログラムはプロキシ・ポートを介して複数のピアと通信する。しかし、プログラムは、本発明の分散型Webホスティング・メカニズムの精神から逸脱することなく他の構成で実装することができる。例えば、このプログラムは、他のブラウザ・コンポーネントとともにインストールされるブラウザ内のコンポーネント（例えば、ツールバー）とすることができる。また、ブラウザのオプションの取り出しブロックとして実装することも可能である。任意の構成について、このプログラムにより実行される取り出しプロセスは、図6に例示されている流れ図に関連して詳細に説明される。

30

【0055】

そこで、ブロック602で、`www.xyz.com/personal/trip2004/picture1.jpg`などのURL（Uniform Resource Locator）を受け入れる。このURLは、分散型Webコンデンサを含むものとして認識される。処理は、ブロック604に続く。

【0056】

40

ブロック604で、このURLに関連付けられているWebコンテンツの一部を保持するピアのリストが取得される。URLは、GUID（Global Unifier Identifier）などの大きな整数値に変換される。このURLを保持しているピアのリストは、GUIDに関連付けられているピアのリストである。一実施形態では、ピアのリストを決定することは、ルート・パスを指定する1つの列（カラム）とGUIDを指定する第2の列（カラム）を持つローカルのGUIDテーブルをチェックすることにより行われる。このローカルGUIDテーブルは、クライアントが分散型Webサイトにアクセスすると必ず作成され、更新されるようにできる。その後、Webコンテンツの一部を格納しているコンピュータのリストを取得するため、GUIDテーブル・リストがコンピュータ・クラスタ全体を通して送信される。他の実施形態では、GUIDを保持するピアのリストは、分散型ハ

50

ッシュ・テーブル (DHT) 方式を使用してルックアップを実行することにより決定される。GUIDに関連付けられているピアを識別するDHT手法は、当業ではよく知られている。処理は、ブロック606に続く。

【0057】

ブロック606で、このURLに関連付けられているサマリ・ファイルが取得される。一般に、階層型Webサイトのディレクトリ毎にサマリ・ファイルが用意される。つまり、サマリ・ファイルは、Webサイトの構造についての提示 (glimpse) を与える。サマリ・ファイルを使用することにより、クライアントは、Webページ/ファイルを構成するパケットの個数およびインデックスを決定し、そうして、目的のWebページ/ファイルに対する正しいパケットを取り出すことができる。要求されたファイルを取り出すために、その要求されたファイルが配置されているディレクトリに関連付けられているサマリ・ファイルを取得する。そこで、正しいサマリ・ファイルを特定するために、Webサイト階層がトラバースされる。サマリ・ファイルは、ファイル/子ディレクトリ毎に1エンタリを含む。サマリ・ファイルの各エンタリは、ファイル/子ディレクトリの名前を識別し、変更プロパティ、ファイル/フォルダ識別子、およびファイルの長さを含む。さらに、各エンタリは、Webページ/ファイルを構成するパケットを識別するイレージャコーディング・パケット識別子を含むことができる。変更プロパティは、関連するファイル/子ディレクトリが最後に更新された時点のタイム・スタンプである。変更プロパティがピア毎に異なる場合、古いタイム・スタンプを持つピアに格納されているアイテムは取り出されない。最初に、サマリ・ファイルがチェックされ、ローカルに置かれているかどうかを調べられる。サマリ・ファイルは、このディレクトリの下の子ファイルのどれかがすでにクライアントによりアクセスされていた場合にローカルに置かれる。すでにアクセスされていた場合、サマリ・ファイルはピアから取り出す必要はない。アクセスされていなかった場合、サマリ・ファイルは、ピアの1つまたは複数から取り出され、後続の要求のためローカルのクライアントに格納される。したがって、漸進的なキャッシュ・プロセスがある。サマリ・ファイル自体は、複数のパケットに分割され、イレージャ符号化されたメッセージに符号化されるようにできる。取り出し時に、クライアントは、サマリ・ファイルに関連付けられているイレージャ符号化されたメッセージを取り出し、この符号化メッセージを複数のパケットに復号化し、サマリ・ファイルを組み立て直す。サマリ・ファイルが利用可能になると、処理はブロック608に続く。

【0058】

ブロック608で、複数のピアから取得されたサマリ・ファイルが検証される。Webページ/ファイルの古い (つまり、以前の) バージョンを保持していることをサマリ・ファイルが示しているピアは、取り出しプロセスには関わらない。サマリ・ファイルが検証された後、処理はブロック610に続く。

【0059】

ブロック610で、ネゴシエーション・プロセスが利用可能なピアにより実行される。ネゴシエーション・プロセスを使用して、Webページ/ファイルおよび上述のサマリ・ファイルの両方を取り出す。一実施形態では、ネゴシエーション・プロセスにより、利用可能なピアから実際のWebサイト・ファイルを取得する。他の実施形態では、ネゴシエーション・プロセスにより、利用可能なピアからWebサイト・ファイルを構成するパケットを取得する。さらに他の実施形態では、ネゴシエーション・プロセスにより、利用可能なピアからメッセージまたは符号化メッセージを取得する。この符号化メッセージは、その後、パケットに復号化され、パケットはさらに組み立てられ、オリジナルのWebページ/ファイルが再作成される。

【0060】

一時的に、図7を参照すると、利用可能なピアから欠落していたファイルまたはメッセージを取得する例示的なネゴシエーション・プロセス700を示す、シーケンシャルな流れ図が説明されている。図7の「アイテム」という用語は、要求された情報がWebサイト・ファイル、パケット、メッセージ、または符号化メッセージであることを示すために

10

20

30

40

50

使用されている。左のラインは、クライアント（つまり、Webコンテンツを見たい、分散型WebサイトのURLを送信した、コンピューティング・デバイス）が受信または送信した通信を表す。右のラインは、ピアの1つが受信または送信した通信を表す。ネゴシエーション・プロセス700は、必要なアイテムが取り出されるまで識別されたピア毎に実行することができる。

【0061】

クライアントは、プロキシを介して、どのアイテムが必要かを通知する（702）。これらのアイテムが符号化メッセージの場合、クライアントはパケット識別子およびキーを供給する。それに応じて、ピアは、要求された、欠落していたアイテムに関連付けられている利用可能なアイテムに関する情報を提供する。再び、これらのアイテムが符号化メッセージであれば、ピアは、ピア上にローカルで格納されている符号化メッセージのパケット識別子およびキーを供給する。706で、クライアントは、特定の欠落していたアイテム、即ち通知されたピアがそのアイテムの利用が可能であると通知した特定の欠落アイテム、を求める要求を送信する。708で、そのピアは、特定の欠落していたアイテムをクライアントに送信する。

【0062】

図6に戻ると、ネゴシエーション・プロセスで欠落していたアイテムが取得された後、処理はブロック612に続く。ブロック612で、情報は組み立てられ、関連するWebページが表示される。情報を組み立てることには、符号化メッセージをパケットに復号化し、その後、それらのパケットをURLに関連する特定のWebファイルに組み立てることを必要とする。

【0063】

このようにして、図6および7に関して上で説明したように、1つのWebファイルの異なる部分が、異なる場所で実行している複数のコンピューティング・デバイスから、取得されることが可能になる。通常、どのコンピューティング・デバイスも、Webファイルの完全なコピーを持たない。このことは、特にイレージャ・コーディングが実行される場合、特に当てはまる。したがって、クライアントは、複数の部分を取り出して、それらを1つにまとめ上げて要求されたWebファイルを作成し、Webファイルをブラウザに表示する役目を負う。

【0064】

本発明の分散型ホスティング・メカニズムに関する実験が実行された。これらの実験では、3つのシナリオ、1)複数のピア上にWebサイト全体を複製すること、2)イレージャ・コーディングを使用せずに複数のピア上にWebサイトの一部分を複製すること、および3)イレージャ・コーディングを使用して複数のピア上にWebサイトの一部分を複製すること、をテストした。図8および9は、上の3つの特定のシナリオについて、それぞれ、テスト結果801、802、803、ならびに901、902、および903を例示している。シナリオ2については、それぞれのパケットはk個の断片に分割される。しかし、複製段階では、オリジナルのメッセージ断片が、イレージャ・コーディング無しで、他のピアに送信される。すべてのシナリオについて、コンテンツ（つまり、集計されたコンテンツ複製比C）を分配しホスティングするために同じ量のネットワークおよび記憶資源が使用されると仮定した。もう1つの仮定は、コンピュータ・クラスタ内のピアのそれぞれが、同一の提供帯域幅を持ち、オンラインになっている確率として独立した値を持ってクライアントにサービスを提供するというものであった。コンピュータ・クラスタ（例えば、P2Pネットワーク）内のWebサイトを取り出すのに成功する確率は、図8に示されている。様々なシナリオを使用した際のWebサイトを取り出すための平均速度増大は、図9に示されている。図8および9の両方について、横軸は、ピアがオンラインである確率である。コンテンツ分配のパラメータはC = 8（コンテンツの8コピーがホスティングされている）、k = 16に設定された。次に、これらの図のそれぞれについて詳述する。

【0065】

図8は、様々なシナリオにおいて、オンラインになる確率に対する取り出し失敗率、に関連するテスト結果を例示するグラフである。上述のように、曲線801は、Webサイト全体が複製されるシナリオ1を表し、曲線802は、Webサイト全体のうちの一部がイレージャ・コーディング無しで複製されるシナリオ2を表し、曲線803は、Webサイト全体のうちの一部がイレージャ・コーディングを使用して複製されるシナリオ3を表す。曲線803では、同じ量のネットワークおよび記憶資源を使用する曲線801および802との比較で、取り出しの信頼性が著しく向上していることに留意されたい。実際、ピアがオンラインである確率が0.13よりも大きいと、イレージャ符号化されたコンテンツに関するWebサイトの取り出し失敗率は、Web複製曲線801全体、およびイレージャ・コーディング無しの部分的Web複製曲線802、と比べて数千倍小さかった。

10

【0066】

図9は、様々なシナリオにおいて、オンラインになる確率に対する取り出し速度向上、に関連するテスト結果を例示するグラフである。上述のように、曲線901は、Webサイト全体が複製されるシナリオ1を表し、曲線902は、Webサイト全体のうちの一部がイレージャ・コーディング無しで複製されるシナリオ2を表し、曲線903は、Webサイト全体のうちの一部がイレージャ・コーディングを使用して複製されるシナリオ3を表す。イレージャ符号化曲線903は、それぞれ全Web複製(曲線901)およびイレージャ・コーディング無しの部分的Web複製(曲線902)と比べて、16倍および1~10倍、取り出し速度が高速化されることに注目されたい。

20

【0067】

イレージャ符号化コンテンツ分配、および不均等な重み割り当てを使用する階層型コンテンツ編成によるP2P Webホスティングシステムが設計された。Webサイトは、7つのピア上に複製された。オリジナルのWebサイトでは、228メガバイトを消費した。複製時に、それぞれのピアは、そのWebサイトのうちの60メガバイト分をホスティングすることに同意しており、その結果、平均複製比は0.26となった。Webファイルの重み付けは不等なので、実際のWebファイルに対するピア複製比は0.25から1.0と変動する。Webページ時に、クライアントは、7つのピアから同時にWebコンテンツを取り出し、イレージャ符号化メッセージを復号化し、Webページをレンダリングした。

30

【0068】

したがって、説明したように、本発明の分散型ホスティング・メカニズムは、Webコンテンツにアクセスする取り出し速度を向上させ、またコンテンツを取り出す信頼性を向上させる。さらに、分散型ホスティング・メカニズムは、それぞれの個別ピア上に複製されるコンテンツの量を低減する。そのため、ピアでは、記憶領域の追加または帯域幅の増大に関わる膨大な費用は発生しない。本発明のWebホスティング・メカニズムは、ホスト・コンポーネント、ピア・コンポーネント、およびクライアント・コンポーネントの3つのコンポーネントを含む。ホスト・コンポーネントは、Webページ/ファイルの一部(つまり、部分集合)をピアに分配する。ピア・コンポーネントは、Webページ/ファイルの分配された部分を受け取り、その後、クライアントから要求があったら再分配する。クライアント・コンポーネントでは、複数のピアからオリジナル・メッセージまたはイレージャ符号化メッセージの形でWebページ/ファイルの一部を取り出す作業を調整する。その後、オリジナル・メッセージは組み立てられ、要求されたWebページ/ファイルを形成する。符号化メッセージはパケットにイレージャ復号化され、パケットはさらに組み立てられ、要求されたWebページ/ファイルが形成される。

40

【0069】

本明細書全体を通して、「一実施形態」、または「一実施例」と記述されている場合、これは、特定の説明されている機能、構造、または特性が本発明の少なくとも一実施形態に含まれることを意味する。したがって、このようなフレーズを使用した場合、1つの実施形態というよりも複数の実施形態を指している。さらに、説明されている機能、構造、

50

または特性は、1つまたは複数の実施形態において、適当な方法により組み合わせることができる。

【0070】

ただし、当業者であれば、本発明は、特定の詳細の1つまたは複数を使用せずに、または他の方法、資源、材料などを使用して、実践できることは理解できるであろう。一方、よく知られている構造、資源、またはオペレーションについては、単に本発明の態様をわかりにくくすることを避けるために、詳細に示したり、あるいは説明したりしてはいない。

【0071】

複数の実施例およびアプリケーション例を例示し、説明したが、本発明は、上述の正確な構成および資源に限定されないことが理解されるであろう。請求の範囲から逸脱することなく、本明細書に開示されている本発明の方法およびシステムの配置、オペレーション、および詳細に対し、様々な修正、変更、および変形を加えられることは、当業者には明白なことであろう。

【図面の簡単な説明】

【0072】

【図1】本明細書で説明されている手法およびメカニズムを実装するために使用することができるコンピューティング・デバイスの図である。

【図2】図2に示されている2つ以上のコンピュータデバイスが本明細書で説明されている分散型Webホスティング・メカニズムを実装するように構成されているネットワークの図である。

【図3】複製されたWebコンテンツの部分集合を図2に示されている複数のコンピューティング・デバイスのうちの1つに分配する分配プロセスを例示する流れ図である。

【図4】図3に示されている分配プロセスで使用するのに好適なWebコンテンツにサイト重みを割り当てることを例示するツリー図である。

【図5】図3に示されている分配プロセスで使用するのに好適なメッセージにWebファイルを変換することを図形で例示するブロック図である。

【図6】図2に示されている複数のコンピューティング・デバイスのうちの1つからWebコンテンツの部分集合を取り出す取り出しプロセスを例示する流れ図である。

【図7】図6に示されている取り出しプロセスで使用するのに好適な欠落している項目を取得するネゴシエーション・プロセスを例示する逐次的流れ図である。

【図8】様々なシナリオにおいて取り出し失敗率対オンラインになる確率に関連するテスト結果を例示するグラフである。

【図9】本発明の分散型Webホスティング・メカニズムにより、様々なシナリオにおいて取り出し速度向上対オンラインになる確率に関連するテスト結果を例示するグラフである。

【符号の説明】

【0073】

- 100 コンピューティング・デバイス
- 102 処理ユニット
- 104 システム・メモリ
- 105 オペレーティング・システム
- 106 プログラム・モジュール
- 107 プログラム・データ
- 108 基本構成
- 109 取り外し可能記憶装置
- 110 取り外し不可能記憶装置
- 112 入力デバイス
- 114 出力デバイス
- 116 通信接続

10

20

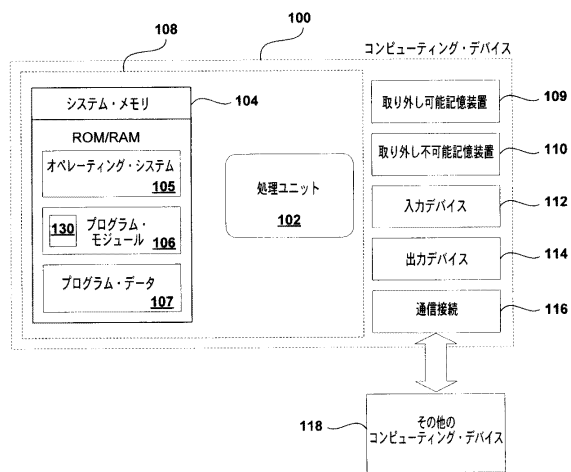
30

40

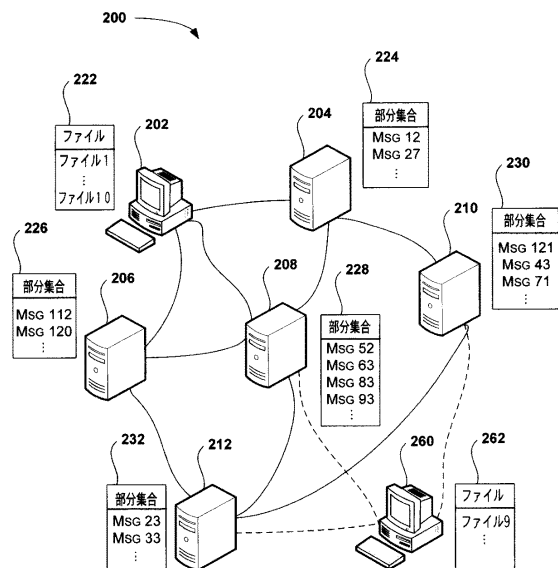
50

- 118 その他のコンピューティング・デバイス
- 200 ネットワーク(コンピュータ・クラスタ)
- 202 ~ 212 コンピューティング・デバイス
- 260 コンピューティング・デバイス
- 801、802、803 テスト結果
- 901、902、903 テスト結果

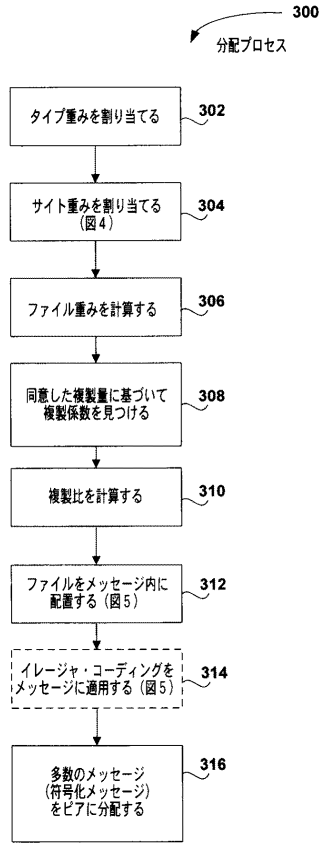
【図1】



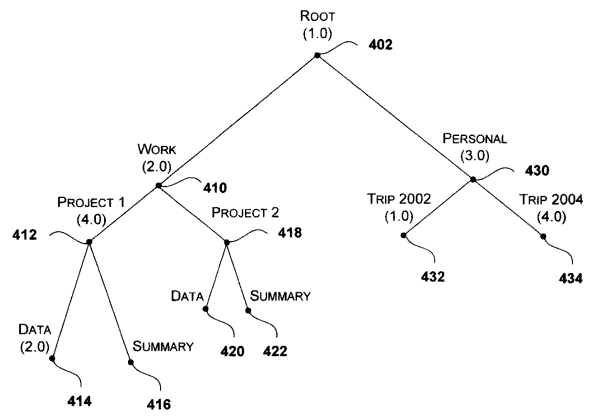
【図2】



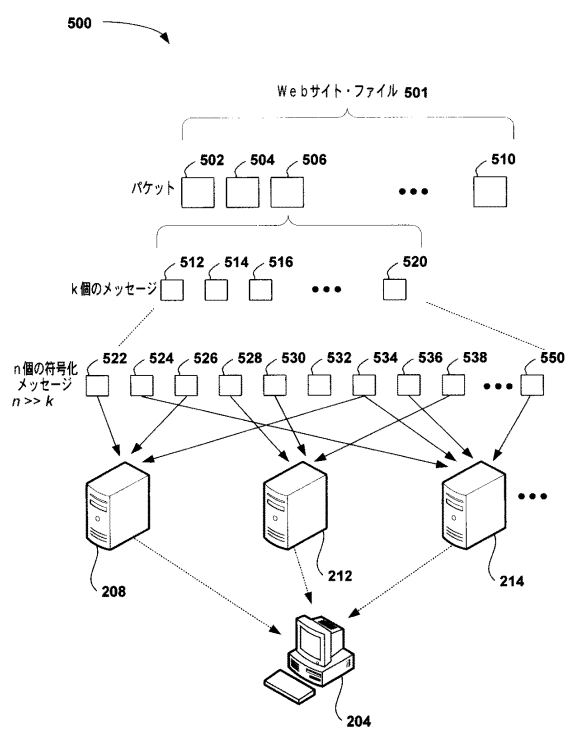
【 図 3 】



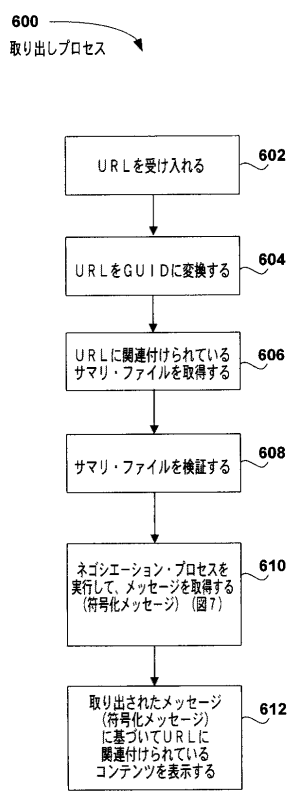
【 図 4 】



【 図 5 】



【 図 6 】



フロントページの続き

(72)発明者 ジン リ

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ
マイクロソフト コーポレーション内

審査官 鈴木 理絵子

(56)参考文献 特開2002-032280(JP,A)
特開2003-216521(JP,A)
特開平10-093446(JP,A)
特開2002-278859(JP,A)
国際公開第01/090943(WO,A1)
特表2002-520735(JP,A)
特開2003-067279(JP,A)
国際公開第03/071800(WO,A1)
特開2002-063064(JP,A)
特開2001-333032(JP,A)
特開平10-198590(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/00
G06F 15/00