

【公報種別】特許法第17条の2の規定による補正の掲載

【部門区分】第6部門第3区分

【発行日】平成21年8月27日(2009.8.27)

【公開番号】特開2008-140357(P2008-140357A)

【公開日】平成20年6月19日(2008.6.19)

【年通号数】公開・登録公報2008-024

【出願番号】特願2007-70697(P2007-70697)

【国際特許分類】

G 06 F 12/00 (2006.01)

G 06 F 17/30 (2006.01)

【F I】

G 06 F 12/00 520 A

G 06 F 17/30 170 A

G 06 F 17/30 414 B

【手続補正書】

【提出日】平成21年7月9日(2009.7.9)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

記憶装置に入力された複数の電子文書を含む文書集合中の個々の電子文書から単語あるいは連続する複数の文字からなる索引語を切り出して記憶装置に格納する工程と、

それぞれの索引語が各電子文書中に出現する回数を取得する工程と、

前記文書集合中の各電子文書にそれぞれ異なる文書識別番号を割り当てる工程と、

各々の索引語について得られた文書識別番号と索引語出現回数の組であるポスティングをデータ圧縮する工程とをコンピュータに実行させると共に、

前記ポスティングをデータ圧縮する工程は、各ポスティング中の電子文書の文書識別番号を可変長のバイト列で表現し、wを1以上8以下の与えられた整数とするとき、前記バイト列中のwビットによって1から($2^w - 1$)までの索引語出現回数を表現し、 2^w 以上の索引語出現回数を追加のバイト列を用いて表現する

ことを特徴とするプログラム。

【請求項2】

請求項1記載のプログラムにおいて、xを、x+wが1以上8以下となる与えられた整数とする場合、前記ポスティングを表現するバイト列中のxビットに追加情報を書き込むことを特徴とするプログラム。

【請求項3】

請求項1又は2記載のプログラムにおいて、

前記ポスティングをデータ圧縮する工程は、各バイト中の特定の位置のビットにより電子文書の文書識別番号を表現するバイト列の最終バイトか否かを表現すると共に、前記追加のバイト列は、各バイト中の前記特定の位置のビットにより電子文書の文書識別番号を表現する最終バイトではないことを表現し、前記特定の位置のビットに隣接するビットにより当該追加のバイト列の最終バイトか否かを表現する

ことを特徴とするプログラム。

【請求項4】

複数の電子文書を含む文書集合を記憶装置に入力する工程と、

前記入力された文書集合中の個々の電子文書から単語あるいは連続する複数の文字からなる索引語を切り出す工程と、

それぞれの索引語が各電子文書中に出現する回数を取得する工程と、

前記文書集合中の各電子文書にそれぞれ異なる文書識別番号を割り当てる工程と、

各々の索引語について得られた文書識別番号と索引語出現回数からなるポスティングをデータ圧縮する工程と、

データ圧縮されたポスティングを含む転置インデックスを出力する工程とを有し、

前記ポスティングをデータ圧縮する工程では、各ポスティング中の電子文書の文書識別番号を可変長のバイト列で表現し、wを1以上8以下の与えられた整数とするとき、索引語出現回数が1から(2^w-1)までのときは当該索引語出現回数を前記バイト列中のwビットによって表現し、索引語出現回数が2^w以上のときは当該索引語出現回数を追加のバイト列を用いて表現する

ことを特徴とする転置インデックス作成方法。

【請求項5】

請求項4記載の転置インデックス作成方法において、xをx+wが1以上8以下となる与えられた整数とする場合、前記ポスティングを表現するバイト列中のxビットに追加情報を書き込む

ことを特徴とする転置インデックス作成方法。

【請求項6】

請求項4又は5記載の転置インデックス作成方法において、

前記ポスティングをデータ圧縮する工程は、各バイト中の特定の位置のビットにより電子文書の文書識別番号を表現するバイト列の最終バイトか否かを表現すると共に、前記追加のバイト列は、各バイト中の前記特定の位置のビットにより電子文書の文書識別番号を表現する最終バイトではないことを表現し、前記特定の位置のビットに隣接するビットにより当該追加のバイト列の最終バイトか否かを表現する

ことを特徴とする転置インデックス作成方法。

【請求項7】

それぞれに識別番号が割り当てられた複数の電子文書からなる文書集合の全ての索引語について、前記索引語が出現する電子文書の識別番号と前記電子文書中の当該索引語の出現回数の組であるポスティングを格納した転置インデックスを用いた検索手段による検索方法であって、

前記転置インデックスは、各ポスティング中の電子文書の識別番号が可変長のバイト列としてコンピュータにより作成され、wを1以上8以下の与えられた整数とするとき、コンピュータは、前記バイト列中のwビットによって1から(2^w-1)までの出現回数を表現し、2^w以上の出現回数を追加のバイト列を用いて表現し、前記電子文書の識別番号を表現するバイト列は、各バイト中の特定の位置のビットにより当該バイト列の最終バイトか否かを表現し、前記追加のバイト列は、各バイト中の前記特定の位置のビットにより文書識別番号の最終バイトではないことを表現し、前記特定の位置のビットに隣接するビットにより当該追加のバイト列の最終バイトか否かを表現しているとき、

入力部から第1の索引語と第2の索引語を受け付ける工程と、

前記転置インデックス中の前記第1の索引語の転置リストに含まれる電子文書の識別番号を取得する工程と、

前記取得した識別番号が前記第2の索引語の転置リストに含まれているかどうかを2分探索によって判定する工程と、

前記取得した識別番号のうち前記第2の索引語の転置リストに含まれている識別番号またはその識別番号に対応する電子文書を出力する工程とを有し、

前記2分探索に際しては、前記第2の索引語の転置リストの所望の読み位置のバイトに移動し、当該バイトが前記電子文書の識別番号を表現するバイト列の最終バイトか否かを判定し、最終バイトでない場合には最終バイトが見つかるまで前記転置リストの先頭の方向に1バイトずつ読み進み、前記電子文書の識別番号を表現するバイト列の最終バイト中

の出現回数を表現する w ビットの情報から前記出現回数を表現する追加のバイト列が存在するか否かを判定し、追加のバイト列が存在しないと判定されたときは当該バイトの次のバイトをポスティングの先頭としてポスティングを読み取り当該ポスティング中の識別番号を取得し、追加のバイト列が存在すると判定されたときは追加のバイト列の最終バイトが見つかるまで前記転置リストの下流側に 1 バイトずつ読み進み、追加のバイト列の最終バイトの次のバイトをポスティングの先頭としてポスティングを読み取り当該ポスティング中の識別番号を取得する処理を実行する

ことを特徴とする検索方法。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】発明の名称

【補正方法】変更

【補正の内容】

【発明の名称】プログラム、転置インデックスの格納方法及び検索方法