

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5240475号  
(P5240475)

(45) 発行日 平成25年7月17日 (2013. 7. 17)

(24) 登録日 平成25年4月12日 (2013. 4. 12)

(51) Int. Cl.

F I

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 1 7 0 A

G 0 6 F 17/30 3 5 0 C

G 0 6 F 17/30 4 1 5

請求項の数 30 (全 21 頁)

(21) 出願番号 特願2009-509952 (P2009-509952)  
 (86) (22) 出願日 平成19年4月24日 (2007. 4. 24)  
 (65) 公表番号 特表2009-535747 (P2009-535747A)  
 (43) 公表日 平成21年10月1日 (2009. 10. 1)  
 (86) 国際出願番号 PCT/US2007/067319  
 (87) 国際公開番号 W02007/130818  
 (87) 国際公開日 平成19年11月15日 (2007. 11. 15)  
 審査請求日 平成22年4月22日 (2010. 4. 22)  
 (31) 優先権主張番号 11/381, 214  
 (32) 優先日 平成18年5月2日 (2006. 5. 2)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者 508137947  
 エクセジー・インコーポレイテッド  
 アメリカ合衆国、ミズーリ・63127、  
 セント・ルイス、サウス・ガイヤー・ロー  
 ド・3668、スイート・300  
 (74) 代理人 100103920  
 弁理士 大崎 勝真  
 (74) 代理人 100114188  
 弁理士 小野 誠  
 (74) 代理人 100140523  
 弁理士 渡邊 千尋  
 (74) 代理人 100119253  
 弁理士 金山 賢教  
 (74) 代理人 100124855  
 弁理士 坪倉 道明

最終頁に続く

(54) 【発明の名称】 近似パターン合致の方法および装置

(57) 【特許請求の範囲】

【請求項 1】

複数のデータシンボルを含むデータ文字列が複数のパターンのうちのいずれかとの近似合致を含むか否かの決定を容易化するために、データ文字列内の複数のデータサブ文字列を処理する方法であって、

複数の可能性のある合致を検出するためにフィルタ回路にデータサブ文字列について問い合わせをするステップであって、各可能性のある合致がデータサブ文字列とパターン断片との間の可能性のある合致を表し、複数のパターンの各パターンが複数の対応するパターン断片を含み、フィルタ回路がパターン断片でプログラムされている、問い合わせをするステップと、

複数のパターンセットを決定するために、検出された可能性のある合致を削減ステージに適用するステップであって、各パターンセットが、検出された可能性のある合致と対応し、(1) 対応する可能性のある合致を生じたパターン断片に対応するパターンを表すデータおよび(2) 該パターンに関連する許容可能な誤りを表すデータを含む、適用するステップと、

配信されたパターンセット内の許容可能な誤りデータを考慮して、配信されたパターンセット内のパターンのいずれかと、データ文字列の少なくとも一部が近似合致するか否かを決定するために、決定されたパターンセットを近似合致エンジンに配信するステップとを含む、方法。

【請求項 2】

フィルタ回路がブルームフィルタ回路である、請求項 1 に記載の方法。

【請求項 3】

配信されたパターンセット内の許容可能な誤りデータを考慮して、データ文字列部分と配信されたパターンセット内のパターンとの間で何らかの近似合致が存在するか否かを近似合致エンジンを用いて判定するステップをさらに含む、請求項 1 または 2 に記載の方法。

【請求項 4】

少なくとも複数のパターンであるパターンを複数のパターン断片にスライスするステップと、

パターン断片を用いてブルームフィルタ回路をプログラムするステップとをさらに含む、請求項 2 に記載の方法。

【請求項 5】

再構成可能なハードウェアを用いて、判定するステップを実行するステップをさらに含む、請求項 3 に記載の方法。

【請求項 6】

問い合わせするステップおよび適用するステップをパイプライン方式で実行するステップをさらに含む、請求項 1 から 5 のいずれか一項に記載の方法。

【請求項 7】

再構成可能なハードウェアを用いて、問い合わせするステップおよび適用するステップを実行するステップをさらに含む、請求項 1 から 6 のいずれか一項に記載の方法。

【請求項 8】

ブルームフィルタ回路が複数の並列ブルームフィルタ回路を含み、並列ブルームフィルタ回路の各々が、各並列ブルームフィルタ回路についてのパターン断片長が他の並列ブルームフィルタ回路についての他のパターン断片長と異なり、問い合わせステップが、複数のデータサブ文字列を並列ブルームフィルタ回路に並列に同時に提供するステップを含む、請求項 4 に記載の方法。

【請求項 9】

ブルームフィルタ回路が、データ文字列の少なくとも一部を記憶するシフトレジスタをさらに含み、方法が、シフトレジスタを通してデータ文字列のデータシンボルを流すステップと、問い合わせするステップのためにシフトレジスタからデータサブ文字列を読むステップとをさらに含む、請求項 2、4 又は 8 のいずれか一項に記載の方法。

【請求項 10】

各並列ブルームフィルタ回路が、複数のビットベクトルを記憶するように構成されており、各ビットベクトルが複数のビットを含み、各ビットが値を有し、問い合わせするステップが、

少なくとも 1 つのハッシュキーを生成するために、各データサブ文字列を少なくとも 1 つのハッシュ関数に適用するステップと、

少なくとも 1 つの生成されたハッシュキーに基づいて、並列ブルームフィルタ回路からビットベクトルを取り出すステップと、

少なくとも 1 つの生成されたハッシュキーに基づいて、取り出されたビットベクトルの複数のビットの位置を選択するステップと、

選択されたビットの位置にあるビットの値に基づいて、該データサブ文字列とパターン断片との間に可能性のある合致が存在するか否かを判定するステップとをさらに含む、請求項 8 に記載の方法。

【請求項 11】

各データサブ文字列を少なくとも 1 つのハッシュ関数に適用するステップが、複数のハッシュキーを生成するために複数のハッシュ関数に各データサブ文字列を適用するステップを含み、ビットベクトルを取り出すステップが、生成されたハッシュキーに基づいて並列ブルームフィルタ回路からビットベクトルを取り出すステップを含み、ビットの位置を選択するステップが、複数の他の生成されたハッシュキーに基づいて、取り出されたビッ

10

20

30

40

50

トベクトルの複数のビットの位置を選択するステップを含む、請求項 10 に記載の方法。

【請求項 12】

各データサブ文字列を少なくとも 1 つのハッシュ関数に適用するステップが、単一のハッシュキーを生成するために単一のハッシュ関数に各データサブ文字列を適用するステップを含み、ビットベクトルを取り出すステップが、生成されたハッシュキーの一部に基づいて並列ブルームフィルタ回路からビットベクトルを取り出すステップを含み、ビットの位置を選択するステップが、生成されたハッシュキーの別の部分に基づいて、取り出されたビットベクトルの複数のビットの位置を選択するステップを含む、請求項 10 に記載の方法。

【請求項 13】

削減ステージが、(1) 各パターン断片識別子が問い合わせするステップの結果として生成されたデータによってインデックス付けされる、複数のパターン断片識別子を記憶し、(2) 各パターン識別子がパターン断片識別子によってインデックス付けされる、複数のパターン識別子を記憶し、(3) 各パターンセット対がパターン識別子によってインデックス付けされる、複数のパターンセット対を記憶するように構成されており、各パターンセット対が、(a) パターンを表すデータおよび (b) 該パターンに関連する許容可能な誤りを表すデータを含み、可能性のある合致を適用するステップが、

問い合わせするステップの結果として生じたデータに基づいて、記憶されているパターン断片識別子の少なくとも 1 つを取り出すステップと、

各取り出されたパターン断片識別子に基づいて、記憶されているパターン識別子の少なくとも 1 つを取り出すステップと、

各取り出されたパターン識別子に基づいて、記憶されているパターンセット対の少なくとも 1 つを取り出すステップとを含み、取り出された少なくとも 1 つのパターンセット対が、パターンセットの少なくとも一部を規定する、請求項 1 から 12 のいずれか一項に記載の方法。

【請求項 14】

パターン断片識別子を取り出すステップが、少なくとも 1 つの生成されたハッシュキーに対応するデータに基づいて、記憶されているパターン断片識別子の少なくとも 1 つを取り出すステップを含む、請求項 13 に記載の方法。

【請求項 15】

削減ステージが、(1) 各パターン識別子がパターン断片によってインデックス付けされる、複数のパターン識別子を記憶し、(2) 各パターンセット対がパターン識別子によってインデックス付けされる、複数のパターンセット対を記憶するように構成されており、各パターンセット対が、(a) パターンを表すデータおよび (b) 該パターンに関連する許容可能な誤りを表すデータを含み、可能性のある合致を適用するステップが、

可能性のある合致を生じたデータサブ文字列に基づいて、記憶されているパターン識別子の少なくとも 1 つを取り出すステップと、

各取り出されたパターン識別子に基づいて、記憶されているパターンセット対の少なくとも 1 つを取り出すステップとを含み、取り出された少なくとも 1 つのパターンセット対が、パターンセットの少なくとも一部を規定する、請求項 1 から 12 のいずれか一項に記載の方法。

【請求項 16】

複数のデータシンボルを含むデータ文字列が複数のパターンのうちのいずれかとの近似合致を含むか否かの決定を容易化するために、データ文字列内の複数のデータサブ文字列を処理するシステムであって、

複数の可能性のある合致を検出するためにデータサブ文字列によって問い合わせされるように構成されているフィルタ回路であって、各可能性のある合致がデータサブ文字列とパターン断片との間の可能性のある合致を表し、複数のパターンの各パターンが複数の対応するパターン断片を含み、フィルタ回路がパターン断片でプログラムされている、フィルタ回路と、

10

20

30

40

50

フィルタ回路と通信する削減ステージであって、削減ステージが、(1) 検出された可能性のある合致を処理し、(2) 処理に応じて複数のパターンセットを決定し、各パターンセットが可能性のある合致に対応して(a) 対応する可能性のある合致を生じたパターン断片に対応するパターンを表すデータと、(b) 該パターンに関連する許容可能な誤りを表すデータとを含み、(3) 決定されたパターンセット内の許容可能な誤りデータを考慮して、データ文字列の少なくとも一部が、決定されたパターンセット内のいずれかパターンと近似合致するか否かを判定するために、近似合致エンジンに決定されたパターンセットを出力するように構成されている、削減ステージとを含む、システム。

【請求項 17】

フィルタ回路がブルームフィルタ回路である、請求項 16 に記載のシステム。

10

【請求項 18】

削減ステージと通信しており、決定されたパターンセット内の許容可能な誤りデータを考慮して、データ文字列部分と決定されたパターンセット内のパターンとの間で何らかの近似合致が存在するか否かを判定するように構成されている近似合致エンジンをさらに含む、請求項 16 または 17 に記載のシステム。

【請求項 19】

ブルームフィルタ回路および削減ステージがパイプライン方式で動作するように構成されている、請求項 17 に記載のシステム。

【請求項 20】

ブルームフィルタ回路および削減ステージが再構成可能なハードウェアを用いて実装されている、請求項 17 または 19 に記載のシステム。

20

【請求項 21】

近似合致エンジンが再構成可能なハードウェアを用いて実装されている、請求項 18 に記載のシステム。

【請求項 22】

パターン断片が、異なる長さの複数のパターン断片である、請求項 16 から 21 のいずれか一項に記載のシステム。

【請求項 23】

同じパターンに対応するパターン断片が互いに重なり合わない、請求項 16 から 21 のいずれか一項に記載のシステム。

30

【請求項 24】

ブルームフィルタ回路が複数の並列ブルームフィルタ回路を含み、各並列ブルームフィルタ回路についてのパターン断片長が他の並列ブルームフィルタ回路についての他のパターン断片長と異なり、ブルームフィルタ回路が、複数のデータサブ文字列を並列ブルームフィルタ回路に並列に同時に提供するようにさらに構成されている、請求項 17、19 または 20 のいずれか一項に記載のシステム。

【請求項 25】

ブルームフィルタ回路が、データ文字列の少なくとも一部を記憶するシフトレジスタをさらに含み、ブルームフィルタ回路が、(1) シフトレジスタを通してデータ文字列のデータシンボルを流し、(2) 並列ブルームフィルタ回路への配信のためにシフトレジスタからデータサブ文字列を読むようにさらに構成されている、請求項 24 に記載のシステム。

40

【請求項 26】

各並列ブルームフィルタ回路が、複数のビットベクトルを記憶するように構成されており、各ビットベクトルが複数のビットを含み、各ビットが値を有し、ブルームフィルタ回路が、(1) 少なくとも 1 つのハッシュキーを生成するために、各データサブ文字列を少なくとも 1 つのハッシュ関数に適用し、(2) 少なくとも 1 つの生成されたハッシュキーに基づいて、並列ブルームフィルタ回路からビットベクトルを取り出し、(3) 少なくとも 1 つの生成されたハッシュキーに基づいて、取り出されたビットベクトルの複数のビットの位置を選択し、(4) 選択されたビットの位置にあるビットの値に基づいて、該デー

50

タサブ文字列とパターン断片との間に可能性のある合致が存在するか否かを判定するようにさらに構成されている、請求項 24 に記載のシステム。

【請求項 27】

各並列ブルームフィルタ回路が、(1)複数のハッシュキーを生成するために複数のハッシュ関数に各データサブ文字列を適用し、(2)生成されたハッシュキーに基づいて並列ブルームフィルタ回路からビットベクトルを取り出し、(3)複数の他の生成されたハッシュキーに基づいて、取り出されたビットベクトルの複数のビットの位置を選択するようにさらに構成されている、請求項 26 に記載のシステム。

【請求項 28】

各並列ブルームフィルタ回路が、(1)単一のハッシュキーを生成するために単一のハッシュ関数に各データサブ文字列を適用し、(2)生成されたハッシュキーの一部に基づいて並列ブルームフィルタ回路からビットベクトルを取り出し、(3)生成されたハッシュキーの別の部分に基づいて、取り出されたビットベクトルの複数のビットの位置を選択するようにさらに構成されている、請求項 26 に記載のシステム。

【請求項 29】

削減ステージが第1のテーブル、第2のテーブルおよび第3のテーブルを含み、第1のテーブルは、各パターン断片識別子がブルームフィルタ回路によって生成されたデータによってインデックス付けされる、複数のパターン断片識別子を記憶するように構成されており、第2のテーブルは、各パターン識別子がパターン断片識別子によってインデックス付けされる、複数のパターン識別子を記憶するように構成されており、第3のテーブルは、各パターンセット対がパターン識別子によってインデックス付けされる、複数のパターンセット対を記憶するようにさらに構成されており、各パターンセット対が、(a)パターンを表すデータおよび(b)該パターンに関連する許容可能な誤りを表すデータを含み、削減ステージが、(1)ブルームフィルタ回路によって生じたデータに基づいて、第1のテーブルから記憶されているパターン断片識別子の少なくとも1つを取り出し、(2)各取り出されたパターン断片識別子に基づいて、第2のテーブルから記憶されているパターン識別子の少なくとも1つを取り出し、(3)各取り出されたパターン識別子に基づいて、第3のテーブルから記憶されているパターンセット対の少なくとも1つを取り出すようにさらに構成されており、少なくとも1つの取り出されたパターンセット対が、パターンセットの少なくとも一部を規定する、請求項 16 から 28 のいずれか一項に記載のシステム。

【請求項 30】

削減ステージが第1のデータ構造および第2のデータ構造を含み、第1のデータ構造は、各パターン識別子がパターン断片によってインデックス付けされる、複数のパターン識別子を記憶するように構成されており、第2のデータ構造は、各パターンセット対がパターン識別子によってインデックス付けされる、複数のパターンセット対を記憶するように構成されており、各パターンセット対が、(a)パターンを表すデータおよび(b)該パターンに関連する許容可能な誤りを表すデータを含み、削減ステージが、(1)可能性のある合致を生じたデータサブ文字列に基づいて、第1のデータ構造から記憶されているパターン識別子の少なくとも1つを取り出し、(2)各取り出されたパターン識別子に基づいて、第2のデータ構造から記憶されているパターンセット対の少なくとも1つを取り出すようにさらに構成されており、少なくとも1つの取り出されたパターンセット対が、パターンセットの少なくとも一部を規定する、請求項 16 から 28 のいずれか一項に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、パターンの大きなセットを用いた近似パターン合致の分野に関する。特に、本発明は、多数のグループのパターンを用いた近似パターン合致のためのスケーラブルフィルタ回路および削減ステージに関する。

## 【背景技術】

## 【0002】

近似パターンまたは文字列合致は、多数の重要なアプリケーションで生じる重要な問題である。限定ではないが、これらは、コンピュータ利用生命工学、データベースおよびコンピュータ通信を含み得る。このタスクは、通常、特定の数の誤りを許しながらの、特定のパターンまたはパターンセットの間での合致の検索を含む。一例として、2つの誤りを許しながら、単語「queueing」を検索することを望むことが可能である。これは、1つの文字が挿入されている単語「queueing」、および1つの文字が置換されて1つの文字が削除されている「cueing」等の結果を返す可能性がある。特定の数の誤りを許可することによって、これは、通常のとおりの変化または誤りを捕捉して検索し、それでも所望のパターンの発見を可能にする。近似パターン合致は、複雑なタスクであるだけでなく、途方もない量のコンピュータリソースを必要とする。

10

## 【0003】

通常、高速のフィルタリングステップがあり、次に全部の近似合致機能を実行する検証ステップが続く。この従来技術のフィルタ技術の一例は、図1を参照すると示されており、図番10で概略的に示されている。この通常の手法は、図番12で示されるようなパターン「P」を、重複しない一連のサブパターンである $k+1$ 個のパターン断片にスライスし、テキストとパターン断片との間の正確な合致を検索する。この場合、「k」は、最大編集距離 $ed(T_{i \dots j}, P)$ であり、この限定でない例では図番14で示されるように数字2(2)で示されている、許容可能な誤りの数と等しい。

20

## 【0004】

データ文字列 $T_{i \dots j}$  16は、次に、データ文字列16のうちの少なくとも1つのサブ文字列が、パターン「P」12と関連する重複しないサブパターンのうちの少なくとも1つと合致する出現に関して解析される。この手法は、以下の特性に依存する：

a. 文字列 $S = T_{a \dots b}$ が最大kの誤りでパターンPと合致し、 $P = p_1 \dots p_j$  (重複しないサブパターンの配列)である場合、Sのいくつかのサブ文字列は複数の $p_i$ のうちの少なくとも1つと最大

## 【数1】

$$\lfloor k/j \rfloor$$

30

の誤りで合致する。

b.  $ed(T_{i \dots j}, P) = k$ である文字位置 $i \dots j$ が存在する場合、 $T_{j-m+1 \dots j}$ は、Pのうちの少なくとも $m-k$ 個の文字を含み、ここでmはパターンの(文字の)大きさである。

c. それゆえに、Pを $k+1$ 個の断片(重複しないサブパターン)にスライスする場合、断片のうちの少なくとも1つが正確に合致するはずである。

## 【0005】

それゆえに、「P」12を、誤りの総数「k」14に一個(1)加えた数の重複しないサブパターン断片にスライスする場合、重複しないサブパターン断片のうちの少なくとも1つが正確に合致するはずである。図1の例に示されるように、データ文字列 $T_{i \dots j}$  16が、重複しないサブパターンの $k+1$ または三個(3)の断片に分割されている。それゆえに、この三個(3)の断片は、図番18で示される「abra」、図番20で示される「cada」、および図番22で示される「bra」である。この例では、文字「br」が置換され、文字「b」が削除されている2つの誤りと共に、図番20で示される「cada」が正確に合致している。

40

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0006】

可能性のある合致パターンの大きなセットを用いた、大量の入力データを利用するパタ

50

ーン合致に関する高速で費用効率が高い機構に対する著しい必要性が存在する。

【課題を解決するための手段】

【 0 0 0 7 】

本発明の一態様では、少なくとも1つのサーチエンジンを用いて、各々のパターンが所定の許容可能な誤りを有する1つまたは複数のパターンと合致するデータセグメントについてデータストリームを検査する方法が開示されている。この方法は、各々のパターンが関連する許容可能な誤りを有する1つまたは複数のパターンを検出するように各々の並列フィルタ機構が構成されている複数の並列フィルタ機構を用いて、複数のパターンのシンボルの組み合わせに関してデータストリームをフィルタリングすることと、複数の並列フィルタ機構を用いて複数の可能性のあるパターン断片の合致を検出することと、複数の並列フィルタ機構から、各々の合致パターンが関連する許容可能な誤りを有する、複数の可能性のある合致パターンを特定することと、削減ステージを用いて、特定された複数の可能性のある合致パターンを、各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンのセットへと削減することと、関連するデータ、および各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンの削減されたセットを検証ステージに提供することと、関連するデータおよび可能性のある合致パターンの削減されたセットを利用する近似合致エンジンを含む検証ステージを用いて、複数のパターンのシンボルの組み合わせおよび関連する許容可能な誤りからデータストリームでのパターン合致の存在を検証することを含む。

10

【 0 0 0 8 】

20

本発明の別の態様では、少なくとも1つのサーチエンジンを用いて、各々のパターンが所定の許容可能な誤りを有する1つまたは複数のパターンと合致するデータセグメントについてデータストリームを検査する方法が開示されている。この方法は、並列ブルーム ( Bloom ) フィルタのセット、並列ブルームフィルタ配列のセットまたは単一のハッシュジェネレータを利用する並列ブルームフィルタ配列のセットから成るグループであり、各々のパターンが関連する許容可能な誤りを有する1つまたは複数のパターンを検出するように各々の並列フィルタ機構が構成されている、複数の並列フィルタ機構を用いて、複数のパターンのシンボルの組み合わせに関してデータストリームをフィルタリングすることと、複数の並列フィルタ機構を用いて複数の可能性のあるパターン断片の合致を検出することと、複数の並列フィルタ機構から、各々の合致パターンが関連する許容可能な誤りを有する、複数の可能性のある合致パターンを特定することと、特定された複数の可能性のある合致パターンを、各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンのセットへと削減することと、関連するデータ、および各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンの削減されたセットを検証ステージに提供することと、関連するデータおよび可能性のある合致パターンの削減されたセットを利用する近似合致エンジンを含む検証ステージを用いて、複数のパターンのシンボルの組み合わせおよび関連する許容可能な誤りからデータストリームでのパターン合致の存在を検証することを含む。

30

【 0 0 0 9 】

本発明のさらに別の態様では、少なくとも1つのサーチエンジンを用いて、各々のパターンが所定の許容可能な誤りを有する1つまたは複数のパターンと合致するデータセグメントについてデータストリームを検査する方法およびシステムが開示されている。この方法は、少なくとも1つの検索エンジンを用いて、誤検出誤りを有する1つまたは複数のパターン断片と合致するデータセグメントについてデータストリームを検査するために、単一のハッシュ値から複数のハッシュ値を抽出するために単一のハッシュジェネレータを利用することと、複数の並列ブルームフィルタ配列を用いて複数のハッシュ値を利用することを含む。

40

【 0 0 1 0 】

本発明のさらに別の態様では、少なくとも1つのサーチエンジンを用いて、各々のパターンが所定の許容可能な誤りを有する1つまたは複数のパターンと合致するデータセグメ

50

ントについてデータストリームを検査するシステムが開示されている。このシステムは、各々のパターンが関連する許容可能な誤りを有する1つまたは複数のパターンを検出するように各々の並列フィルタ機構が構成されている複数の並列フィルタ機構を利用し、複数のパターンのシンボルの組み合わせに関してデータストリームをフィルタリングし、複数の可能性のあるパターン断片の合致を検出し、各々の合致パターンが関連する許容可能な誤りを有する、複数の可能性のある合致パターンを特定するフィルタステージと、特定された複数の可能性のある合致パターンを、各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンのセットへと削減する削減ステージと、複数のパターンのシンボルの組み合わせからデータストリームでのパターン合致の存在を検証するために、関連するデータおよび可能性のある合致パターンの削減されたセットおよび関連する許容可能な誤りを受信および利用する、近似合致エンジンを含む検証ステージとを含む。

10

**【0011】**

本発明のさらに別の態様では、少なくとも1つのサーチエンジンを用いて、各々のパターンが所定の許容可能な誤りを有する1つまたは複数のパターンと合致するデータセグメントについてデータストリームを検査するシステムが開示されている。このシステムは、並列ブルームフィルタのセット、並列ブルームフィルタ配列のセットまたは単一のハッシュジェネレータを利用する並列ブルームフィルタ配列のセットから成るグループであり、各々のパターンが関連する許容可能な誤りを有する1つまたは複数のパターンを検出するように各々の並列フィルタ機構が構成され、複数のパターンのシンボルの組み合わせに関してデータストリームをフィルタリングし、複数の可能性のあるパターン断片の合致を検出し、各々の合致パターンが関連する許容可能な誤りを有する、複数の可能性のある合致パターンを特定する複数の並列フィルタ機構と、特定された複数の可能性のある合致パターンを、各々の合致パターンが関連する許容可能な誤りを有する、可能性のある合致パターンのセットへと削減する削減ステージと、複数のパターンのシンボルの組み合わせからデータストリームでのパターン合致の存在を検証するために、関連するデータおよび可能性のある合致パターンの削減されたセットおよび関連する許容可能な誤りを受信および利用する、近似合致エンジンを含む検証ステージとを含む。

20

**【0012】**

例示的であり、限定ではないが、本発明の可能性のあるアプリケーションの例は、コンピュータ通信ネットワークに関する侵入検知システム(IIDS)、コンピュータ計算による生物学および遺伝学、構造化および非構造化テキストに対するテキスト検索、および光学文字走査(OCR)からのテキスト検索を含む。

30

**【0013】**

本発明の付加的な態様は、限定ではないが、各パターンが、例えば何万パターン以上の多数のパターン断片を含み得るその許容可能な誤りを特定し得る、複数のパターンを用いた近似合致に関するフィルタリング技術と、並列合致動作を実行するために、および多種多様の(パターン長、許容可能な誤り)組み合わせをサポートするために、1つが各パターン断片長のためのものである、正確な合致エンジンの並列セットを利用することと、各パターンに特定の数の誤りを許容することと、並列ハードウェア検索実装に対する従順さおよびそのような実装が高速の検索結果を提供し得ることと、テキストの範囲に関して可能性のある合致パターンの数を限定することによって検証ステージを単純化し、検索エンジンが付加的な可能性のある検索結果をより短い周期で処理することを可能にして、全システムが非常に高速で動作しながら能力を向上させることと、各正確な合致エンジンに対してブルームフィルタ配列を利用することと、1つだけのハッシュ関数ジェネレータを使用することによって各ブルームフィルタ配列を効果的に実装することとを含む。

40

**【0014】**

本発明のさらに別の態様は、検証ステージでの検索の範囲が可能性のある合致パターンのより小さいセットへと削減される削減ステージである。これらの技術は、各パターンおよびその許容可能な誤りが一度だけ記憶されることを可能にする、断片およびパターンの

50



間の間接層を使用する。データ構造を単純化し、それらをハードウェア実装に従順にする第1の例示的な技術が存在する。この技術は、可能性のある合致断片に関する断片識別子を取り出すために、ピンインデックスを使用する参照を含む。第2の参照は、断片を含むパターンに関するパターン識別子を取り出すために断片識別子を使用する。第3の参照は、検証エンジンによって考慮されるパターンおよび関連する許容可能な誤りの対を取り出すためにパターン識別子を使用する。断片を含むパターンに関するパターン識別子を解決するために、正確な合致エンジンでの合致を生じたテキスト断片を使用する第2の例示的であるが、限定ではない技術が存在する。パターン識別子は、検索エンジンによって考慮されるパターンおよび関連する許容可能な誤りの対を取り出すために使用される。

【0015】

10

これらは、本発明の無数の態様のうちのただ一部であり、本発明に関連付けられている無数の態様の全てを含むリストとみなされるべきではない。

【0016】

本発明をより理解するために、添付の図面を参照されたい。

【発明を実施するための最良の形態】

【0017】

以下の詳細な説明では、本発明の完全な理解をもたらすために多数の特定の詳細が示されている。しかしながら、当業者は、本発明がこれらの特定の詳細なしで実施され得ることを理解するであろう。別の例では、よく知られている方法、手順および構成要素は、本発明を分かりにくくさせるので詳細には説明されていない。

20

【0018】

本発明は、パターンの大きなセットを用いた近似パターン合致のためのスケーラブルなフィルタ回路である。フィルタ回路は、記憶されているパターンのセットの間で可能性のある合致を検証し、各パターンが許容可能な誤りの数および文字の入力ストリームを特定する。許容可能な誤りの数は、所定のものであるか、または単一文字の付加、削除および置換の数を数える一般的な編集距離尺度を使用して特定される。可能性のある合致が検出されると、合致パターンまたは複数の可能性のある合致パターン同様に、入力データストリームでの単数または複数の位置が特定される。本発明は、可能性のある合致パターンの総数から、前に特定された可能性のある合致パターン（単数または複数）を利用し、入力データストリームのデータセグメント（単数または複数）での近似合致を検索する検証ス

30

【0019】

単一パターンに関する検索の方法論は、以下のように規定され得る：「k」の許容可能な誤りでテキスト「T」でのパターン「P」のインスタンスを特定する、ここで「k」は最大編集距離である。編集距離は、全ての誤りが、必ずではないが通常は、同じ重み付けを共有する、単一文字の挿入、削除および置換の数として規定されている。一般に、転置等の別の種類の誤りが距離尺度に含まれてもよく、各種類の誤りは固有の重み付けが割り当てられ得る。

【0020】

「m」がパターンの文字としての大きさであり、「n」がテキストの文字としての大きさであると仮定すると、誤りレベルは、特定のパターンに対して、許容可能誤りの数とパターンの大きさの比率、即ち「 $\alpha = k / m$ 」として規定され得る。誤りレベルならびにパターンおよびテキストが構成されるアルファベットの大きさ「 $\sigma$ 」は、合致が見つかる確率に影響を及ぼす。無作為に構成されているテキストおよび無作為に構成されているパターンを仮定しての合致確率の式「 $f(m, k)$ 」は、以下の式：

40

【数2】

$$f(m, k) = \left( \frac{e^2}{\sigma(1 - \alpha)^2} \right)^{m(1 - \alpha)}$$

50

で提供される。ここで「e」は自然対数の底である。

【0021】

フィルタリングアルゴリズムは種々の近似パターン合致問題に対してより良い成績を達成すると考えられている。一般的な手法は、可能性のある合致を特定するためにテキストの小セクションに対し単純な検索を実行するものである。可能性のある合致が見つかり、テキストのその範囲は、それが実際に特定パターンに対する近似合致であるか否かを確かめるために検査される。一般に、検証は、あらゆる近似パターン合致アルゴリズムで実行され、多かれ少なかれフィルタリング動作に結び付いている可能性がある。

【0022】

本発明のシステムの概略的な模式図が図2に示されており、図番30によって示されている。高いデータ転送速度でフィルタ回路に供給可能な種々の他のデータソースと同様に、通信リンク、ディスク、レイド(redundant array of independent disks, RAID)またはストレージエリアネットワーク(SAN)のいずれかを介して大量のデータが入力として提供され得る。このデータ入力は図番32で示され、図番34で示されるネットワーク入力を介して同様に提供され得る。高速データは、図番36で示される高帯域相互接続を介して高速で提供され得る。この高速データは、可能性のある合致に対して入力データを入力パターンのセットに関してスキャンするフィルタ回路38を次に通される。次に、フィルタ回路38で合致を取り出すデータセグメントを処理する際に、検証ステージ42によって考慮される必要がある可能性のある合致パターンのセットを削減する削減ステージ40がフィルタ回路38と検証ステージ42との間に存在する。検証ステージ42は、入力パターン46のセットに対し、合致が存在するか否かを検証するための全ての適切な動作を実行する。次に、検索結果が、図番44で示されるように提供される。所定の入力パターン46がフィルタ回路38、削減ステージ40および検証ステージ42に提供される。

【0023】

テキストの特定の窓では、正確な合致を検索することが可能であり、パターン断片のいずれもが所定数「r」のパターンによって特定される。特定のパターン「i」が「k<sub>i</sub>」個の誤りを許容する場合、パターン断片の総数は式：

【数3】

$$p = \sum_{i=1}^r (k_i + 1)$$

によって示される。

【0024】

このフィルタリング手法は、ブルームフィルタの並列セット、並列ブルームフィルタ配列のセット、または単一ハッシュ関数ジェネレータを利用するブルームフィルタ配列のセットを含み得る、並列フィルタ(parallel filter)機構と共に利用される。本発明の基本的なハードウェア実装の概略図である図3で示されるように、本発明が図番50で示されており、図番54で示される多数のブルームフィルタを含む。従来のブルームフィルタ54では、「b」個のハッシュ値を使用してセットに要素が挿入され、要素がキーとして利用され、そこで各ハッシュ値がB-ビットベクトルでのビット位置を特定する。bビット位置の各々でのビットは1に設定されるのが好ましい。ビットが既に1に設定されている場合、変更はされない。特定の要素が、ブルームフィルタ54によって示されるセットの構成要素であるか否かを検査するために、要素および同じbハッシュ関数がbハッシュ値を計算するのに利用される。ベクトルの全てのbビットが1に設定されている場合、その要素がセットの構成要素であることが宣言される。

【0025】

ブルームフィルタ54は、検出漏れを生じない。要素がセットの構成要素である場合、要素がセットに挿入される際にB-ビットベクトルのbビット位置が1に設定される。付

10

20

30

40

50

加的な要素のセットへの挿入は、ベクトル内のいずれのビットをもリセットしない。しかしながら、ブルームフィルタ 54 は、所定の確率で検出漏れを生じる。この確率は、式：  
【数 4】

$$f = \left(1 - e^{-\frac{pb}{B}}\right)^b$$

で計算され得る。以下の関係が成り立つ場合：

【数 5】

$$b = \frac{B}{p} \ln 2 \text{ の場合 : } f = (1/2)^b$$

10

【0026】

複数のパターンでの近似合致フィルタリングと、各パターンがその許容可能な誤りを特定することとは、種々の長さのパターン断片のセットを生じる。好ましくは、必ずではないが、ブルームフィルタ 54 は、可能性のある各パターン断片長に対して 1 つのブルームフィルタ回路 54 の固定長要素を記憶している。それゆえに、可能性のあるパターン断片長の範囲は、ある範囲内に制限される。

【0027】

$l_{min}$  が最小パターン断片長である場合、 $l_{min}$  は、パターン m の大きさを最大編集距離 k に 1 を加えたもので割った値より小さいか等しい：

20

【数 6】

$$l_{min} \leq \left\lfloor \frac{m}{k+1} \right\rfloor$$

【0028】

$l_{max}$  が最大断片長である場合、 $l_{max}$  は、パターン m の大きさを最大編集距離 k に 1 を加えたもので割った値より大きい等しい：

【数 7】

$$\left\lceil \frac{m}{k+1} \right\rceil \leq l_{max}$$

30

【0029】

ブルームフィルタ 54 の各ブルームフィルタがパターン断片長に対応する際に必要であるブルームフィルタ 54 の総数は：

$$l_{max} - l_{min} + 1$$

である。

【0030】

好ましい手法は、図 3 の図番 50 で概略的に示されているように、ブルームフィルタ 54 の各々に並行して問い合わせることである。種々のストライドのテキスト窓が各ブルームフィルタ 54 への入力キーとして選択可能であるように、ブルームフィルタ 54 の各々は、パターン断片長に対応している。

40

【0031】

ブルームフィルタ 54 のいずれかが検出の合致 56 を生じると、データのセグメントまたはテキスト窓 58、入力ストリームデータのセグメントの位置 59 および付加的な合致メタデータ 60 が削減ステージ 40 に送信される。削減ステージ 40 で利用される技術は、可能性のある合致パターンを著しく、例えば、10000 超を 10 未満へと、削減することが可能である。

【0032】

50

削減ステージ40から渡された結果は、近似合致検索エンジンを含む検証ステージ42に進む。検証ステージ42によって考慮される候補パターンの数を削減することによって、検証ステージが所与の時間でより多くの可能性のある検索結果を処理することが可能になり、従って、高速で動作しながら全システムの能力面を向上させることが可能になる。

【0033】

通常、ブルームフィルタ配列54は、セットメンバーシップ問い合わせに必要な、例えば、ランダムアクセスメモリ(RAM)等のメモリアクセスの数を最小化する。さらに、ブルームフィルタ配列54は、 $B$ -ビットベクトルを大きさ「 $q = B / W$ 」の「 $W$ 」個のベクトルに分割し、ここで「 $q$ 」はメモリのワードサイズである。好ましくは、プリフィルタハッシュ関数を使用して、記憶されている要素が「 $W$ 」個のベクトル(メモリワード)上に均等に分配されていることが好ましい可能性がある。これは、「 $W$ 」個の「 $q$ 」-ビットブルームフィルタの配列を生じる。

【0034】

プログラミングの際に「 $q$ 」-ビットブルームフィルタ54のビットが設定され、次に、問い合わせの際に「 $b$ 」個のハッシュ関数を使用して検査される。ブルームフィルタ配列54への問い合わせは、1個のメモリが「 $q$ 」-ビットベクトルを読んで取ってくることを必要とする。図4で示されて図番70で概略的に示されるように、レジスタ80およびビット選択回路82を使用することによって、「 $b$ 」個のハッシュ関数によって特定されるビット位置の検査がオンチップでパイプライン方式で実行され得る。

【0035】

このアプリケーションでは、キー(即ち、パターン断片)が図番72によって示されている。このキーは、「 $w$ 」個のベクトルのリストの特定の「 $q$ 」-ビットベクトルを特定するために、プレフィルタハッシュ関数73によって使用される。「 $w$ 」個のベクトルのリストの特定の「 $q$ 」-ビットベクトルは、例えば、RAM等のメモリの列74で示されており、そこで具体的で例示的なベクトルが図番76で特定されている。これらの問い合わせは、レジスタ80内でのビット位置を特定する、図番78で示される一連の「 $b$ 」個のハッシュ関数で検査され、次に、合致検出機能82に提供される。「 $b$ 」個のビット位置の全てが1に設定されている場合、このキーは可能性のある合致パターンに関するパターン断片である。

【0036】

図5に示されるように、ブルームフィルタ配列54を実装するために必要なロジックの量が最小化され得るのが好ましい。このロジックは、図番90で概略的に示されている。図番92で示される単一のハッシュ関数が存在する。単一の乱数値の生成が存在する。この乱数値からのビットのサブセットがプレフィルタハッシュアドレスおよび図番94で示される「 $b$ 」個のフィルタビット位置を構成するために利用される。「 $w$ 」個のベクトルのリストの特定の「 $q$ 」-ビットベクトルが、例えば、RAM等のメモリの列96で示されており、そこで具体的で例示的なベクトルが図番98で特定されている。この特定のベクトル98は、レジスタ100に渡され、次に、合致パターン検出機能102に渡される。このビット選択値94は、サイズが少なくとも $\log_2(W) + (b * \log_2(q))$ ビットである必要がある。ハッシュ関数92の $H_3$ クラスは、このアプリケーションに十分に広範な値を提供し得る例示的で限定ではないハッシュ関数例である。

【0037】

利用可能な再構成可能なハードウェアの例示的で限定ではない例は、Xilinx(登録商標)VirtexII(登録商標)4000シリーズFPGAを含むFPGAs、即ち、フィールドプログラマブルゲートアレイを含む。Xilinx社は、95124-3400、カリフォルニア州、サンノゼ、ロジックドライブ2100に事業所を有するデラウェア州の会社である。VirtexII(登録商標)のデバイスシリーズの内蔵メモリの例示的で限定ではない例は、120個の18キロバイトブロックのランダムアクセスメモリ(ブロックRAMs)を含む。これらのブロックRAMは、例示的で限定ではないが、最大ワード長36ビット $\times$ 512ワードを用いて種々の大きさのワードに構成され得る

10

20

30

40

50

。

## 【 0 0 3 8 】

誤検出の確率に関する前式を利用し、一定のハッシュ性能を仮定すると、18キロバイトブロックRAMを用いて実装されているブルームフィルタ配列54は、bビット位置の数が4と等しい場合、3,194個の要素のセットの誤りの確率を0.063と示し得る。bビット位置の数が3と等しい場合、能力は4,259個の要素に増加するが、誤検出の可能性が0.125に増大する。

## 【 0 0 3 9 】

前述したように、1つのブルームフィルタ配列54が各固有のパターン断片長に対して必要である。さらに、データ入力速度の歩調をとる必要がある並列回路の数の検討材料も存在する。例示的で限定ではない例では、周期(64-ビットインターフェース)当たり8個の新規のASCII文字を受け入れるシステムは、並列に動作する回路の8個のインスタンスを必要とする。V i r t e x I I 4 0 0 0 (登録商標)の場合、各回路インスタンスに対して最大限で15個のブロックRAMが利用可能である。インターフェースバッファに対してブロックRAMリソースを利用可能にするために、これは、ブロックRAMをより小さい数、例えば、14個のブロックRAMに制限することをもたらす。この例示的で限定ではないリソース割り当ては、パターン断片の長さ、パターン断片および許容可能な誤りの組み合わせとに制約を生じる可能性がある。これは、次に、最大誤りレベルに制限を生じる。「m」がパターンの大きさと等しく、「p」がパターン断片の数に等しく、「k」が許容可能な誤りの数または編集距離である場合、以下の式：

$$m_{min} = p_{min} (k + 1)$$

$$m_{max} = p_{max} (k + 1)$$

## 【 数 8 】

$$\alpha_{max} = \frac{k}{m_{min}} = \frac{k}{p_{min}(k+1)}$$

が適用可能であり、ここで  $\alpha$  = 誤りの数をパターンの大きさに除算した比率である。

## 【 0 0 4 0 】

以下のテーブル1は、「 $\alpha$ 」、即ち、誤りの数をパターンの大きさに除算した比率が1以下である場合のものである。

【 表 1 】

テーブル1		
k	$m_{min}$	$m_{max}$
0	1	14
1	2	28
2	3	42
3	4	56
...	...	...

## 【 0 0 4 1 】

以下のテーブル2は、「 $\alpha$ 」、即ち、誤りの数をパターンの大きさに除算した比率が1/2以下である場合のものである。

10

20

30

40

【表 2】

テーブル2		
k	$m_{\min}$	$m_{\max}$
0	2	15
1	4	30
2	6	45
3	8	60
...	...	...

【 0 0 4 2 】

以下のテーブル 3 は、「 $k$ 」、即ち、誤りの数をパターンの大きさに除算した比率が  $1/3$  以下である場合のものである。

【表 3】

テーブル3		
k	$m_{\min}$	$m_{\max}$
0	3	16
1	6	32
2	9	48
3	12	64
...	...	...

【 0 0 4 3 】

それゆえに、第 2 の例およびテーブル 2 で、「 $k$ 」、即ち、誤りの数をパターンの大きさに除算した比率が  $1/2$  以下である場合で、パターンが誤りを許容しない、即ち、「 $k = 0$ 」である場合、少なくとも 2 個の文字が存在し、15 個以下の文字であるはずである。1 個の誤りを許容するパターン、即ち、「 $k = 1$ 」は、少なくとも 4 個の文字で、30 個以下の文字を含むはずである。広範な種類の許容パターンの大きさおよび許容可能な誤りが利用され得るが、少なくとも 2 個の文字が存在して 15 個以下の文字であるパターンが、英語での大部分のテキスト検索に対して有効な制約であると考えられており、しかしながら、これが限定として解釈されるべきではない。

【 0 0 4 4 】

断片が許容可能な長さの範囲に様に分布していると仮定することによって、大まかな能力の推定が展開され得る。例示的で限定ではない例では、各ブルームフィルタ配列 54 が約 3000 個のパターン断片の能力を有する場合、システムは 42000 個のパターン断片の総合能力を有する。各パターンが 3 個のパターン断片に分割され得ると仮定される場合、システムは 14000 個のパターンの能力を有する。

【 0 0 4 5 】

図 3 に示されるように、一度、1 つまたは複数のパターン断片長 56 に関して可能性のある合致が検出されると、「 $r$ 」個のパターンのうちの 1 つに関する近似合致が存在するか否かを判定するために検証ステージ 42 によってテキストのその範囲が検査される必要がある。「 $r$ 」はおよそ 10000 パターン以上であり得るので、検証ステージ 42 が実行する必要のある検索範囲を削減する必要性が存在する。これこそ、検証ステージ 42 が考慮する、可能性のある合致パターン（パターンセット）のセットを低減する重要な役割を削減ステージ 40 が提供することである。パラメータが制約の範囲内にある限り、許容可能な誤りの数が各パターンに対して特定され得るという仮定が存在する。

【 0 0 4 6 】

図 3 に示されるようなこの特許出願全体で、図 3 に示されるようにフィルタステージ 38、削減ステージ 40 および / または検証ステージ 42 は、少なくとも 1 つの再構成可能なロジックデバイス、例えば、フィールドプログラマブルゲートアレイ（「FPGA」）または少なくとも 1 つの集積回路、例えば、特定用途向け集積回路（「ASIC」）を使

10

20

30

40

50

用し得る。

【 0 0 4 7 】

削減ステージ 4 0 を実行する 2 つの例示的だが限定ではない手法が存在する。第 1 の手法は、図 6 の図番 1 2 0 で示されるパターンセットを解決するために利用されるデータ検索を単純化することである。これは、データ文字列がシフトレジスタ 5 2 に入ることを可能にし、シフトレジスタ 5 2 は次にブルームフィルタ配列 5 4 を含むことが好ましいフィルタステージまたは回路 3 8 に渡す。その目的はこの手法を用いて、フィルタステージまたは回路 3 8 のブルームフィルタ配列 5 4 によって計算されたハッシュ値のうちの、全てではない場合には一部を、テーブル 1 2 8 へのインデックス、例えば、ピンインデックス ( Bin Index ) 1 2 4 として利用することである。

10

【 0 0 4 8 】

例えば、パターン断片 1 2 1 がフィルタステージまたは回路 3 8 によって受信可能であり、ピンインデックス 1 2 4 を含むハッシュ値 1 2 6 として受信される。この第 1 のテーブル 1 2 8 へのエントリは、パターン断片に対する識別子 1 2 7、例えば、ピンインデックス 1 2 4 を含むハッシュ値 1 2 6 をマッピングする断片 ID を含む。例えば、例示的な識別子、例えば、例示的なパターン断片 1 2 1 と関連付けられている断片 ID<sub>1</sub> および断片 ID<sub>4</sub> が存在する。パターン断片に対する識別子 1 2 7、例えば、断片 ID は、各パターン断片に割り当てられている固有の 2 進数タグである。

【 0 0 4 9 】

第 2 のテーブル 1 3 2 は、1 つまたは複数のパターン識別子を可能性のある合致パターンのセットに関してインデックス付けするために、パターン断片に対する識別子 1 2 8、例えば、断片 ID を利用する。1 つまたは複数のパターンが特定のパターン断片を特定し得るので、第 2 のテーブル 1 3 2 のエントリは、1 つまたは複数のパターン識別子、例えば、P I D s を含む。パターン識別子、例えば、P I D s は、各パターンに関連付けられている固有の 2 進数タグである。例えば、例示的な識別子 1 3 1 は、2 つのパターン 1 3 7 および 1 3 8 に関連する。

20

【 0 0 5 0 】

第 3 のテーブル 1 3 4 は、1 つまたは複数の ( パターン、許容可能な誤り ) 対をインデックス付けするためにパターン識別子、例えば、P I D を利用する。次に、各々の合致パターンが所定の許容可能な誤りを有する、識別されたセットまたは複数の可能性のある合致パターンが、例えば、( パターン、許容可能な誤り ) 対 1 3 7 および 1 3 8 が、図番 1 3 6 で示されるように生成される。

30

【 0 0 5 1 】

可能性のある合致 1 3 6 のこのパターンセットが検証ステージ 4 2 に渡され、それは合致パターンおよび関連する所定の許容可能な誤りに関する近似合致エンジン 1 4 2 による評価を含む。周期当たりの参照数がメモリによってサポートされる参照数を超過しない限りは、そこでは、パターン断片識別子 1 2 8、例えば、断片 ID s、パターン識別子 1 3 2、例えば、P I D s およびパターン 1 3 4 に関するテーブルの 1 つのコピーだけを必要とする。

【 0 0 5 2 】

40

削減ステージ 4 0 に関する第 2 および好適な方法論が存在し、図 7 の図番 1 5 0 で概略的に示されている。これは、データ文字列が、好ましくはブルームフィルタ配列 5 4 を含むフィルタステージまたは回路 3 8 に渡されることを可能にする。この手法を用いると、ブルームフィルタ配列 5 4 での合致を生じる実際のデータセグメントまたは断片が、それらのパターン断片を特定するパターン 1 6 2 に対してパターン識別子を解決するために利用される。合致を生じるデータセグメントまたは断片が、図番 1 5 6、1 5 8 および 1 6 0 で示されるデータ構造の 1 つまたは複数のエントリを特定するために使用され、データ構造は、関連するパターン断片を特定するパターンに関するパターン識別子、例えば、P I D s を記憶している。適切なデータ構造の例示的であるが、限定ではない例は、ハッシュテーブルおよび平衡探索木である。この方法論は、1 つまたは複数のデータ構造を含み

50

得る。図 7 の例示的であるが、限定ではない例 1 5 0 では、1 つのデータ構造が各パターン断片長に対して割り当てられている。

【 0 0 5 3 】

例示的であるが、限定ではない例では、ブルームフィルタ配列 5 4 での合致を生じる 2 つのデータセグメントまたは断片が図番 1 2 1 および 1 2 3 で特定される。データ断片 1 2 1 は、データ構造のエントリ 1 5 7 を削減ステージ 4 0 の一部として特定する。これらのデータ構造 1 5 6、1 5 8 および 1 6 0 は、多種多様の異なる構造、例えば、決定木、ハッシュテーブル等などを含み得る。これらの参照（単数または複数）の結果は、パターンセットのパターンに対するパターン識別子のセットである。従来技術の手法と同様に、これらのパターン識別子、例えば、P I D s、1 5 6、1 5 8 および 1 6 0 は、検証ステップの前にテーブル 4 0 からパターンおよび関連する許容可能な誤りを取り出すために利用される。前に参照したエントリ 1 5 7 は、可能性のある合致 1 6 4 のセットを生じさせるために、パターン 1 6 3 および 1 6 5 と関連する許容可能な誤りとを特定する。可能性のある合致および関連する所定の許容可能な誤り 1 6 4 のセットは、次に、検証ステージ 4 2 の近似合致エンジン 1 6 6 によって評価される。

【 0 0 5 4 】

この手法は、パターン識別子のセットを、2 つのステップの代わりに単一のステップで解決し、さらに、ブルームフィルタ配列 5 4 によって生じる誤検出誤りを除去する。さらに、パターン識別子構造（単数または複数）のエントリの位置を決定するために実際のデータセグメントが利用されているため、明示的な合致が実行される。テーブルにエントリが存在しない場合、誤検出が検出され、検出ステージ 4 2 にその特定のパターン断片長に対するパターン識別子、例えば、P I D s は渡されない。その代償は、データ構造がより複雑である可能性があること、および実装によってはその実装がよりリソース消費型になることである。

【 0 0 5 5 】

それゆえに、これは、ハードウェア実装の影響を受け易い複数のパターンに対する近似パターン合致に関するフィルタ回路 3 8 および削減ステージ 4 0 に関するスケーラブルな設計である。何千ものパターンに加え、複数のフィルタ回路が、周期当たりの複数の入力シンボルをサポートすることが可能である。高性能フィルタ回路 3 8 および削減ステージ 4 0 を利用することによって、検証ステージ 4 0 に配置された性能要求が解析され得る。この解析を目的として、出願人は、全てのパターンが同じ数「k」の許容可能な誤りを特定すると仮定する。検証ステージでの有効負荷は、合致の確率および可能性のある合致パターンのセットの期待される大きさとして決定される。合致の確率は、単純に、「r ( k + 1 )」個の断片のうちのいずれかがテキスト窓で合致を生じる確率と、ブルームフィルタ配列 5 4 の誤検出確率との総和であり、ここで「r」は「r」個のパターン断片の所定の数であり、「k」は所定の誤りの数である。

【 0 0 5 6 】

「L」がブルームフィルタ配列 5 4 の数であり、「」個の文字のアルファベットに対して無作為のテキストを仮定する場合、これらの断片のいずれかの合致の確率は、

【数 9】

$$E[\text{合致}] = \frac{r(k+1)}{\sigma \left\lfloor \frac{m}{k+1} \right\rfloor}$$

になる。

【 0 0 5 7 】

ブルームフィルタ配列の誤検出確率の付加は、

10

20

30

40



【数 1 0】

$$E[\text{合致}] = \frac{r(k+1)}{\sigma^{\lfloor \frac{m}{k+1} \rfloor}} + Lf$$

を提供する。

【0 0 5 8】

例示的だが、限定ではない例示値  $L = 14$ 、 $m = 40$ 、 $r = 14000$  および  $f = 0.0034$  を利用することによって、合致確率は最小断片長に非常に敏感である。例えば、 $m = 5$  および  $k = 1$  (2 個の文字の最小断片長) である場合、合致確率が 1 である。パターンの大きさが 6 (3 個の文字の最小パターン断片長) に増大されると、合致確率は 0.079 である。この結果は、最小パターン断片の大きさが 3 個以上の文字であるべきであることを示唆している。この状況では、12 周期毎に 1 つの合致が期待される。

10

【0 0 5 9】

図 6 に示される削減ステージ 40 の方法論を、一連のインデックス参照と共に利用することによって、もし少なくとも 1 つのブルームフィルタ配列 54 が合致を生じる場合、合致を生じるブルームフィルタ配列 54 の期待数、即ち、ピンインデックス参照の期待数は：

【数 1 1】

$$E[\text{ピン}] \leq 1 + \frac{L}{\sigma^{\lfloor \frac{m}{k+1} \rfloor}} + Lf$$

20

である。

【0 0 6 0】

パターン断片が「L」個のブルームフィルタ配列 54 (パターン断片長) 上に一様に分布し、ピンの上に一様に分布しているとの仮定の下で、図 6 に示されるように、ピン当たりのパターン断片識別子 128 の期待数は：

【数 1 2】

$$E[\text{断片ID/ピン}] \leq \frac{r(k+1)}{LW \left( \frac{B}{W} \right)^b}$$

30

である。ここで「B」は、ブルームフィルタ配列 54 を実装するために使用されているメモリの大きさであり、「W」は、ブルームフィルタ配列 54 のワードの数であり、「b」は、ブルームフィルタ配列 54 の各ブルームフィルタで使用されているハッシュ関数の数である。

【0 0 6 1】

最後に、所与のパターン断片 131 を特定するパターンの期待数は、図 6 に示されるように：

40

【数 1 3】

$$E[\text{パターン/断片ID}] \leq 1 + \frac{r}{\sigma^{\lfloor \frac{m}{k+1} \rfloor}}$$

を含む。

【0 0 6 2】

それゆえに、少なくとも 1 つのブルームフィルタ配列 54 が合致するという条件で、期待されるパターンの大きさは：

50

【数 1 4】

$$E[\text{パターン}] \leq \left( 1 + \frac{L}{\sigma^{\lfloor \frac{m}{k+1} \rfloor}} + Lf \right) \left( \frac{r(k+1)}{LW \left( \frac{B}{W} \right)^b} \right) \left( 1 + \frac{r}{\sigma^{\lfloor \frac{m}{k+1} \rfloor}} \right)$$

である。

【0063】

一様なテキスト、断片長の一様な分布、一様な分布（優良なハッシュ関数）と、 $L = 14$ 、 $m = 40$ 、 $r = 14000$ 、 $b = 3$  および  $W = 512$  とを仮定すると、 $m$  および  $k$ （パターン長および許容可能な誤り）の実際的な値に対して期待されるパターンセットの大きさは10未満である。アルファベットの大きさおよび/またはパターンの大きさが増大するに従って、期待されるパターンセットの大きさは、急速に1に近づく。

10

【0064】

図7に示される他の削減ステージ手法を、テキストセグメントをインデックスとして用いて使用する場合の期待されるパターンセットの大きさの式は同様であり、若干小さいパターンセットサイズを生じる。12の周期毎に1つのパターンが合致する前の結果と組み合わせ、検証ステージ42の平均スループットに対する控えめの制約は、毎周期当たり1つのパターンで約1個の合致である。

20

【0065】

従って、新規の発明の複数の実施形態が示され、説明されてきた。前述の説明から明らかであるように、本発明の所定の態様は、本明細書で例示されている例の特定の詳細によって制限されず、それゆえに、当業者が他の修正および適用またはそれらの同等物を創案することが意図されている。前述の明細書で使用されている用語「有する」、「有している」、「含む」および「含んでいる」ならびに類似の用語は、「必要とする」の意味ではなく、「任意選択で」または「含んでよい」の意味として使用されている。一方、本構成の多くの変更、修正、変形および他の使用および用途が、明細書および添付の図面を熟考した後に当業者に明白になるであろう。本発明の精神および範疇から逸脱しないそのような変更、修正、変形および他の使用ならびに用途の全ては、添付の特許請求の範囲によってのみ制限される本発明によって網羅されているとみなされる。

30

【図面の簡単な説明】

【0066】

【図1】従来技術の近似パターン合致技術の例示的な概略図である。

【図2】データソース、フィルタ回路、削減ステージおよび検証ステージを含む本発明の例示的なブロック図である。

【図3】合致検出機能、削減ステージならびに近似パターン合致を有する検証ステージを備えるブルームフィルタを含む、本発明の例示的であり、限定ではないブロック図である。

【図4】ブルームフィルタ配列を使用する第1のフィルタリングステージ処理技術の例示的なブロック図である。

40

【図5】単一のハッシュ関数ジェネレータと共にブルームフィルタを使用する第2のフィルタリングステージ処理技術の例示的なブロック図である。

【図6】第1の削減ステージ処理技術の例示的なブロック図である。

【図7】第2の削減ステージ処理技術の例示的なブロック図である。

【図 1】

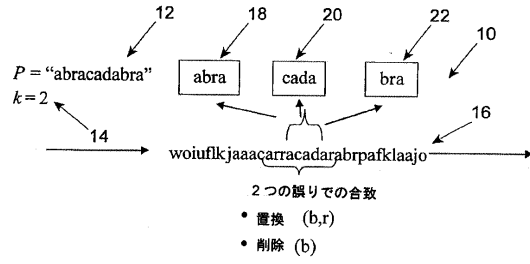


FIG. 1

【図 2】

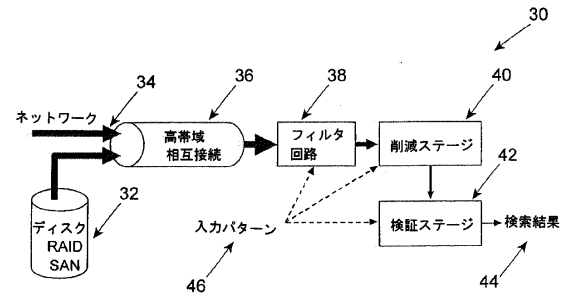


FIG. 2

【図 3】

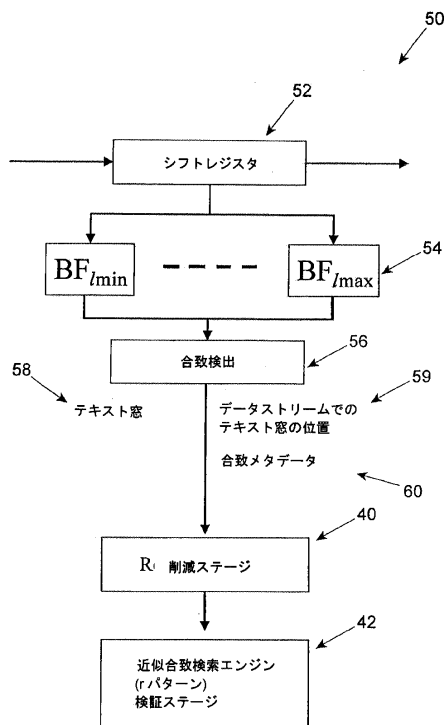


FIG. 3

【図 4】

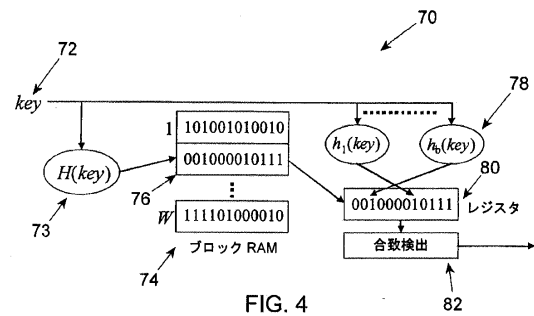


FIG. 4

【図 5】

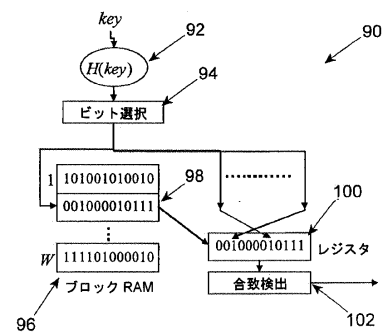


FIG. 5

【 図 6 】

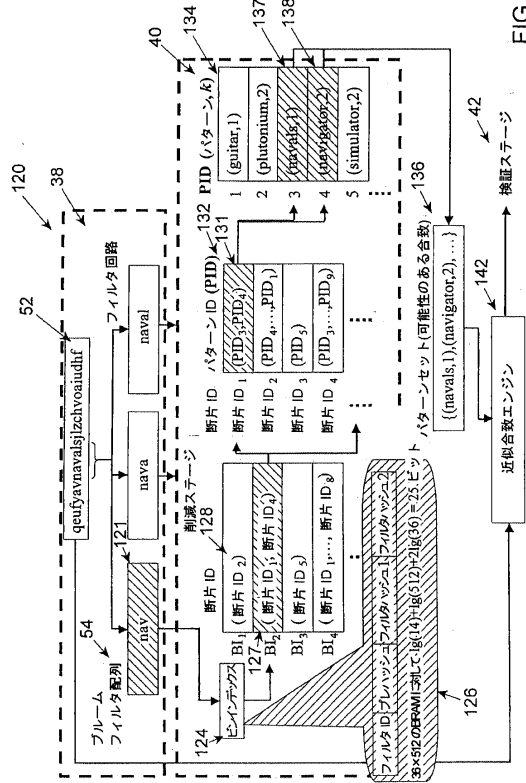


FIG. 6

【圖 7】

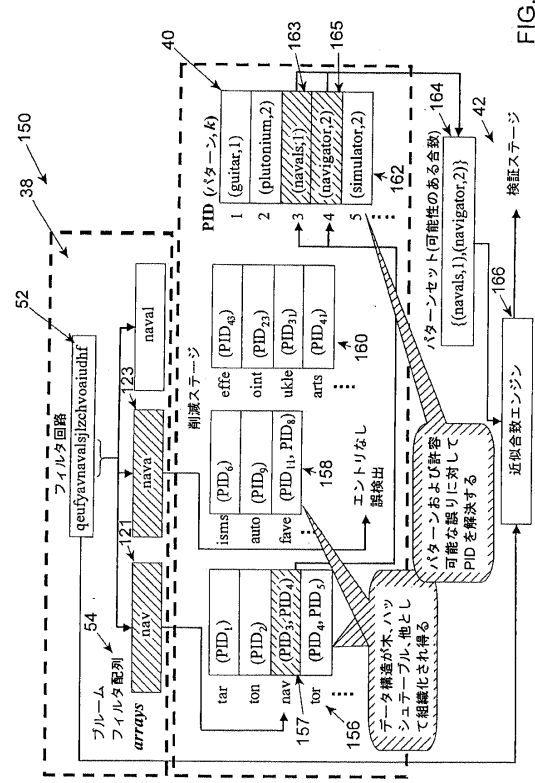


FIG. 7

---

フロントページの続き

(74)代理人 100062007

弁理士 川口 義雄

(72)発明者 テイラー, デイビッド・エドワード

アメリカ合衆国、ミズーリ・63118、セント・ルイス、ミズーリ・アベニュー・3448

審査官 早川 学

(56)参考文献 特表2008-532177(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 17/30