

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2014-532220

(P2014-532220A)

(43) 公表日 平成26年12月4日(2014.12.4)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30 180Z	5B084
G06F 13/00 (2006.01)	G06F 13/00 540F	

審査請求 有 予備審査請求 未請求 (全 21 頁)

(21) 出願番号 特願2014-532240 (P2014-532240)
 (86) (22) 出願日 平成24年12月13日 (2012.12.13)
 (85) 翻訳文提出日 平成26年3月28日 (2014.3.28)
 (86) 国際出願番号 PCT/CN2012/086575
 (87) 国際公開番号 W02013/087005
 (87) 国際公開日 平成25年6月20日 (2013.6.20)
 (31) 優先権主張番号 201110415749.9
 (32) 優先日 平成23年12月13日 (2011.12.13)
 (33) 優先権主張国 中国 (CN)

(71) 出願人 507231932
 北大方正集▲团▼有限公司
 PEKING UNIVERSITY F
 OUNDER GROUP CO., L
 TD
 中華人民共和国北京市▲海▼淀区成府路2
 98号中▲関▼村方正大厦5▲層▼
 5 Floor, Zhongguanc
 un Founder Building
 , No. 298, Chengfu R
 oad, Haidian Distri
 ct, Beijing 100871,
 China

最終頁に続く

(54) 【発明の名称】 ネットコメントの収集方法およびシステム

(57) 【要約】

本発明は、ネットコメントの収集方法およびシステム開示した。当該方法は、前記ウェブページのエン트리ーリンクアドレスに対応するウェブページにN個ネットコメントを有するか否かを判断する。ここで、Nは正の整数であるステップと；前記N個ネットコメントを有する場合、前記N個ネットコメントにN個ネットコメントに収集の条件を満たすM個ネットコメントを有するか否かを判断する。ここで、前記MはNより小さいまたは大きい正の整数であるステップと；前記収集の条件を満たすM個ネットコメントを有する場合、前記M個ネットコメントを収集するステップとを備える。

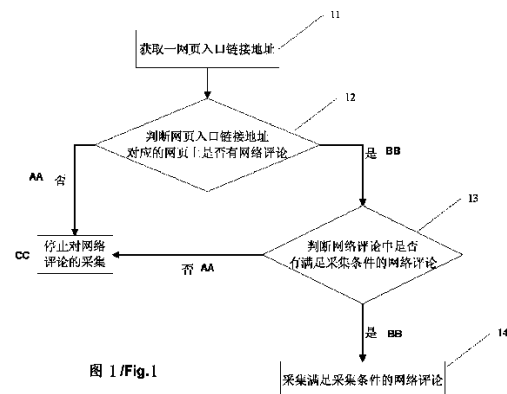


图 1/fig.1

11 OBTAIN A WEB PAGE ENTRY LINK ADDRESS
 12 DETERMINE WHETHER A WEB PAGE CORRESPONDING TO THE WEB PAGE ENTRY
 LINK ADDRESS HAS NETWORK COMMENTS
 13 DETERMINE WHETHER ANY NETWORK COMMENT SATISFYING A COLLECTION
 CONDITION IS AMONG THE NETWORK COMMENTS
 14 COLLECT THE NETWORK COMMENT SATISFYING THE COLLECTION CONDITION
 AA NO
 BB YES
 CC STOP NETWORK COMMENT COLLECTING

【特許請求の範囲】

【請求項 1】

ウェブページのエン트리リンクアドレスを取得するステップと、
 前記ウェブページのエン트리リンクアドレスに対応するウェブページに N 個ネットコメントを有するか否かを判断するステップと、
 前記 N 個ネットコメントを有する場合、前記 N 個ネットコメントに N 個ネットコメントに収集の条件を満たす M 個ネットコメントを有するか否かを判断するステップと、
 前記収集の条件を満たす M 個ネットコメントを有する場合、前記 M 個ネットコメントを収集するステップとを備え、
 前記 N は正の整数であり、
 前記 M は N より小さいまたは大きい正の整数であることを特徴とするネットコメントの収集方法。

10

【請求項 2】

前記ウェブページのエン트리リンクアドレスの取得は、
 前記 N 個ネットコメントがコメントしたトピックがあるトピックページを取得するステップと、
 前記トピックページの特性コードを取得するステップと、
 前記トピックがあるチャンネルの特性コードを取得するステップと、
 前記トピックページの特性コードと前記チャンネルの特性コードを接合し、ウェブページのエン트리リンクアドレスを取得するステップとを備えることを特徴とする請求項 1 に記載のネットコメントの収集方法。

20

【請求項 3】

前記方法は、周期的に前記トピックページのエン트리リンクアドレスを更新することをさらに備えることを特徴とする請求項 2 に記載のネットコメントの収集方法。

【請求項 4】

前記方法は、前記ウェブページでのネットコメントが予定時間を越えても更新がない場合、前記ウェブページのエン트리リンクアドレスを削除することをさらに備えることを特徴とする請求項 1 に記載のネットコメントの収集方法。

【請求項 5】

請求項 1 に記載した前記 N 個ネットコメントに収集の条件を満たす M 個ネットコメントを有するか否かの判断は、具体的に、N と P の差分値を計算し、もし N が P より大きければ、新たに増えたネットコメントを有することを示すステップをさらに備え、
 前記新たに増えたネットコメントの個数は N と P の差分値 M であり、P は前回前記ページにアクセスした際のネットコメントの個数であることを特徴とする請求項 1 に記載のネットコメントの収集方法。

30

【請求項 6】

前記方法は、前記ページの目下のページに備えるネットコメントの個数 L を計算し、もし前記 L が M より小さければ、ページングするページ数を計算し、かつ前記ページ数に対応するページングのリンクを抽出ことさらに備え、

前記 L は正の整数であることを特徴とする請求項 5 に記載のネットコメントの収集方法

40

【請求項 7】

前記方法は、前記 N 個ネットコメントでの各ネットコメントと前記 P 個ネットコメントでの各ネットコメントをそれぞれ比較し、もし比較結果が異なれば、前記比較結果が異なる M 個ネットコメントを抽出することをさらに備えることを特徴とする請求項 5 に記載のネットコメントの収集方法。

【請求項 8】

請求項 1 に記載した前記 N 個ネットコメントに収集の条件を満たす M 個ネットコメントを有するか否かの判断は、具体的に、前記 N 個ネットコメントでの各ネットコメントと前記 P 個ネットコメントでの各ネットコメントをそれぞれ比較し、もし比較結果が異なれば

50

、比較結果が異なるM個ネットコメントは収集の条件を満たすネットコメントであることを確認するステップを備えることを特徴とする請求項1に記載のネットコメントの収集方法。

【請求項9】

前記方法は、抽出した前記M個ネットコメント内容を前記ウェブページと異なるストレージユニットに保存することをさらに備えることを特徴とする請求項1に記載のネットコメントの収集方法。

【請求項10】

ウェブページのエン트리リンクアドレスを取得する、エン트리リンク取得コンポーネントと、

前記ウェブページのエン트리リンクアドレスに対応するウェブページにN個ネットコメントを有するか否かを判断する、第1判断コンポーネントと、

前記N個ネットコメントを有する場合、前記N個ネットコメントにN個ネットコメントに収集の条件を満たすM個ネットコメントを有するか否かを判断する、第2判断コンポーネントと

前記収集の条件を満たすM個ネットコメントを有する場合、前記M個ネットコメントを収集する、内容収集コンポーネントとを備え、

前記Nは正の整数であり、

前記MはNより小さいまたは大きい正の整数であることを特徴とするネットコメントの収集システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は情報検索とデータ統合技術分野に関し、特にネットコメントの収集方法およびシステムに関する。

【背景技術】

【0002】

本出願は、2011年12月13日に中国特許局に提出し、出願番号が201110415749.9であり、発明名称が「ネットコメントの収集方法およびシステム」である中国特許出願を基礎である優先権を主張し、その開示の総てをここに取り込む。

【0003】

現在、インターネット技術の急速な発展に伴い、インターネットはすでに世界で最大のデータベースとなり、人類のほぼ全ての領域を網羅し、人々の情報獲得、情報交換における重要なプラットフォームとなっている。人々が情報を閲覧しやすくするため、インターネットに基づいた情報検索技術もより深い研究と充足した発展を得た。さらに、ネットワーク世論分析、パーティカル検索の評価等のような情報検索に基づいた関連応用も出現した。これら応用技術は全て、まずウェブページをローカルにダウンロードし、不純物を取り除いて分析に必要な内容を抜き出し、最後にこの基礎の上でさらに分析を行う。

【0004】

インターネットに公表する情報について、ネットワークユーザーが情報をブラウジング後に自身の考えを發表し、当該情報に対するコメントを形成する。目下のインターネットの普遍性、広汎性および即効性により、ネットコメントは大衆のある事柄への見方を一定程度代表すると言える。これは、世論分析に大きな意義と応用空間を持たせる。

【発明の概要】

【発明が解決しようとする課題】

【0005】

よって、ネットコメントはすでに多様な応用の重要なデータソースの一つとなり、ネットコメントデータソースの収集は最も基本的な条件である。だが、従来技術において、ネットコメントの収集研究はほぼ白紙であり、ネットコメントに対する効率的、全面的な収集技術に欠ける。

10

20

30

40

50

【課題を解決するための手段】

【0006】

本発明に係る実施形態は、効率的、全面的にネットコメントを収集するネットコメントの収集方法およびシステムを提供する。

【0007】

一方で本発明に係るネットコメントの収集方法は、ウェブページのエン트리リンクアドレス(Entry link address)を取得するステップと；前記ウェブページのエン트리リンクアドレスに対応するウェブページにN個ネットコメントを有するかどうかを判断し、Nが正の整数であるステップと；前記N個ネットコメントを有する場合、前記N個ネットコメントにN個ネットコメントに収集の条件を満たすM個ネットコメントを有するかどうかを判断し、前記MはNより小さいまたは大きい正の整数であるステップと；前記収集の条件を満たすM個ネットコメントを有する場合、前記M個ネットコメントを収集するステップとを備える。

10

【0008】

好ましくは、前記ウェブページのエン트리リンクアドレスの取得は、具体的に、前記N個ネットコメントがコメントしたトピックがあるトピックページを取得するステップと；前記トピックページの特徴コードを取得するステップと；前記トピックがあるチャネルの特徴コードを取得するステップと；前記トピックページの特徴コードと前記チャネルの特徴コードを接合し、ウェブページのエン트리リンクアドレスを取得するステップとを備える。

20

【0009】

好ましくは、周期的に前記トピックページのエン트리リンクアドレスを更新する。

【0010】

好ましくは、前記ウェブページでのネットコメントが予定時間を越えても更新がない場合、前記ウェブページのエン트리リンクアドレスを削除する。

【0011】

好ましくは、上述の前記N個ネットコメントに収集の条件を満たすM個ネットコメントを有するかどうかの判断は、具体的に、NとPの差分値を計算し、もしNがPより大きければ、新たに増えたネットコメントを有することを示す。かつ、前記新たに増えたネットコメントの個数はNとPの差分値Mである。ここで、Pは前回前記ページにアクセスした際のネットコメントの個数であるステップを備える。

30

【0012】

好ましくは、前記ページの目下のページに備えるネットコメントの個数Lを計算する。もし前記LがMより小さければ、ページングするページ数を計算する。かつ前記ページ数に対応するページングのリンクを抜き出す。ここで、Lは正の整数である。

【0013】

好ましくは、前記N個ネットコメントでの各ネットコメントと前記P個ネットコメントでの各ネットコメントをそれぞれ比較する。もし比較結果が異なれば、前記比較結果が異なるM個ネットコメントを抽出する。

【0014】

好ましくは、上述の前記N個ネットコメントに収集の条件を満たすM個ネットコメントを有するかどうかの判断は、具体的に、前記N個ネットコメントでの各ネットコメントと前記P個ネットコメントでの各ネットコメントをそれぞれ比較する。もし比較結果が異なれば、比較結果が異なるM個ネットコメントは収集の条件を満たすネットコメントであることを確認するステップを備える。

40

【0015】

好ましくは、抽出した前記M個ネットコメント内容を前記ウェブページと異なるストレージユニットに保存する。

【0016】

また一方で本発明に係るネットコメントの収集システムは、ウェブページのエントリー

50

リンクアドレスを取得する、エントリーリンク取得コンポーネントと；前記ウェブページのエントリーリンクアドレスに対応するウェブページにN個ネットコメントを有するか否かを判断する。ここで、Nは正の整数である、第1判断コンポーネントと；前記N個ネットコメントを有する場合、前記N個ネットコメントにN個ネットコメントに収集の条件を満たすM個ネットコメントを有するか否かを判断し、前記MはNより小さいまたは大きい正の整数である、第2判断コンポーネントと；前記収集の条件を満たすM個ネットコメントを有する場合、前記M個ネットコメントを収集する、内容収集コンポーネントとを備える。

【発明の効果】

【0017】

10

本発明の有益な効果を以下に記述する。

【0018】

本発明に係る実施形態はネットコメント収集システムを用いてネットコメントを収集し、ネットコメントのエントリーリンクアドレスの取得および収集の条件の設定により、全面的にネットコメントを収集する技術効果を果たす。

【0019】

さらに、比較コンポーネントを用いて、今回抽出した全コメントでの各コメントと前回抽出した全コメントでの各コメントの比較を実現できる。そして、内容抽出コンポーネントを用いて、比較結果が異なるコメントのみを抽出するため、全面的にネットコメントを収集する基礎において効率的な収集の効果を果たすことができる。

20

【図面の簡単な説明】

【0020】

【図1】本発明の実施形態における収集方法のフロー図。

【図2】本発明の図1における収集方法の詳細なフロー図。

【図3】本発明の図1における収集方法の詳細なフロー図。

【図4】本発明の第1実施形態における収集システムアーキテクチャを示す図。

【図5】本発明の第2実施形態における収集システムアーキテクチャを示す図。

【図6】本発明の第3実施形態における収集システムアーキテクチャを示す図。

【図7】本発明の第4実施形態における収集システムアーキテクチャを示す図。

【図8】本発明の別の実施形態における収集システムアーキテクチャを示す図。

30

【発明を実施するための形態】

【0021】

図1は、本発明に係るネットコメントを収集するネットコメントの収集方法はであり、以下のステップ11からステップ14を備える。

【0022】

ステップ11において、ウェブページのエントリーリンクアドレスを取得する。

【0023】

ステップ12において、ウェブページのエントリーリンクアドレスに対応するウェブページにN個ネットコメントを有するか否かを判断する。ここでNは正の整数である。

【0024】

40

ステップ13において、N個ネットコメントを有する場合、N個ネットコメントにN個ネットコメントに収集の条件を満たすM個ネットコメントを有するか否かを判断する。ここで、前記MはNより小さいまたは大きい正の整数である。

【0025】

ステップ14において、収集の条件を満たすM個ネットコメントを有する場合、M個ネットコメントを収集する。

【0026】

ここで、図2に示すように、ステップ11は具体的にさらに以下のステップ111からステップ114を備える。

【0027】

50

ステップ 1 1 1 において、N 個ネットコメントがコメントしたトピックがあるトピックページを取得する。

【0028】

ステップ 1 1 2 において、トピックページの特徴コードを取得する。

【0029】

ステップ 1 1 3 において、トピックがあるチャンネルの特徴コードを取得する。

ステップ 1 1 4 において、トピックページの特徴コードとチャンネルの特徴コードを接合 (S p l i c i n g) し、ウェブページのエンターリンクアドレスを取得する。

【0030】

本発明において、トピックページはニュースがあるページでも良く商品情報があるページでも良い。ここではニュースウェブページを例に挙げ、本実施形態を詳細に説明する。実際には、トピックページは他の情報があるページでも良いが、本発明ではこれを制限しない。

10

【0031】

本実施形態において、ニュースにコメントするコメントページのエンターリンクアドレスは、ニュースページのスクリプトにおける特徴コードにより特定ルールに従い接合後に取得する。例えば、ニュースに対するネットコメントページのエンターリンクアドレスは、ニュースページのスクリプトにより当該ニュースを識別する特徴コード、当該ニュースがあるチャンネル識別する特徴コード、さらにドメイン名および一部の要素 (例えば目下の時間) を加えて接合してできる。前記特徴コードを取得し、かつ個性的なルールを設定し、指定モデルに基づき、ネットコメントページのエンターリンクアドレスをマッチングする。

20

【0032】

引き続き図 2 に示すように、ステップ 1 1 はさらに以下のステップ 1 1 5 を備える。

【0033】

ステップ 1 1 5 において、周期的にウェブページのエンターリンクアドレスを更新する。

【0034】

ステップ 1 1 5 において、ニュースウェブページのホームページバックグラウンドはニュースを再編集する可能性があり、同じ内容のニュースウェブページリンクには変化が生じる。即ち、ニュースの識別およびニュースがあるチャンネルの特徴コードには変化が生じ、ネットコメントのエンターリンクもこれに伴い変化する。新しいネットコメント内容は変化後のネットコメントのエンターリンクによりロードする。さらに、これより前に抽出したネットコメントのエンターリンクアドレスが指定するページには新しいコメントの更新は無いことを意味する。よって、もし元々記録したネットコメントのエンターリンクを引き続き使用し、アクセスすれば、新たに更新したコメント内容を取得できない。故に当該状況において、周期的に目下記録したニュースページリンクを更新する。もしリンクアドレスが変化すれば、サイトは自動的に変化後のニュースウェブページにジャンプする。こうして、新たに獲得したニュースウェブページに基づき、ネットコメントのエンターリンク改めて抽出し、引き続き収集できることは明らかである。即ち、ニュースウェブページのエンターリンクアドレスが更新される場合、ステップ 1 1 1 にジャンプし、実行する。そうでなければ、本フローを終了する。

30

40

【0035】

図 3 に示すように、ステップ 1 3 の具体的なステップは、ステップ 1 3 1 からステップ 1 3 3 を備える。

【0036】

ステップ 1 3 1 において、ウェブページから目下のネットコメントの個数 N を抽出し、N と P の差分値 M を計算する。ここで、P は前回アクセスした当該リンクが抽出したネットコメント個数である。

【0037】

50

ステップ 1 3 2 において、M が 0 より大きいかなかを判断する。

【 0 0 3 8 】

ステップ 1 3 3 において、ステップ 1 3 2 の結果が、M が 0 より大きい場合、M 個ネットコメントを抽出する。

【 0 0 3 9 】

ここで、ステップ 1 3 1 におけるウェブページからの目下のネットコメントの個数 N の抽出は、正規表現によりウェブページから抽出しても良く、他の方法を使用し、抽出しても良いが、本発明はこれを制限しない。最初にネットコメントを収集する場合、P は 0 と等しい。

【 0 0 4 0 】

引き続き図 3 に示すように、ここでステップ 1 3 3 は具体的に以下のステップ 1 3 3 1 からステップ 1 3 3 3 を備える。

【 0 0 4 1 】

ステップ 1 3 3 1 において、ページにおける目下のページに備えられるネットコメントの個数 L を計算する。ここで、L は M より小さいまたは等しい正の整数である。

【 0 0 4 2 】

ステップ 1 3 3 2 において、L が M より小さいかなかを判断する。

【 0 0 4 3 】

ステップ 1 3 3 3 において、ステップ 1 3 3 2 の結果が、L が M より小さい場合、ページングするページ数を計算する。かつ、ページ数に対応するページングのリンクを抽出する。

【 0 0 4 4 】

ここで、ステップ 1 3 3 3 において、ページングの計算公式は：

【 0 0 4 5 】

【 数 1 】

$$P_{count} = \text{ceil}\left(\frac{C_{Update} - C_{Current}}{N_{perpage}}\right)$$

10

20

30

【 0 0 4 6 】

ここで、 P_{count} は、ページングするページ数を示し、 P_{Update} (即ち、M) は、コメント更新数を示し、 $C_{Current}$ (即ち、L) は、目下のウェブページコメント個数を示し、 $N_{perpage}$ は、単数のウェブページコメント数を示す。

【 0 0 4 7 】

引き続き図 3 に示すように、ステップ 1 3 3 はさらに以下のステップ 1 3 3 4 およびステップ 1 3 3 5 を備える。

【 0 0 4 8 】

ステップ 1 3 3 4 において、N 個ネットコメントでの各ネットコメントと P 個ネットコメントでのネットコメントが同じかなかを判断する。

【 0 0 4 9 】

ステップ 1 3 3 5 において、ステップ 1 3 3 4 の結果が、N 個ネットコメントでの各ネットコメントと P 個ネットコメントでの各ネットコメントが同じである場合、比較結果が異なる M 個ネットコメントを抽出する。

【 0 0 5 0 】

ステップ 1 3 3 5 において、抽出した M 個ネットコメント内容は、コメントウェブページの異なるストレージユニットに保存される。ストレージユニットに保存されたネットコメントは集中ブラウジングしやすく、ユーザー収集後のネットコメントを応用しやすい。

【 0 0 5 1 】

本実施形態において、ニュースには即効性があり、一定時間を越えたニュースは意味が

40

50

無いと認識される。同様に、ニュースの附属であるニュースコメントもニュースの失効に伴い失効する。前記原因に基づき、もしネットコメントが予定時間を越えても更新がない場合、当該ニュースコメントリンクを削除し、引き続いて更新はしない。こうして、システムリソースを節約し、より高い作業効率を有することができる。

【0052】

別の実施形態において、N個ネットコメントに収集の条件を満たすM個ネットコメントを有するか否かを判断する場合、前記実施形態におけるNとPの差分値Mを計算する方法を用いなくても良い。つまり、N個ネットコメントでの各ネットコメントとP個ネットコメントでの各ネットコメントをそれぞれ直接比較する。もし比較結果が異なれば、前記比較結果が異なるM個ネットコメントを抽出する。このような収集方法を用いるのは、ニュースウェブページのホームページバックグラウンドが不定期にネットコメントを削除するためである。例えば、システムの最初の収集は15ネットコメントを有し、2回目の収集感覚では、一部の原因によりホームページバックグラウンドは15コメントを全て削除し、同時に30の新しいコメントを加える。つまり1つのウェブページでは15コメントしか表示できないため、ネットコメントの第1ページと第2ページのネットコメントは全て新しいと認識できる。収集周期に達する場合、今回収集した30コメントと前回の15コメントを比較する。こうして、比較の結果が今回収集した30コメントと前回の15コメント全てが異なる。故に、今回30の新しいコメントを収集する。さらに、今回収集した30ネットコメント内容はコメントウェブページの異なるストレージユニットに保存される。ストレージユニットに保存されたネットコメントは集中ブラウジングしやすく、ユーザー収集後のネットコメントを応用しやすい。

10

20

【0053】

本発明の第1実施形態に係るネットコメントデータの収集システムは、図4に示すようにである。図4は本実施形態におけるシステムアーキテクチャであり、当該システムは、エン트리リンク取得コンポーネント10、第1判断コンポーネント20、第2判断コンポーネント30および内容収集コンポーネント40を備える。エン트리リンク取得コンポーネント10は、ウェブページのエン트리リンクアドレスを取得する。第1判断コンポーネント20は、ウェブページのエン트리リンクアドレスに対応するウェブページにN個ネットコメントを有するか否かを判断する。第2判断コンポーネント30は収集の条件を満たすM個ネットコメントを有するか否かを判断する。内容収集コンポーネント40は、ネットコメントを収集する。

30

【0054】

ここで、エン트리リンク取得コンポーネント10は、第1獲得ユニット101、第2獲得ユニット102、第3獲得ユニット103および接合ユニット104を備える。第1獲得ユニット101は、N個ネットコメントがコメントしたトピックがあるトピックページを取得する。第2獲得ユニット102は、トピックページの特徴コードを取得する。第3獲得ユニット103は、トピックがあるチャンネルの特徴コードを取得する。接合ユニット104は、トピックページの特徴コードとチャンネルの特徴コードを接合し、ウェブページのエン트리リンクアドレスを取得する。

40

【0055】

第2判断コンポーネント30による収集の条件を満たすM個ネットコメントを有するか否かの判断は具体的に、ウェブページからN個ネットコメントを抽出し、NとPの差分値Mを計算する。ここで、Pは前回アクセスした当該リンクが抽出したネットコメント個数であるステップをさらに備える。さらに、Mが0より大きいかが否かを判断する。もしMが0より大きければ、M個ネットコメントは収集の条件を満たすコメントであることを示す。第2実施形態において、第1実施形態と異なる点は、システムが周期的にウェブページのエン트리リンクアドレスを更新する、エン트리リンクアドレス更新コンポーネント50をさらに備えることである。本実施形態において、エン트리リンクアドレス更新コンポーネント50は、エン트리リンク取得コンポーネント10と共に運用でき、更新したネットコメントの速やかな収集を実現する。

50

【0056】

第3実施形態において、第1、第2実施形態と異なる点は、システムがウェブページのネットコメントの無更新が予定時間を越えているか否かを判断する。もし超えていれば、ウェブページのエンターリンクアドレスを削除する、ネットコメントページ更新コンポーネント60をさらに備えることである。本実施形態において、ネットコメントページ更新コンポーネント60は、第1判断コンポーネント20と共に運用でき、システム収集効率を高め、いまだ更新しないネットコメントは収集を放棄できる。

【0057】

第2と第3実施形態はそれぞれ図5と図6に示すようにである。実際には、2つの実施形態を結合して使用でき、収集は全面的なネットコメントの収集を実現すると同時にシステムの収集効率を高める。第4実施形態において、第1、第2および第3実施形態と異なる点は、内容収集コンポーネント40がページング抽出コンポーネント401、比較コンポーネント402、内容抽出コンポーネント403およびディスクI/Oコンポーネント404をさらに備えることである。ページング抽出コンポーネント401は、ページングするページ数を計算する。かつ、ページ数に対応するページングのリンクを抽出する。比較コンポーネント402は、前記N個ネットコメントでの各ネットコメントと前記P個ネットコメントでの各ネットコメントをそれぞれ比較する。内容抽出コンポーネント403は、比較結果が異なる場合、前記比較結果が異なるネットコメントを抽出する。ディスクI/Oコンポーネント404は、抽出したネットコメント内容をウェブページの異なるストレージユニットに保存する。本実施形態は図7に示すようにである。

10

20

【0058】

本発明に係る別のネットコメントデータの収集システムは図8に示すようにである。図8は、本実施形態におけるシステムアーキテクチャである。

【0059】

本実施形態と第1実施形態が異なる点は、本実施形態が比較コンポーネント402と内容抽出コンポーネント403を備えないことである。図8に示すように、本実施形態のシステムは、エンターリンク取得コンポーネント80、第1判断コンポーネント81、第2判断コンポーネント82および内容収集コンポーネント83を備える。エンターリンク取得コンポーネント80は、ウェブページのエンターリンクアドレスを取得する。第1判断コンポーネント81は、ウェブページのエンターリンクアドレスに対応するウェブページにネットコメントを有するか否かを判断する。第2判断コンポーネント82は、収集の条件を満たすネットコメントを有するか否かを判断する。内容収集コンポーネント83は、ネットコメントを収集する。

30

【0060】

ここで、エンターリンク取得コンポーネント80は、第1獲得ユニット801、第2獲得ユニット802、第3獲得ユニット803および接合ユニット804を備える。第1獲得ユニット801は、N個ネットコメントがコメントしたトピックがあるトピックページを取得する。第2獲得ユニット802は、トピックページの特徴コードを取得する。第3獲得ユニット803は、トピックがあるチャネルの特徴コードを取得する。接合ユニット804は、トピックページの特徴コードとチャネルの特徴コードを接合し、ウェブページのエンターリンクアドレスを取得する。

40

【0061】

第2判断コンポーネント82は、前記N個ネットコメントでの各ネットコメントと前記P個ネットコメントでの各ネットコメントをそれぞれ比較する。もし比較結果が異なれば、比較結果が異なるM個ネットコメントは収集の条件を満たすネットコメントであることを確認する。

【0062】

内容収集コンポーネント83は、ページング抽出コンポーネント831およびディスクI/Oコンポーネント832を備える。ページング抽出コンポーネント831は、ページングするページ数を計算する。かつ、ページ数に対応するページングのリンクを抽出する

50

。ディスクI/Oコンポーネント832は、抽出したネットコメント内容をウェブページの異なるストレージユニットに保存する。

【0063】

本実施形態において、エンターリンク取得コンポーネント80は、第2実施形態におけるエンターリンクアドレス更新コンポーネント84と結合して共に応用でき、比較的全面的なネットコメントの収集を実現する。第1判断コンポーネント81は、第3実施形態におけるネットコメントページ更新コンポーネント85と結合して共に応用でき、全面的、効率的にネットコメントの収集を実現する。

【0064】

前記第1、第2、第3、第4および別の実施形態におけるシステムは、本発明が提供したネットコメント収集方法の実施形態における方法およびその各種変化の形式の記述に基づき、実施できる。明細書を簡潔にするため、ここでは説明を繰り返さない。

10

【0065】

本発明の実施形態は、ネットコメント収集システムを用いて、ネットコメントを収集し、ネットコメントのエンターリンクアドレスの取得および収集の条件の設定により全面的なネットコメントを収集する技術効果を果たす。

【0066】

さらに、比較コンポーネントを用いて、今回抽出した全コメントでの各コメントと前回抽出した全コメントでの各コメントの比較を実現できる。そして、内容抽出コンポーネントを用いて、比較結果が異なるコメントのみを抽出する。よって、全面的にネットコメントを収集する基礎の上に効率的な収集の効果を果たすことができる。

20

【0067】

以上は本発明の実施形態の方法、装置（システム）、およびコンピュータプログラム製品のフロー図および/またはブロック図によって、本発明を記述した。理解すべきことは、コンピュータプログラム指令によって、フロー図および/またはブロック図における各フローおよび/またはブロックと、フロー図および/またはブロック図におけるフローおよび/またはブロックの結合を実現できる。プロセッサはこれらのコンピュータプログラム指令を、汎用コンピュータ、専用コンピュータ、組込み式処理装置、或いは他のプログラム可能なデータ処理装置設備の処理装置器に提供でき、コンピュータ或いは他のプログラム可能なデータ処理装置のプロセッサは、これらのコンピュータプログラム指令を実行し、フロー図における一つ或いは複数のフローおよび/またはブロック図における一つ或いは複数のブロックに指定する機能を実現する。

30

【0068】

これらのコンピュータプログラム指令は又、コンピュータ或いは他のプログラム可能なデータ処理装置を特定方式で動作させるコンピュータ読取記憶装置に記憶できる。これによって、指令を含む装置は当該コンピュータ読取記憶装置内の指令を実行でき、フロー図における一つ或いは複数のフローおよび/またはブロック図における一つ或いは複数のブロックに指定する機能を実現する。

【0069】

これらコンピュータプログラム指令はさらに、コンピュータ或いは他のプログラム可能なデータ処理装置設備に実装もできる。コンピュータプログラム指令が実装されたコンピュータ或いは他のプログラム可能設備は、一連の操作ステップを実行することによって、関連の処理を実現し、コンピュータ或いは他のプログラム可能な設備において実行される指令によって、フロー図における一つ或いは複数のフローおよび/またはブロック図における一つ或いは複数のブロックに指定する機能を実現する。

40

【0070】

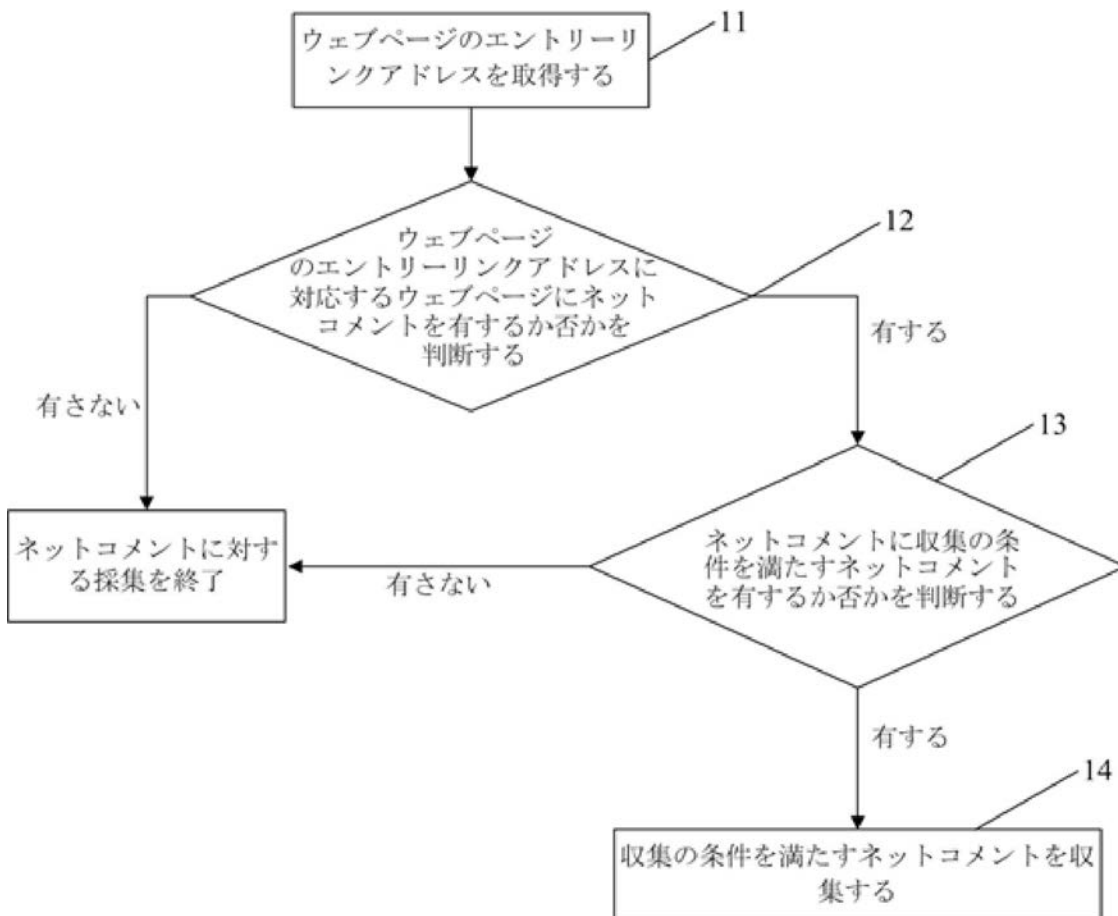
上述した実施形態に記述された技術的な解決手段を改造し、或いはその中の一部の技術要素を置換することもできる。そのような、改造と置換は本発明の各実施形態の技術の範囲から逸脱するとは見なされない。

【0071】

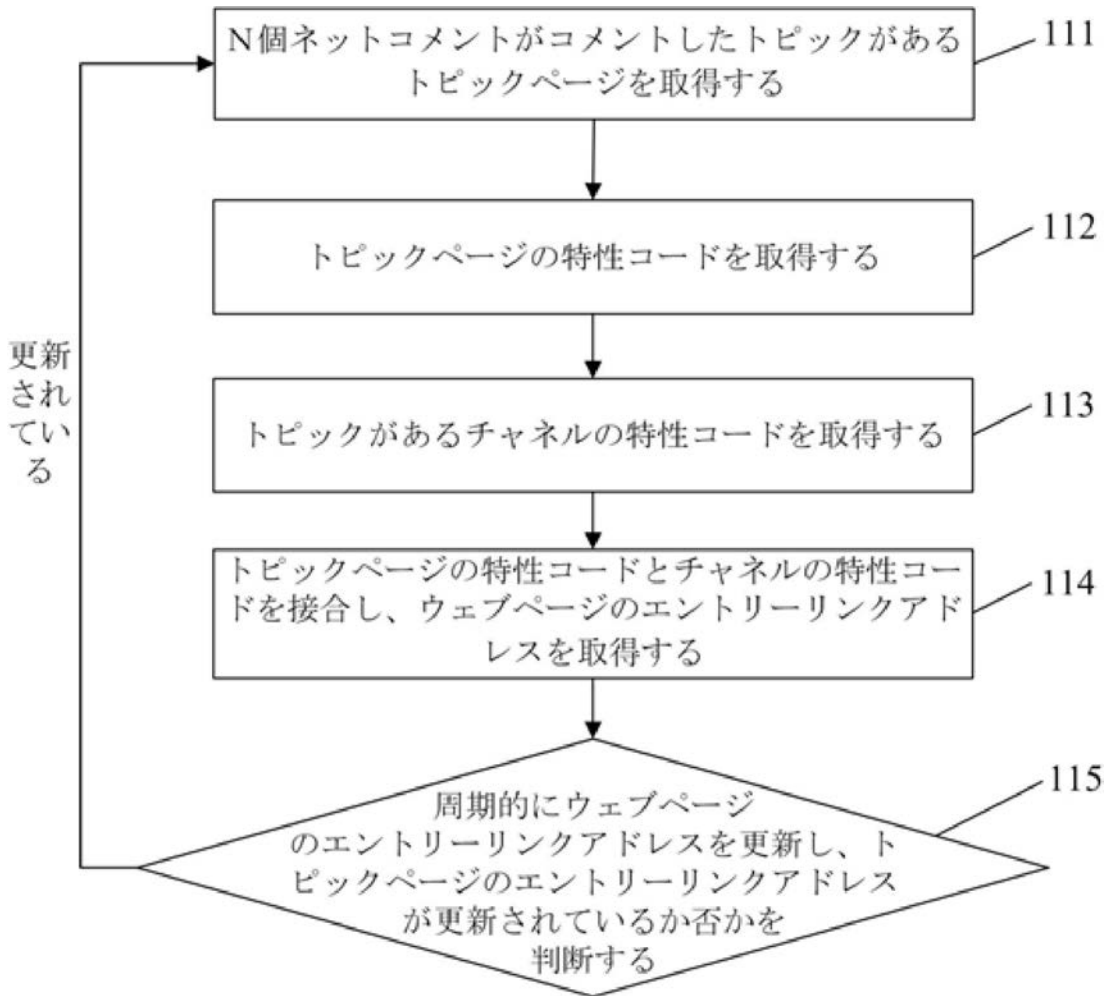
50

無論、当業者によって、上述した実施形態に記述された技術的な解決手段を改造し、或いはその中の一部の技術要素を置換することもできる。そのような、改造と置換は本発明の各実施形態の技術の範囲から逸脱するとは見なされない。そのような改造と置換は、すべて本発明の請求の範囲に属する。

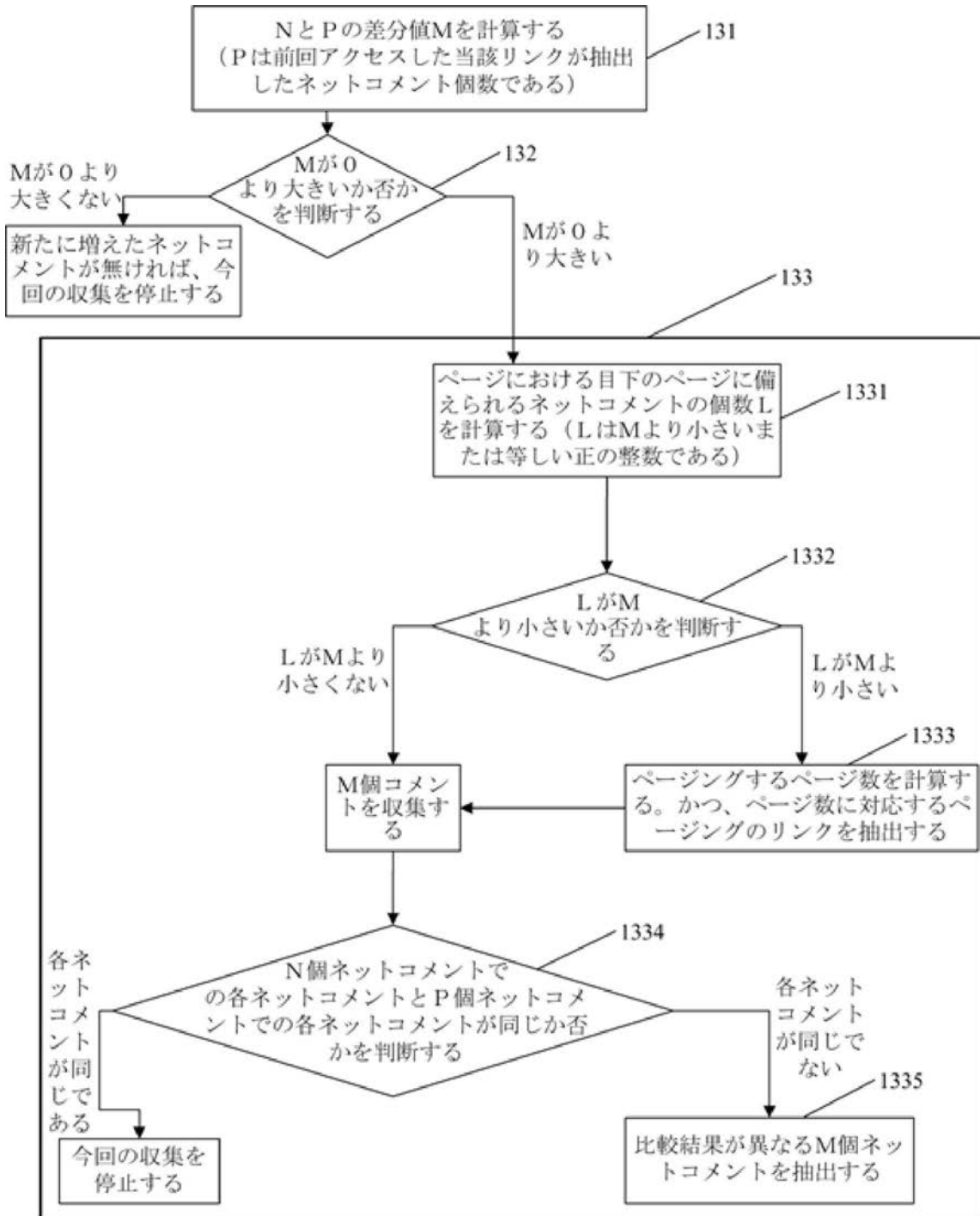
【図1】



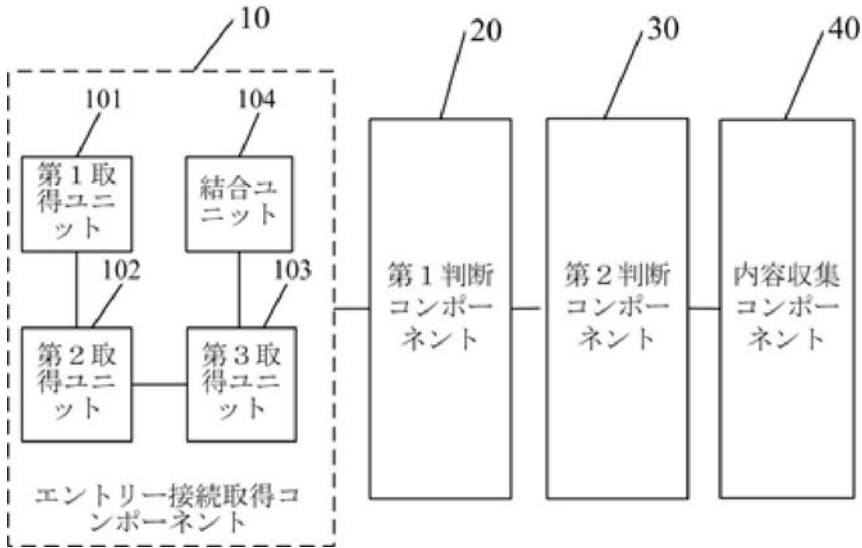
【図2】



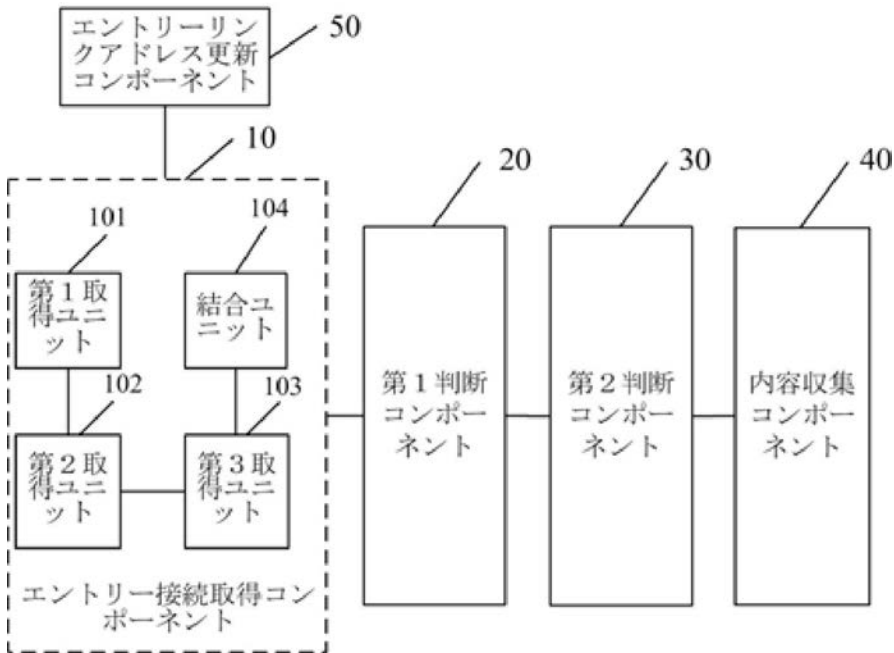
【図3】



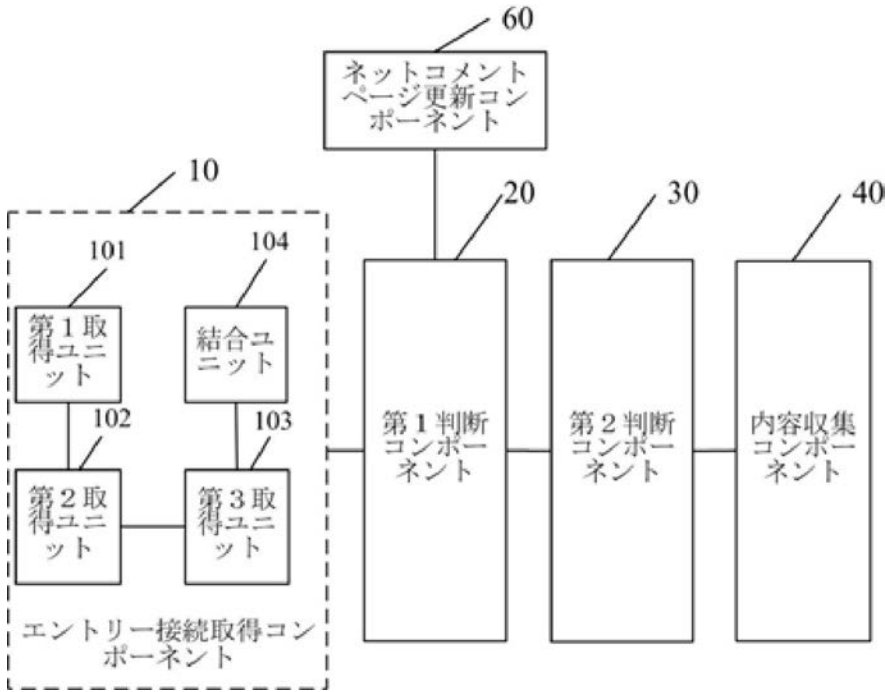
【図4】



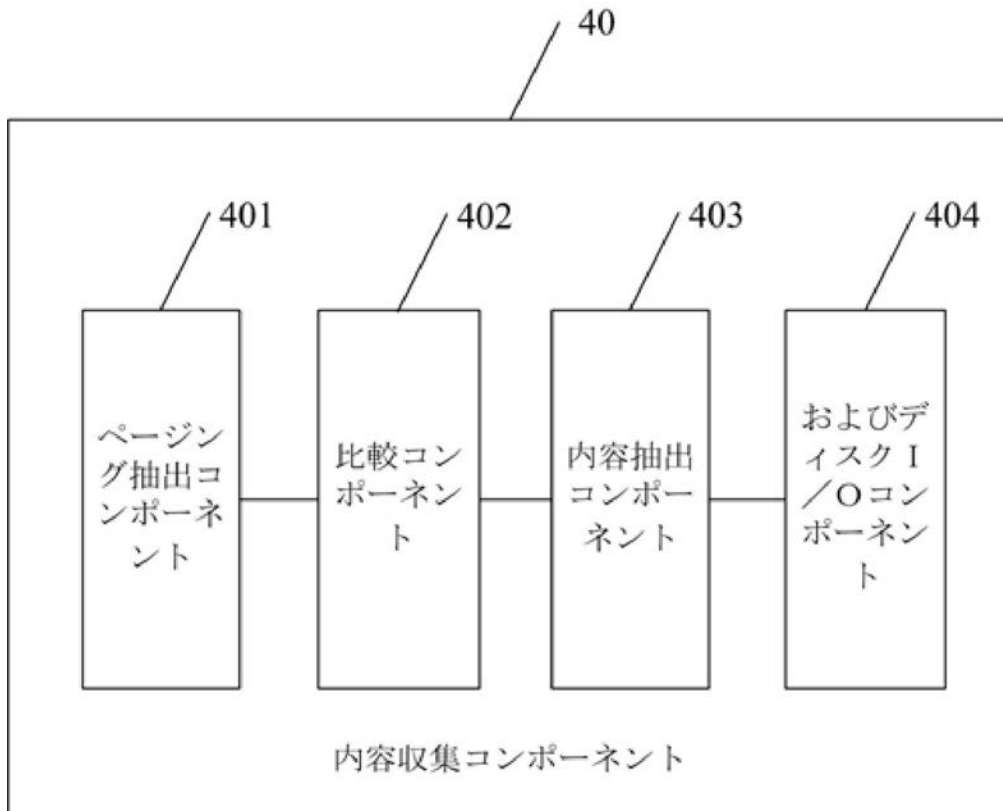
【図5】



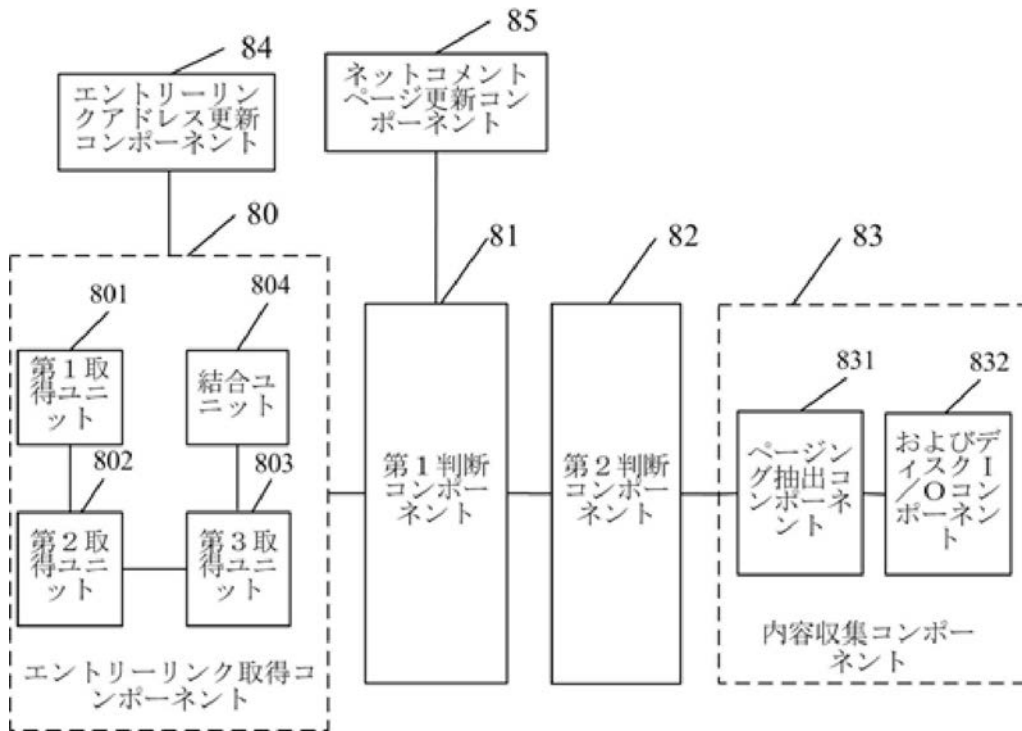
【図6】



【図7】



【 図 8 】



【 国际调查报告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/CN2012/086575		
A. CLASSIFICATION OF SUBJECT MATTER				
G06F 17/30 (2006.01) i				
According to International Patent Classification (IPC) or to both national classification and IPC				
B. FIELDS SEARCHED				
Minimum documentation searched (classification system followed by classification symbols)				
IPC: G06F17/-; G06F15/-; H04L12/-; H04M1/-; H04Q7/-; H04M3/-;				
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched				
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)				
CNABS, SIPOABS, DWPI: web, network, data, information, link, criticism, acquisition, review				
C. DOCUMENTS CONSIDERED TO BE RELEVANT				
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.		
X	CN 101515269 A (CHINESE ACAD SCI AUTOMATION INST) 26 August 2009 (26.08.2009) description, page 2, line 8 to page 7, line 1, figure 1 and the abstract	1-2, 5-10		
Y		3-4		
Y	CN 101178713 A (TENCENT SCI&TECHNOLOGY SHENZHEN CO LTD) 14 May 2008 (14.05.2008) description, page 2, line 2, line 15 to line 23	3-4		
A	US 6707470 B1 (NEC COPRORATION) 16 March 2004 (16.03.2004) , the whole document	1-10		
A	CN 101436196 A (UNIV BEIJING POSTS&TELECOM) 20 May 2009 (20.05.2009) description, page 2, line 8 to page 8, line 24, figure 1 and the abstract	1-10		
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.				
* Special categories of cited documents: <table style="width: 100%; border: none;"> <tr> <td style="width: 50%; vertical-align: top;"> "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed </td> <td style="width: 50%; vertical-align: top;"> "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family </td> </tr> </table>			"A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family
"A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family			
Date of the actual completion of the international search 26 February 2013 (26.02.2013)		Date of mailing of the international search report 21 March 2013 (21.03.2013)		
Name and mailing address of the ISA State Intellectual Property Office of the P. R. China No. 6, Xitucheng Road, Jimenqiao Haidian District, Beijing 100088, China Facsimile No. (86-10) 62019451		Authorized officer GONG , JinLing Telephone No. (86-10) 62413492		

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/CN2012/086575

Patent Documents referred in the Report	Publication Date	Patent Family	Publication Date
CN 101515269 A	26.08.2009	CN 101515269 B	31.08.2011
CN 101178713 A	14.05.2008	None	
US 6707470 B1	16.03.2004	JP 2000330856 A	30.11.2000
CN 101436196 A	20.05.2009	CN 101436196 B	08.12.2010

国际检索报告		国际申请号 PCT/CN2012/086575
A. 主题的分类		
G06F17/30(2006.01)i		
按照国际专利分类(IPC)或者同时按照国家分类和 IPC 两种分类		
B. 检索领域		
检索的最低限度文献(标明分类系统和分类号)		
IPC: G06F17/-,G06F15/-,H04L29/-,H04L12/-,H04M1/-,H04Q7/-,H04M3/-.		
包含在检索领域中的除最低限度文献以外的检索文献		
在国际检索时查阅的电子数据库(数据库的名称, 和使用的检索词(如使用))		
CNABS,SIPOABS,DWPI; 网络评论, 采集, 网页, 链接, 评论, web, network, data, information, link, Criticism, acquisition, review,		
C. 相关文件		
类 型*	引用文件, 必要时, 指明相关段落	相关的权利要求
X	CN101515269A (中国科学院自动化研究所)26.8月2009(26.08.2009), 说明书摘要, 说明书第2页第8行至说明书第7页第一行, 附图1	1-2,5-10
Y		3-4
Y	CN101178713A (腾讯科技(深圳)有限公司)14.5月2008(14.05.2008) 说明书第2页第2行, 第15-23行	3-4
A	US6707470B1(NEC CORPORATION)16.3月2004(16.03.2004),全文	1-10
A	CN101436196A(北京邮电大学),20.5月2009(20.05.2009), 说明书摘要, 说明书第2页第8行至说明书第8页最后一行, 附图1	1-10
<input type="checkbox"/> 其余文件在 C 栏的续页中列出。 <input checked="" type="checkbox"/> 见同族专利附件。		
* 引用文件的具体类型: “A” 认为不特别相关的表示了现有技术一般状态的文件 “E” 在国际申请日的当天或之后公布的在先申请或专利 “L” 可能对优先权要求构成怀疑的文件, 或为确定另一篇引用文件的公布日而引用的或者因其他特殊理由而引用的文件(如具体说明的) “O” 涉及口头公开、使用、展览或其他方式公开的文件 “P” 公布日先于国际申请日但迟于所要求的优先权日的文件		
“T” 在申请日或优先权日之后公布, 与申请不相抵触, 但为了理解发明之理论或原理的在后文件 “X” 特别相关的文件, 单独考虑该文件, 认定要求保护的发明不是新颖的或不具有创造性 “Y” 特别相关的文件, 当该文件与另一篇或者多篇该类文件结合并且这种结合对于本领域技术人员为显而易见时, 要求保护的发明不具有创造性 “&” 同族专利的文件		
国际检索实际完成的日期 26.2月2013(26.02.2013)		国际检索报告邮寄日期 21.3月2013(21.03.2013)
ISA/CN 的名称和邮寄地址: 中华人民共和国国家知识产权局 中国北京市海淀区蓟门桥西土城路6号100088 传真号: (86-10)62019451		受权官员 龚锦玲 电话号码: (86-10) 62413492

国际检索报告
关于同族专利的信息国际申请号
PCT/CN2012/086575

检索报告中引用的 专利文件	公布日期	同族专利	公布日期
CN101515269A	26.08.2009	CN101515269B	31.08.2011
CN101178713A	14.05.2008	无	
US6707470B1	16.03.2004	JP2000330856A	30.11.2000
CN101436196A	20.05.2009	CN101436196B	08.12.2010

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC

(71)出願人 507232478

北京大学

PEKING UNIVERSITY

中華人民共和国北京市 海淀区 頤和園路5号

No. 5, Yiheyuan Road, Haidian District, Beijing 100871, China

(71)出願人 507232456

北京北大方正 電子有限公司

BEIJING FOUNDER ELECTRONICS CO., LTD.

中華人民共和国北京市 海淀区上地五街9号方正大厦

Founder Building, No. 9, Shangdiwu Street, Haidian District, Beijing 100085, China

(74)代理人 110001243

特許業務法人 谷・阿部特許事務所

(72)発明者 ジャン タオ

中華人民共和国 100871 ベイジン ハイディアン チェンファー ロード ナンバー 29
8 ジョングアンツン ファウンダー ビルディング 5 フロア

(72)発明者 ユー シアオミン

中華人民共和国 100871 ベイジン ハイディアン チェンファー ロード ナンバー 29
8 ジョングアンツン ファウンダー ビルディング 5 フロア

(72)発明者 ヤン ジエンウー

中華人民共和国 100871 ベイジン ハイディアン チェンファー ロード ナンバー 29
8 ジョングアンツン ファウンダー ビルディング 5 フロア

Fターム(参考) 5B084 BB12