

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4451908号  
(P4451908)

(45) 発行日 平成22年4月14日 (2010. 4. 14)

(24) 登録日 平成22年2月5日 (2010. 2. 5)

(51) Int. Cl.

F I

G 0 6 F 17/21 (2006.01)

G 0 6 F 17/21 5 7 0 L

請求項の数 27 (全 57 頁)

(21) 出願番号	特願2008-50694 (P2008-50694)	(73) 特許権者	503260918
(22) 出願日	平成20年2月29日 (2008. 2. 29)		アップル インコーポレイテッド
(62) 分割の表示	特願2007-182633 (P2007-182633)		アメリカ合衆国 9 5 0 1 4 カリフォル
原出願日	平成8年9月13日 (1996. 9. 13)		ニア州 クパチーノ インフィニット ル
(65) 公開番号	特開2008-192163 (P2008-192163A)	(74) 復代理人	100155284
(43) 公開日	平成20年8月21日 (2008. 8. 21)		弁理士 井原 光雅
審査請求日	平成20年2月29日 (2008. 2. 29)	(74) 代理人	100077481
審判番号	不服2008-32772 (P2008-32772/J1)		弁理士 谷 義一
審判請求日	平成20年12月25日 (2008. 12. 25)	(74) 代理人	100088915
(31) 優先権主張番号	08/527, 438		弁理士 阿部 和夫
(32) 優先日	平成7年9月13日 (1995. 9. 13)	(72) 発明者	エドバーグ, ピーター, ケイ.
(33) 優先権主張国	米国 (US)		アメリカ合衆国 9 7 4 0 3 オレゴン州
(31) 優先権主張番号	08/527, 827		ユージン フェアマウント ブールバー
(32) 優先日	平成7年9月13日 (1995. 9. 13)		ド 1 8 9 2
(33) 優先権主張国	米国 (US)		最終頁に続く

(54) 【発明の名称】 ユニコード・コンバータ

(57) 【特許請求の範囲】

【請求項 1】

ユニコード変換のためにコンピュータ上でソース文字列をターゲット文字列に変換する方法であって、前記コンピュータは、ソース文字列のターゲット文字列への変換を制御するコンバータ手段と、前記ソース文字列をスキャンして前記ソース文字列のテキスト要素を判定するスキャナ手段と、複数のマッピング・テーブルを格納するメモリ手段とを含み、前記方法は、

(a) バッファ手段が、第1の文字コード化を有するソース文字列を受け取ることと、

(b) 前記スキャナ手段が、前記バッファ手段が受け取った前記ソース文字列をテキスト要素に順次分割することであって、各テキスト要素は、前記ソース文字列の1または複数の文字を含み、前記テキスト要素の少なくとも1つは、前記ソース文字列の複数の文字を含むことと、

(c) 前記スキャナ手段が、前記分割 (b) の後または間に前記テキスト要素についての属性情報を得ることと、

(d) 前記コンバータ手段が、前記テキスト要素のそれぞれについて第2の文字コード化に関連付けられた変換コードをマッピング・テーブルでルックアップすることであって、前記テキスト要素のそれぞれについての変換コードを有するマッピング・テーブルでルックアップすること (d) は、要求された変種を識別するオペレーションと、前記属性情報、前記要求された変種、および前記テキスト要素の長さに基づいて前記メモリ手段に格納された複数のマッピング・テーブルの1つを選択するオペレーションと、前記マッピン

10

20

グ・テーブルの選択された１つから変換コードをルックアップすることを含むことと、  
 (e) 前記コンバータ手段が、前記第２の文字コード化のターゲット文字列を形成するように前記テキスト要素の前記変換コードを結合することと  
 を備えることを特徴とする方法。

【請求項２】

請求項１に記載の方法であって、  
 前記変換コードは、前記第２の文字コード化中の１または複数の文字からなり、  
 前記得られた属性情報は、方向、クラス、優先順位、サブセット、コンテキスト、シメトリック・スワッピング状態の少なくともいずれか２つを含むことを特徴とする方法。

【請求項３】

請求項１に記載の方法であって、前記テキスト要素は、互いに隣接し、複数の文字を含む前記テキスト要素のそれぞれについて、前記文字は、前記ソース文字列において隣接することを特徴とする方法。

【請求項４】

請求項１に記載の方法であって、前記マッピング・テーブルは、レギュラ・マッピングとフォールバック・マッピングを含み、

前記ルックアップすること(d)は、前記レギュラ・マッピングを使用して前記マッピング・テーブルが前記テキスト要素についての変換コードを含んでいないとき、前記フォールバック・マッピングを使用して前記テキスト要素のそれぞれについての変換コードを判定することを特徴とする方法。

【請求項５】

請求項１に記載の方法であって、前記文字のそれぞれは、関連する文字クラスを有し、  
 前記分割すること(b)は、前記ソース文字列内の前記文字の文字クラスに少なくとも一部基づくことを特徴とする方法。

【請求項６】

請求項１に記載の方法であって、前記分割すること(b)は、  
 (b1) 前記ソース文字列から次のソース文字を獲得することと、  
 (b2) 獲得した前記ソース文字を現在のテキスト要素に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかをスキャナ・テーブルに基づいて判定することと、

(b3) 前記判定すること(b2)に従い、獲得した前記ソース文字を前記現在のテキスト要素または前記新たな次のテキスト要素に配置することと、

(b4) 前記ソース文字列がテキスト要素に完全に配置されるまで、(b1)から(b3)までを繰り返すことと

を備えることを特徴とする方法。

【請求項７】

請求項６に記載の方法であって、前記判定すること(b2)は、  
 (i) 前記スキャナ・テーブルにおいて前記ソース文字に関連付けられた属性情報をルックアップすることであって、前記属性情報は少なくともクラス・インディケータを含むことと、

(ii) 前記クラス・インディケータに基づき、獲得した前記ソース文字を前記現在のテキスト要素に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかを判定することと

を備えることを特徴とする方法。

【請求項８】

請求項６に記載の方法であって、前記判定すること(b2)は、  
 (i) 前記スキャナ・テーブルにおいて前記ソース文字に関連付けられた属性情報をルックアップすることであって、前記属性情報は少なくともクラス・インディケータを含むことと、

(ii) 複数の状態を有するステート・マシンを提供することであって、前記ステート

10

20

30

40

50

・マシンは、前記クラス・インディケータおよび前記ステート・マシンの現在の状態に基づき、獲得した前記ソース文字を前記現在のテキスト要素に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかを判定するのに使用されることと、

( i i i ) 前記ステート・マシンの前記現在の状態を更新することと  
を備えることを特徴とする方法。

【請求項 9】

請求項 1 に記載の方法であって、前記第 1 の文字コード化および前記第 2 の文字コード化の 1 つは、ユニコード標準に準拠することを特徴とする方法。

【請求項 10】

請求項 1 に記載の方法であって、前記結合すること ( e ) は、

( e 1 ) 前記第 2 の文字コード化についてターゲット文字サイズを判定することと、

( e 2 ) ルックアップされた前記テキスト要素の前記変換コードを、前記ターゲット文字サイズの単位で前記ターゲット文字列にコピーすることにより前記ターゲット文字列を形成することと

を備えることを特徴とする方法。

【請求項 11】

請求項 10 に記載の方法であって、ルックアップ ( d ) された前記変換コードは、間接シーケンスまでのオフセットを指定することを特徴とする方法。

【請求項 12】

請求項 1 に記載の方法であって、

( f ) 前記分割 ( b ) の後であるが前記ルックアップ ( c ) の前に、前記コンバータ手段またはスキャナ手段が、ある文字が前記テキスト要素に存在する場合、前記テキスト要素のそれぞれの前記ある文字を再配列することをさらに備えることを特徴とする方法。

【請求項 13】

請求項 12 に記載の方法であって、前記再配列すること ( f ) は、異なる文字クラスの加重値を用いて実行されることを特徴とする方法。

【請求項 14】

ユニコードについてソース文字列をターゲット文字列に変換するコード変換システムであって、

第 1 の文字コード化を有する前記ソース文字列について第 2 の文字コード化を有する前記ターゲット文字列への変換を制御するコンバータと、

前記コンバータに動作するように結合され、前記ソース文字列をテキスト要素に分割するスキャナであって、各テキスト要素は前記ソース文字列の 1 または複数の文字を含み、前記テキスト要素の少なくとも 1 つは前記ソース文字列の複数の文字を含むスキャナと、

前記ソース・コード化のテキスト要素についてターゲット・コード化を格納するマッピング・テーブルであって、複数のマッピング部分を含み、前記テキスト要素の属性情報および長さに応じて前記マッピング部分の適切な 1 つが前記テキスト要素のそれぞれに利用されるマッピング・テーブルと、

前記コンバータおよび前記マッピング・テーブルに動作するように結合され、前記テキスト要素のそれぞれについて第 2 の文字コード化に関連付けられた変換コードを前記マッピング・テーブルでルックアップするルックアップ・ハンドラと

を備えたことを特徴とするコード変換システム。

【請求項 15】

請求項 14 に記載のコード変換システムであって、

前記コンバータに動作するように結合され、特定の場合にフォールバック変換コードを提供するフォールバック・ハンドラであって、前記ルックアップ・ハンドラが 1 または複数のテキスト要素について変換コードを提供できないとき、前記フォールバック変換コードは、前記テキスト要素の文字と等価ではないが、類似するグラフィカル上の外観を有する、前記ターゲット・コード化における 1 または複数のコード・ポイントを含む、フォールバック・ハンドラをさらに備えたことを特徴とするコード変換システム。

## 【請求項 16】

請求項 15 に記載のコード変換システムであって、

前記入力文字列中の個々の文字を現在のテキスト要素内に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかの判定において前記スキャナを支援するスキャナ・テーブル手段をさらに備えたことを特徴とするコード変換システム。

## 【請求項 17】

請求項 14 に記載のコード変換システムであって、

前記スキャナに動作するよう結合され、前記入力文字列中の個々の文字を現在のテキスト要素内に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかの判定において前記スキャナを支援するスキャナ・テーブルをさらに備えたことを特徴とするコード変換システム。

10

## 【請求項 18】

請求項 17 に記載のコード変換システムであって、前記ソース文字列の前記文字は、関連する文字クラスを有し、

前記スキャナ・テーブルは、要素の配列を備え、前記配列は文字クラスによりインデックス付けされていることを特徴とするコード変換システム。

## 【請求項 19】

請求項 14 に記載のコード変換システムであって、前記ソース文字列中の前記文字は、ユニコード文字であることを特徴とするコード変換システム。

## 【請求項 20】

20

請求項 14 に記載のコード変換システムであって、前記ターゲット文字列中の前記文字は、ユニコード文字であることを特徴とするコード変換システム。

## 【請求項 21】

ユニコード変換のためにコンピュータ上でソース文字列をターゲット文字列に変換する方法であって、前記コンピュータは、ソース文字列のターゲット文字列への変換を制御するコンバータ手段と、前記ソース文字列をスキャンして前記ソース文字列のテキスト要素を判定するスキャナ手段と、複数のマッピング・テーブルを格納するメモリ手段とを含み、前記方法は、

(a) バッファ手段が、第 1 の文字コード化を有するソース文字列を受け取ることと、

(b) 前記スキャナ手段が、前記バッファが受け取った前記ソース文字列をテキスト要素に順次分割することであって、各テキスト要素は、前記ソース文字列の 1 または複数の文字を含み、前記テキスト要素の少なくとも 1 つは、前記ソース文字列の複数の文字を含むことと、

30

(c) 前記スキャナ手段が、前記分割 (b) の後または間に前記テキスト要素についての属性情報を得ることと、

(d) 前記コンバータ手段が、前記テキスト要素のそれぞれについて第 2 の文字コード化に関連付けられた変換コードをマッピング・テーブルでルックアップすることであって、前記テキスト要素のそれぞれについての変換コードを有するマッピング・テーブルでルックアップすること (d) は、(d1) 前記属性情報、および前記テキスト要素の長さに基づいて前記メモリ手段に格納された複数のマッピング・テーブルの 1 つを選択するオペレーションと、(d2) 前記マッピング・テーブルの選択された 1 つから変換コードをルックアップするオペレーションを含むことと、

40

(e) 前記コンバータ手段が、前記第 2 の文字コード化のターゲット文字列を形成するように前記テキスト要素の前記変換コードを結合することとを備えることを特徴とする方法。

## 【請求項 22】

請求項 21 に記載の方法であって、前記分割すること (b) は、

(b1) 前記ソース文字列から次のソース文字を獲得することと、

(b2) 獲得した前記ソース文字を現在のテキスト要素に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかをスキャナ・テーブルに基づいて判定すること

50

と、

(b3) 前記判定すること(b2)に従い、獲得した前記ソース文字を前記現在のテキスト要素または前記新たな次のテキスト要素に配置することと、

(b4) 前記ソース文字列がテキスト要素に完全に配置されるまで、(b1)から(b3)までを繰り返すことと

を備えることを特徴とする方法。

【請求項23】

請求項22に記載の方法であって、前記判定すること(b2)は、

(i) 前記スキャナ・テーブルにおいて前記ソース文字に関連付けられた属性情報をルックアップすることであって、前記属性情報は少なくともクラス・インディケータを含むことと、

10

(ii) 複数の状態を有する状態・マシンを提供することであって、前記状態・マシンは、前記クラス・インディケータおよび前記状態・マシンの現在の状態に基づき、獲得した前記ソース文字を前記現在のテキスト要素に含めるべきか、あるいはまた、新たな次のテキスト要素として始めるかを判定するのに使用されることと、

(iii) 前記状態・マシンの前記現在の状態を更新することと

を備えることを特徴とする方法。

【請求項24】

請求項23に記載の方法であって、

(f) 前記分割(b)の後であるが前記ルックアップ(c)の前に、前記コンバータ手段またはスキャナ手段が、ある文字が前記テキスト要素に存在する場合、前記テキスト要素のそれぞれの前記ある文字を再配列することをさらに備えることを特徴とする方法。

20

【請求項25】

請求項14に記載のコード変換システムであって、

前記マッピング・テーブルは、前記ソース・コード化のあるテキスト要素について1つ以上のターゲット・コード化を含み、前記あるテキスト要素についての前記ターゲット・コード化の特定の1つを前記あるテキスト要素について判定された属性情報によって獲得することを特徴とするコード変換システム。

【請求項26】

請求項14に記載のコード変換システムであって、

30

前記コード変換システムは、

スキャン情報を格納するスキャナ・テーブルと、

前記ソース文字列中の文字についての属性情報を格納する属性テーブルと

をさらに備え、

前記スキャナは、前記スキャナ・テーブルに格納されたスキャン情報および前記属性テーブルに格納された属性情報に従って前記ソース文字列をテキスト要素に分割するように動作することを特徴とするコード変換システム。

【請求項27】

請求項21に記載の方法であって、

前記テキスト要素のそれぞれについて変換コードのマッピング・テーブルでルックアップすること(d)は、(d3) 前記選択(d1)の前に要求された変種を特定するオペレーションをさらに含み、

40

前記複数のマッピング・テーブルの1つを選択すること(d2)は、前記属性情報、前記要求された変種、および前記テキスト要素の長さに基づき行われることを特徴とする方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は記述または表示されるテキストの文字コード間で変換を行うシステムに関し、さらに具体的には、ある文字セットと別の文字セット間で変換を行うコード・コンバータ

50

に関する。

【背景技術】

【0002】

コンピュータや他のエレクトロニック・デバイスはテキストを使用してユーザと対話するのが一般的である。テキストはモニタや他のタイプの表示デバイスから表示されるのが普通である。テキストはコンピュータや他のエレクトロニック・デバイスの内部ではデジタル形式で表現しなければならないので、文字セット・コード化を使用しなければならない。一般的には、文字セット・コード化は文字セットの各文字を固有のデジタル表現でコード化するように働く。文字（コード化される）は文字、数字および種々のテキスト・シンボルに対応し、コンピュータや他のエレクトロニック・デバイスで使用される数値コードが割り当てられている。コンピュータや他のエレクトロニック・デバイスで使用されている最もよく知られた文字セットは情報交換用米国標準コード(American National Standard Code for Information Exchange - ASCII)である。ASCIIはそのコード化のために7ビット・シーケンスを使用している。他の国では、異なる文字セットが使用されている。ヨーロッパでは、支配的な文字コード化標準は国際標準化機構(International Standards Organization - ISO)によって開発されたISO 8859 - X、特にISO 8859 - 1(「Latin - 1」と呼ばれている)である。日本では、支配的な文字コード化標準はJIS X 0208であり、ここでJISとは日本情報標準(Japanese Information Standard)のことであり、日本標準協会(Japan Standards Association - JSA)によって開発されたものである。その他の既存文字セットの例としては、Mac(登録商標)OS Standard Romanコード化(Apple Computer, Inc 開発)、Shift - JIS(日本)、Big 5(台湾)、その他多くがある。

【発明の開示】

【発明が解決しようとする課題】

【0003】

ビジネスとネットワークのグローバル化に伴い、コンピュータや他のエレクトロニック・デバイスが複数の文字コード化を処理する能力をもつことが重要になっている。例えば、同じコンピュータやエレクトロニック・デバイスは、コンピュータや他のエレクトロニック・デバイスと異なる言語で対話することを望んでいる異なる国籍の人達で利用されることがある。このような各言語では、異なる文字セット・コード化が必要になるのが通常である。しかし、文字セットは同じ言語の場合でも、異なることがある。

【0004】

もう1つの要求は、ある文字セット・コードから別の文字セット・コード化への変換ができることである。例えば、ISO 8859 - 1を使用するフランスのユーザは、電子メール・メッセージをISO 8859 - 8を使用するイスラエルのユーザにフランス語で送りたい場合がある。送信者と受信者は異なる文字セット・コード化を使用しているので、メッセージの中の非ASCII文字は、イスラエルのユーザには誤って伝えられることになる。理想的なことは、コンピュータまたはエレクトロニック・デバイスの1つがある文字セットから別の文字セットに変換することである。このことは、限られた範囲でいくつかの文字セット間ですでに達成されているが、最新のコンピュータやエレクトロニック・デバイスでもほとんど不可能である。コード変換を困難にしているのは、異なる文字標準が多数存在し、各国の標準が矛盾し、または不統一であることがよくあるためである。

【0005】

Unicode(登録商標)標準(以下では、単にユニコードまたはユニコード標準という)は国際的文字コード化標準となるために開発されたものである。ユニコード標準の設計者はより効率的で、柔軟性のある文字識別方法を必要として、提供したものである。ユニコード標準には、1990年12月31以前に承認され、公表されたすべての主要国際標準のすべての文字が含まれている他に、以前の標準になかった他の文字も含まれている。これらの文字は重複することなくユニコード標準にコード化されている。ユニコード

標準に含まれるコードは16ビット（または2バイト）幅である。

【0006】

ユニコード標準などの文字セット標準はコード変換を容易にし、テキスト・データに働きかける有用なプロセスのインプリメンテーションを可能にしている。例えば、上記の例によれば、フランス国内のコンピュータや他のエレクトロニック・デバイスはユニコード文字を送信することができ、イスラエル国内のコンピュータや他のエレクトロニック・デバイスは受信したユニコード文字を、イスラエル国内のコンピュータや他のエレクトロニック・デバイスと互換性のあるヘブライ語の文字に変換することができる。

【0007】

以下では、ユニコード標準の概要について説明する。なお、ユニコード標準の詳細は、例えば、「The Unicode Standard, Worldwide Character Encoding, Version 1.0, Addison-Wesley 1991」に記載されている（バージョン1.1にも詳細が記載されている）。これらのどちらのバージョンも、その全体は引用により本明細書の一部を構成するものである。

【0008】

ユニコード・コード化方式の設計は方向性を除き、基本的テキスト処理アルゴリズムの設計から独立している。ユニコード・インプリメンテーションは適当なテキスト処理および/またはレンダリング・アルゴリズムを含んでいるものと想定されている。便宜上、ユニコード標準のすべてのコードは言語および機能のカテゴリ別に分類されているが、ユニコード標準のすべてのコードは等しくアクセス可能になっている。「ユニコード標準のコード・スペースは6つのゾーンに分割されている。すなわち、一般的スクリプト（アルファベットおよび相対的に小さな文字セットをもつ他のスクリプト）、シンボル（記号）、CJK（中国語、日本語および朝鮮語）補足文字、CJK表意文字、プライベート使用および互換性である。一般的スクリプト・ゾーンはラテン、キリル、ギリシア、ヘブライ、アラビア、デーヴァナーガリーおよびタイなどのアルファベットまたは音節スクリプトを含んでいる。シンボル・ゾーンは句読点、数学、化学、読者の注意をひく記号などの、様々な種類の文字を含んでいる。CJK補足ゾーンは句読点、シンボル、カナ文字、Bopomofo文字、および単一および複合ハングル文字を含んでいる。CJK表意文字ゾーンには、20,000個以上の表意文字または他のスクリプトからの文字のためのスペースが用意されている。プライベート使用ゾーンはユーザまたはベンダ固有のグラフィック文字を定義するために使用される。互換性ゾーンはユニコード・コード化で他の規範的表現をもつ、幅広く使用されている企業および国内標準からの文字を含んでいる。」（上記標準バージョン1.0の13ページから引用）。

【0009】

ユニコード標準は文字特性と制御文字を含んでいる。文字特性はユニコード・コード化に含まれるコード・ポイントに関する意味論的知識を必要とする解析、ソート、および他のアルゴリズムで利用すると便利である。ユニコード標準で指定されている文字特性には、ディジット、数字、スペース文字、非スペース・マーク、および方向がある。ユニコード文字は文字特性に基づいてグループ化されている。ディジット、数字およびスペース文字は周知である。非スペース・マーク・グループは非スペース・マークを収容し、方向グループは方向文字を収容している。

【0010】

10

20

30

40

非スペース・マーク（例えば、ギリシアとローマ・スクリプトにおけるアクセント・マーク、アラビアとデーヴァナーガリーにおける母音マーク）は最終的にレンダリングされたテキストには線形的に現れない。ユニコード文字コード・シーケンスでは、このような文字はすべて、その文字で修正されたベース文字、またはその後には音声順（phonetic order）に分節される文字のあとに置かれている（例えば、ローマ文字の“Á”は事前構成形式（precomposed form）でストアされないとき“A<sup>´</sup>”としてストアされる）。これらの文字（つまり、非スペース・マーク）はレンダリングされるとき、その前に置かれたベース文字に対してなんらかの方法で位置付けられることを目的とし、これらの文字自体はスペース位置を占めない。制御文字はユニコード標準でコード化されているが、制御文字自体はグラフィック文字ではない。これらの制御文字は、例えば、水平タブを示したり、フォーマッティング属性または構造などのテキストに関する追加情報を提供したり、方向性フォーマッティングによって制御したりするために使用される。ユニコード標準は双方向文字配列も用意しているので、ユニコード・コード化方式には、方向の変更を指定する文字も含まれている。例えば、ギリシア語、ローマ語およびタイ語では、左から右への方向が支配的になっているのに対し、アラビア語とヘブライ語では、右から左への方向が支配的になっている。制御文字を使用して方向を変更することは、例えば、左から右への文字を右から左への順序で表示したいときに、ユーザまたはシステムによって要求されることがある。

10

20

#### 【 0 0 1 1 】

従来のコード・コンバータには、1つの問題がある。それは、コンバータが1つのソース文字しか1つのターゲット文字に変換できないことである。この種の変換はある種の文字セットでは有効に働くが、多数の非スペース文字セットを含むある種の文字セット（例えば、ユニコード）の場合や、通常他の文字と関連づけられていないマークを結合するときに十分ではない。また、従来のコード・コンバータは複数の文字に関連するシンボル、合字（ligature）または表意文字との間で変換できる機能を備えていない。

30

#### 【 0 0 1 2 】

その結果として、従来のコード・コンバータにはラウンドトリップ忠実度（round trip fidelity）がない。ラウンドトリップ忠実度とは、コード・コンバータが変換し、そのあと変換を戻してオリジナルの入力文字列を再現できる能力のことである。これが重要になるのは、コード・コンバータをある文字コードから別の文字コードに変換するときのハブとして使用するときである。

#### 【 0 0 1 3 】

従って、複数のソース文字から単一ターゲット文字へ、あるいは単一ソース文字から複数のターゲット文字へ変換して、ラウンドトリップ文字コード変換の忠実度を保証するようなコード・コンバータが望まれている。

40

#### 【 0 0 1 4 】

従来のコード・コンバータには、もう1つ問題がある。それは、ソース文字セットの文字をターゲット文字セットの文字に変換するとき方向を考慮に入れないことである。方向を考慮に入れないと、ある文字は左から右へ配列されるのに対し、他の文字は右から左へ配列されるので変換にエラーを生じることになる。これが起こる代表例として、ある与えられたソース・コード化に対して2つの等価文字をもつターゲット文字に変換し、違いが方向だけのときである。このような場合には、正しいターゲット文字をマッピングするために、ソース文字の方向が分かっている必要はない。また、従来のコード・コンバータが不十分であるのは、左から右へ配列されている（方向性）言語だけでなく、右から左

50



へ配列されている言語にも対応する文字を含んでいるある種の文字セット（例えば、ユニコード）を取り扱えるだけの柔軟性がないためである。例えば、文字の配列または方向性がユニコード文字列内で変化するようなユニコード文字は、従来のコード・コンバータが文字列全体の方向が一定であることを前提としているので従来のコード・コンバータでは正しく変換されないことになる。

【 0 0 1 5 】

従って、文字の方向を考慮に入れながら、ソース文字セットの文字をターゲット文字セットの文字に変換することができるようなコード・コンバータが望まれている。

【 0 0 1 6 】

ある種の従来のコード・コンバータには、以下のようなとき起こる別の問題がある。ネットワークからテキストを受信するとき、そのテキストに関連するデータはデータ・ブロックで到着するのが代表的である。データが完全な形で受信されるのは、データを構成するすべてのブロックが受信されたあとだけである。受信データはバッファに置かれ、そこでデータは文字コードの変換を待っている。しかし、バッファは、多くの場合、データ・ストリーム全体を格納する能力もなければ、その１ブロックさえも格納する能力をもっていないこともある。いずれの場合も、スキャナ４０８によって後でテキスト要素と判断されるはずのものの途中でバッファの終わり（つまり、バッファ内の最後の文字）が現れることがある。テキスト要素の途中でバッファの終わりが現れると、スキャナ４０８は、後続文字が最後のテキスト要素に影響したり、その一部であることがあるので最後のテキスト要素を正しく得ることができなくなる。

【 0 0 1 7 】

大部分のコード・コンバータには、別の問題がある。それは、あるソース文字セットの文字をあるターゲット文字セットの文字に変換するときコンテキストを考慮に入れないことである。ある種の文字セット（例えば、ユニコード）では、文字セットにはコンテキスト表現形式が異なるごとに別々の文字コードが含まれているのに対し、別の時では、単一文字コードだけが用意され、コンテキストは表現形式を判断するために使用されている。しかし、従来のコード・コンバータはコードをそのコンテキストに従って変換することができない。コンテキスト・ベースのコード変換が特に問題となるのは、変換の対象となる文字セット（例えば、ユニコード）が文字コンテキストを受け入れるために、いくつかの手法を組み合わせて利用するときである。

【 0 0 1 8 】

従って、文字のコンテキストを考慮に入れながら、ソース文字セットの文字をターゲット文字セットの文字に正確かつ柔軟に変換できるようなコード・コンバータが望まれている。

【課題を解決するための手段】

【 0 0 1 9 】

概要を説明すると、本発明はコード変換および／または切り捨て処理手法に関するものである。

【 0 0 2 0 】

本発明の第１側面では、コード変換手法によれば、ラウンドトリップ忠実度が得られると共に、その結果の文字コードが他のプラットフォームと互換性をもつことが保証される。コード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかにマッピングする能力をもっている。ラウンドトリップ忠実度によれば、ソース・テキストはターゲット・テキストに変換し、その後オリジナル・ソース・テキストに再び戻るように変換することが可能である。互換性は標準ターゲット文字の使用を最大限にし、プライベート文字の使用を最小限にすることにより保証される。コード変換は他の文字セットからユニコード文字へ、ユニコード文字から他の文字セットへ変換するとき利用すると、特に便利である。ユニコード文字シーケンスをターゲット文字セット内の単一文字にマッピングすることは、従来は提供されていなかった。本発明は、オペレーションだけでなく、データ保管においても大幅な柔軟性が得られる

強力な解決法を提供している。本発明は様々な態様で実現することができる。つまり、方法、装置またはシステムとして実現することもコンピュータ可読メディア上に実現することも可能である。

【0021】

ソース文字列をターゲット文字列に変換する方法としては、本発明の第1側面による本発明の実施例は、第1文字コード化をもつソース文字列を受信し、ソース文字列をテキスト要素に順次に分割し（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）、テキスト要素の各々ごとに第2文字コード化に関連する変換コードをマッピング・テーブルで調べ（ルックアップ）、テキスト要素の変換コードを結合して第2文字コード化のターゲット文字列を形成するオペレーションを実行する。

10

【0022】

さらに、本発明の第1側面によれば、マッピング・テーブルにはレギュラ・マッピングとフォールバック・マッピングを含めることが可能であり、ルックアップ・オペレーションは、レギュラ・マッピングを使用するテキスト要素の変換コードがマッピング・テーブルにないとき、フォールバック・マッピングを使用するテキスト要素の各々の変換コードを判断することが可能である。また、好ましいことは、文字の各々にそれぞれに関連する文字クラスをもたせ、少なくともその一部がソース文字列内の文字の文字クラスに基づいてソース文字列の分割が行われるようにすることである。

【0023】

ソース文字列をターゲット文字列に変換するコード変換システムとしては、本発明の第1側面による本発明の実施例は、第1文字コード化をもつソース文字列を第2文字コード化をもつターゲット文字列に変換することを制御するコンバータ、ソース文字列をテキスト要素（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）に分割するスキャナ、ソース・コード化のテキスト要素に対するターゲット・コード化をストアしておくマッピング・テーブルおよびテキスト要素の各々ごとに第2文字コード化に関連する変換コードをマッピング・テーブルで調べる（ルックアップ）ルックアップ・テーブルとを含んでいる。

20

【0024】

本発明の第1側面によるコード変換システムには、さらに、フォールバック・ハンドラとスキャナ・テーブルを含めることが可能である。フォールバック・ハンドラは、ルックアップ・ハンドラが1つまたは複数のテキスト要素の変換コードを提供することができないような、特殊な場合にフォールバック変換コードを提供する。フォールバック変換コードはテキスト要素内の文字と正確には等価ではないが、類似したグラフィック外観をもつターゲット・コード化内の1つまたは複数のコード・ポイントを含んでいる。スキャナ・テーブルは、入力文字列内の個別文字を現在のテキスト要素内に含めるべきか、それとも新しい次のテキスト要素を開始すべきかをスキャナが判断するのを支援するものである。

30

【0025】

入力文字列をスキャンするスキャンング・システムとしては、本発明の第1側面による本発明の実施例は、入力文字を入力文字列から取得する入力デバイス（入力文字列は文字コード化をもち、入力文字列の各文字は文字クラスをもっている）、入力文字の属性を提供する属性テーブルおよびステート・マシンの次の状態と次のアクションの両方を、入力文字の次の状態とステート・マシンの現在状態に従って判断するステート・マシンを含んでいる。スキャンング・システムは入力文字列の入力文字を現在のテキスト要素内に含めるべきか、現在のテキスト要素を終わらせて次のテキスト要素を開始させるかを、ステート・マシンが判断したアクションに基づいて判断する。好ましくは、属性には、入力文字の文字クラスが含まれおり、ステート・マシンは入力文字の文字クラスとステート・マシンの現在状態に基づいて次の状態と次のアクションを判断する。

40

【0026】

ソース文字列をターゲット文字列に変換するためのプログラム命令を収めておくコンピュータ可読メディアとしては、本発明の第1側面による本発明の実施例は、第1文字コー

50

ド化をもつソース文字列を受信することをコンピュータに実行させるように構成されたコンピュータ可読コード、ソース文字列をテキスト要素（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）に分割することをコンピュータに実行させるように構成されたコンピュータ可読コード、テキスト要素の各々ごとに第2文字コード化に関連する変換コードを調べる（ルックアップ）ことをコンピュータに実行させるように構成されたコンピュータ可読コード、およびテキスト要素の変換コードを結合して第2文字コード化のターゲット文字列を形成することをコンピュータに実行させるように構成されたコンピュータ可読コードを含んでいる。

【0027】

本発明の第2側面では、コード変換システムはソース文字コード化からの文字をターゲット文字コード化に変換するとき方向を考慮に入れる。

10

【0028】

本発明の第2側面によるコード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかにマッピングする能力を備えている。変換される文字の方向を判断または解明することにより、コード変換システムは判断または解明された文字の方向を利用して、ターゲット文字・コード化への正しいマッピングが利用されることを保証する。従って、本発明によれば、ソース文字列内の文字の方向性が変化したときでも正しいコード変換が達成される。本発明は様々な態様で実現することが可能である。方法、装置またはシステムとして実現することも、コンピュータ可読メディア上に実現することも可能である。

20

【0029】

ソース文字列をターゲット文字列に変換するコード変換システムとしては、本発明の第2側面による本発明の実施例は、ソース文字コード化をもつ入力文字列をターゲット文字コード化をもつターゲット文字列（入力文字列は複数の文字を含んでいる）に変換することを制御するコンバータ、入力文字列内の文字を判断するスキャナ、ソース・コード化の文字に対するターゲット・コード化をストアしておくマッピング・テーブル、および入力文字列内の文字の各々ごとにターゲット文字コード化に関連する変換コードを、入力文字列内の文字の方向とソース・コード化に基づいてマッピング・テーブルで調べるルックアップ・ハンドラを含んでいる。好ましくは、スキャナはさらに入力文字列をテキスト要素（各テキスト要素は入力文字列の1つまたは複数の文字を含んでいる）に分割し、スキャナはテキスト要素の方向を判断し、マッピング・テーブルはソース・コード化に対するターゲット・コード化をストアしており、ルックアップ・ハンドラはテキスト要素の各々ごとにターゲット文字コード化を、テキスト要素内の文字の方向とソース・コード化に基づいて調べる。この実施例には、前記ルックアップ・ハンドラが1つまたは複数の変換コードを提供できないような、特定な場合にはフォールバック変換コードを提供するフォールバック・ハンドラを含めることも可能である。

30

【0030】

ソース文字列をターゲット文字列に変換する方法としては、本発明の第2側面による本発明の実施例は、第1文字コード化をもつソース文字列を受信し（ソース文字列は複数のソース文字を含んでいる）、ソース文字列のソース文字の方向を判断し、ソース文字の各々ごとに第2文字コード化に関連する変換コードを、第1文字コード化と判断された方向に基づいてマッピング・テーブルで調べ、ソース文字の変換コードを結合して第2文字コード化のターゲット文字列を形成するオペレーションを実行する。

40

【0031】

ソース文字列をターゲット文字列に変換するためのプログラム命令を収めておくコンピュータ可読メディアとしては、本発明の第2側面による本発明の実施例は、第1文字コード化をもつソース文字列を受信することをコンピュータに実行させるように構成されたコンピュータ可読コード、ソース文字列内のソース文字の各々の方向を判断し、ソース文字列をテキスト要素（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）に分割することをコンピュータに実行させるように構成されたコンピュータ可読コード

50

、テキスト要素の各々ごとに第2文字コード化に関連する変換コードを調べることコンピュータに実行させるように構成されたコンピュータ可読コード、およびテキスト要素の変換コードを結合して第2文字コード化のターゲット文字列を形成することをコンピュータに実行させるように構成されたコンピュータ可読コードを含んでいる。

【0032】

本発明の第3側面では、コード変換システムはソース文字コード化からの文字をターゲット文字コード化に変換するときコンテキストを考慮に入れる。

【0033】

本発明の第3側面によるコード変換システムは単一ターゲット文字またはターゲット文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかに変換する能力を備えている。変換される文字のコンテキストを判断することにより、コード変換システムは判断された文字のコンテキストを利用してターゲット文字コード化への正しいマッピングが利用されることを保証する。従って、本発明によれば、文字のコンテキストが異なるターゲット・コード化に導くときでも正しいコード変換が達成される。本発明は様々な態様で実現することが可能である。方法、装置またはシステムとして実現することも、コンピュータ可読メディア上に実現することも可能である。

【0034】

ソース文字列をターゲット文字列に変換するコード変換システムとしては、本発明の第3側面による本発明の実施例は、ソース文字コード化をもつ入力文字列（入力文字列は複数の文字を含んでいる）をターゲット文字コード化をもつターゲット文字列に変換することを制御するコンバータ、入力文字列内の文字の各々のコンテキストを判断するスキャナ、ソース・コード化の文字に対するターゲット・コード化をストアしておくマッピング・テーブル、および入力文字列内の文字の各々ごとにターゲット文字コード化に関連する変換コードを、コンテキストと入力文字列内の文字のソース・コード化に基づいて前記マッピング・テーブルで調べるルックアップ・ハンドラを含んでいる。好ましくは、スキャナはさらに入力文字列をテキスト要素（各テキスト要素は入力文字列の1つまたは複数の文字を含んでいる）に分割する。この実施例には、ルックアップ・ハンドラが1つまたは複数のテキスト要素の変換コードを提供できないような、特定の場合にはフォールバック変換コードを提供するフォールバック・ハンドラを含めることも可能である。

【0035】

ソース文字列をターゲット文字列に変換するコード変換システムとしては、本発明の第3側面による本発明の実施例は、ソース文字コード化をもつソース文字列をターゲット文字コード化をもつターゲット文字列に変換することを制御するコンバータ手段、ソース内の文字の各々のコンテキストを判断するステート・マシン、ソース文字コード化の文字に対するターゲット文字コード化をストアしておくマッピング手段、およびソース文字列内の文字の各々ごとにターゲット文字コード化に関連する変換コードを前記マッピング手段で調べるルックアップ・ハンドラ手段を含んでいる。

【0036】

ソース文字列をターゲット文字列に変換する方法としては、本発明の第3側面による本発明の実施例は、第1文字コード化をもつソース文字列を受信し（ソース文字列は複数のソース文字を含んでいる）、ソース文字列のソース文字の各々のコンテキストを判断し、ソース文字の各々ごとに第2文字コード化に関連する変換コードを、第1文字コード化とソース文字について判断されたコンテキストに基づいてマッピング・テーブルで調べ、ソース文字の変換コードを結合して第2文字コード化のターゲット文字列を形成するオペレーションを実行する。

【0037】

ソース文字列をターゲット文字列に変換するためのプログラム命令を収めておくコンピュータ可読メディアとしては、本発明の第3側面による本発明の実施例は、第1文字コード化をもつソース文字列を受信することをコンピュータに実行させるように構成されたコンピュータ可読コード、ソース文字列内のソース文字の各々のコンテキストを判断し、ソ

ース文字列をテキスト要素（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）に分割することをコンピュータに実行させるように構成されたコンピュータ可読コード、テキスト要素の各々ごとに第2文字コード化に関連する変換コードを調べることがコンピュータに実行させるように構成されたコンピュータ可読コードおよびテキスト要素の変換コードを結合して第2文字コード化のターゲット文字列を形成することをコンピュータに実行させるように構成されたコンピュータ可読コードを含んでいる。

【0038】

本発明の第4側面では、切り捨て処理手法は、変換のために受信されたソース文字列が、そのソース文字列が変換のためにソース文字列を格納している入力バッファの長さを越えたときでも、異なるターゲット文字コード化に正確に変換されることを保証する。

10

【0039】

本発明の第4側面による切り捨て処理は、入力バッファに置かれているソース文字列の一部を切り捨て、その切り捨てられた部分がソース文字列内の後続文字に影響されることなくターゲット文字コード化に変換されるように作用する。当然のことながら、入力バッファに置かれていて、切り捨てが行われた後のソース文字列の部分（切り捨てられた部分）は最大限にしておく、コード変換が効率よく行えるようになる。本発明は、入力ソース文字列がネットワーク経由で受信されるデータであるときに特に便利である。例えば、データはインターネット経由で電子的に転送されるテキスト・データにすることができる。いずれの場合も、本発明は種々の態様で実現することが可能である。方法、装置またはシステムとして実現することも、コンピュータ可読メディア上に実現することも可能である。

20

【0040】

ソース文字列をターゲット文字列に変換するコード変換システムとしては、本発明の第4側面による本発明の実施例は、第1文字コード化をもつソース文字列を第2文字コード化をもつターゲット文字列に変換することを制御するコンバータ、ソース文字列の一部を一度に受け入れるバッファ（ソース文字列は1つまたは複数の文字を含んでいる）、ソース文字列の一部を切り捨てるランケータ、ソース文字列の切り捨てられた部分をテキスト要素（各テキスト要素はソース文字列の切り捨てられた部分の1つまたは複数の文字を含んでいる）に分割するスキャナ、ソース・コード化のテキスト要素に対するターゲット・コード化をストアしておくマッピング・テーブル、およびテキスト要素の各々ごとに第2文字コード化に関連する変換コードをマッピング・テーブルで調べるルックアップ・ハンドラを含んでいる。

30

【0041】

ソース文字列を切り捨ててターゲット文字列に文字変換するための第1方法としては、本発明の第4側面による本発明の実施例は、ソース文字列のバッファ部分をバッファに受け入れ（ソース文字列は2つ以上のバッファ部分を含んでいる）、ソース文字列の後続バッファ部分に影響されることなくターゲット文字列に変換することができる、ソース文字列のバッファ部分のサブパートを判断し、ソース文字列のバッファ部分のサブパートをターゲット・コード化に変換するオペレーションを実行する。好ましくは、第1方法は、変換対象となるバッファ部分の残余部分がサブパートと一緒になったときバッファ部分に等しい場合には、その残余部分を次のバッファ部分またはそのサブパートと一緒に格納するオペレーションも実行する。

40

【0042】

ソース文字列を切り捨ててターゲット文字列に文字変換するための第2方法としては、本発明の第4側面による本発明の実施例は、ソース文字列の一部をバッファに受け入れ、ソース文字列のその部分内のテキスト要素（各テキスト要素はソース文字列の1つまたは複数の文字を含んでいる）を判断し、テキスト要素が完成しているかどうかを判断し、テキスト要素が完成しているときソース文字列の切り捨てられた部分にテキスト要素を挿入し、ソース文字列の部分が完全に考慮されるまでテキスト要素の判断を繰り返し、そのあと文字変換のためにソース文字列の切り捨てられた部分を出力するオペレーションを実行す

50

る。好ましくは、この第2方法は、バッファに受け入れられたソース文字列の次の部分と一緒に使用するために、ソース文字列の残余部分を格納しておくオペレーションも実行する。

【0043】

ソース文字列をターゲット文字列に変換するためのプログラム命令を格納しておくコンピュータ可読メディアとしては、本発明の第4側面による本発明の実施例は、第1文字コード化をもつソース文字列の一部をバッファに受け入れることをコンピュータに実行させるように構成されたコンピュータ可読コード、ソース文字列の一部を切り捨てることをコンピュータに実行させるように構成されたコンピュータ可読コード、ソース文字列の切り捨てられた部分をテキスト要素（各テキスト要素はソース文字列の切り捨てられた部分の1つまたは複数の文字を含んでいる）に分割することをコンピュータに実行させるように構成されたコンピュータ可読コード、テキスト要素の各々ごとに第2文字コード化に関連する変換コードを調べることをコンピュータに実行させるように構成されたコンピュータ可読コード、およびテキスト要素の変換コードを結合して第2文字コード化のターゲット文字列を形成することをコンピュータに実行させるように構成されたコンピュータ可読コードを含んでいる。

10

【0044】

本発明のその他の側面および利点は、本発明の原理を例示している添付図面を参照して以下に詳述している説明の中で明らかにする。

【0045】

20

本発明の理解を容易にするために、以下では、添付図面を参照して本発明を詳しく説明する。なお、図面において、類似する構成要素は類似の参照符号を付けて示されている。

【発明を実施するための最良の形態】

【0046】

以下、図1ないし図29を参照して本発明の実施例について説明する。なお、この分野の精通者ならば理解されるように、これらの図を参照して以下に詳述する説明は本発明の理解を容易にすることを目的とし、本発明は以下に説明する実施例に限定されるものではない。

【0047】

本発明の第1側面によるコード変換システムはソース文字を異なるコード化のターゲット文字に変換するものである。本発明によるコード変換システムはラウンドトリップ忠実度を備えていると共に、その結果の文字コードが他のプラットフォームと互換性をもつことを保証する。コード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかにマッピングする機能を備えている。ラウンドトリップ忠実度により、ソース・テキストはターゲット・テキストに変換し、そのあとオリジナル・ソース・テキストに戻すように再変換することができる。その結果の文字コードが他のプラットフォームと互換性をもつことは、標準ターゲット文字の使用を最大限にし、プライベート文字の使用を最小限にすることによって保証される。コード変換システムは他の文字セットからユニコード文字に変換し、ユニコード文字を他の文字セットに変換するときに利用すると、特に便利である。ユニコード文字シーケンスをターゲット文字セット内の単一文字にマッピングすることは従来にはなかったものである。

30

40

【0048】

コード変換システムはコンピュータ・システムや他のエレクトロニック・デバイスにしてこれらのコード変換オペレーションを行うことができる。このコンピュータ・システムは必要とする目的だけに使用するように構築することも、コンピュータ・プログラムに従って動作する汎用コンピュータにすることも可能である。以下に説明する処理はどのコンピュータ・システムや他のエレクトロニック・デバイスにも適用可能である。具体的には、種々の汎用コンピューティング・マシンは以下に教示する説明に従って書かれたソフトウェアと共に使用することができるが、もっと特殊化されたエレクトロニック・デバイスを構築すると、必要とするオペレーションを実行できるので好都合である。

50

## 【 0 0 4 9 】

本発明の第2側面によるコード変換システムはソース文字コード化からの文字をターゲット文字コード化に変換するとき方向を考慮に入れる。このコード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスにマッピングする機能を備えている。変換される文字の方向を判断または解明することにより、コード変換システムは判断または解明された文字の方向を利用してターゲット文字コード化への正しいマッピングが得られることを保証する。従って、本発明によれば、ソース文字列内の文字の方向が変化するときでも正しいコード変換が達成される。

## 【 0 0 5 0 】

本発明がアラビアまたはヘブライ・ベースの文字セットが使用されるときに特に便利であるのは、これらの文字が右から左への方向をもっているからである。本発明によれば、どちらかの方向を取り扱うときの柔軟性が得られるだけでなく方向を途中で変更できるという機能も得られる。方向を途中で変更できることはアラビア文字および/またはヘブライ文字が左から右への方向が最も普通である他の文字セットと一緒に使用されるような場合に重要である。異なる方向をもつ文字の使用例として、スペース文字がある。ユニコードでは、コード化は1つだけであり、固有の方向をもっていない。他方、MacArabicでは、左から右へのスペース文字と右から左へのスペース文字の2種類がある。

## 【 0 0 5 1 】

変換されるソース文字の方向を判断または解明することは、方向解明手法によって達成されるが、この手法についてはソース文字を異なるコード化のターゲット文字に変換するコード変換システムと関連づけて以下で詳しく説明する。

## 【 0 0 5 2 】

本発明の第3側面によるコード変換システムはソース文字コード化からの文字をターゲット文字コード化に変換するときコンテキストを考慮に入れる。このコード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかにマッピングする機能を備えている。文字のコンテキストを判断することにより、コード変換システムは判断された文字のコンテキストを利用してターゲット・コード化への正しいマッピングが得られることを保証する。従って、本発明によれば、文字のコンテキストがターゲット・コード化へのマッピングに影響するときでも正しいコード変換が達成される。

## 【 0 0 5 3 】

ある種の文字セットでは、文字が使用されているコンテキストによって文字が異なる表現形態をもつことになる。これらの文字の表現形態は絵文字(glyphs)である。異なる表現形態は表示されるとき異なって現れる。文字が使用されているコンテキストによって表現形態が決まる。アラビア文字はそのコンテキストに応じて異なる表現形態をもつ言語である。

## 【 0 0 5 4 】

アラビア・スクリプトはアラビア語を書くために使用される。アラビア・スクリプトはその印刷形体においても筆写体(カーシブ)である。その結果として、同じ文字はそれが隣の文字とどのように結合されるかに応じて異なった形体で書かれることがある。このような文字の例として、アラビア文字“HEH”(ユニコード文字ではu0647とコード化される)がある。例えば、ユニコードでのアラビア文字はDOSアラビア文字の4つの異なる表現形態の1つをコンテキストに応じてマッピングされる場合がある。その結果、本発明によるコード変換システムはユニコード・ストリーム内のアラビア文字のコンテキストを判断し、正しいマッピングが得られるように動作する。

## 【 0 0 5 5 】

変換されるソース文字のコンテキストの判断はコンテキスト処理によって達成されるが、このコンテキスト処理はソース文字を異なるコード化のターゲット文字に変換するコード変換システムと関連づけて以下で詳しく説明する。

## 【 0 0 5 6 】

本発明の第2および第3側面によれば、コード変換システムはソース文字を異なるコード化のターゲット文字に変換する。本発明によるコード変換システムはラウンドトリップ忠実度を提供すると共に、その結果の文字コードが他のプラットフォームと互換性をもつことを保証する。コード変換システムは単一ソース文字または文字シーケンスを単一ターゲット文字またはターゲット文字シーケンスのどちらかにマッピングする機能を備えている。ラウンドトリップ忠実度によると、ソース・テキストはターゲット・テキストに変換し、そのあとオリジナル・ソース・テキストに戻るように再変換することができる。その結果の文字コードと他のプラットフォームとの互換性は標準ターゲット文字の使用を最大限にし、プライベート文字の使用を最小限にすることによって保証される。コード変換システムは他の文字セットからユニコード文字に変換し、ユニコード文字を他の文字セットに変換するときに利用すると、特に便利である。ユニコード文字シーケンスをターゲット文字セット内の単一文字にマッピングすることは従来にはなかったものである。

【0057】

コード変換システムはコンピュータ・システムや他のエレクトロニック・デバイスにしてこれらのコード変換オペレーションを行うことができる。このコンピュータ・システムは必要とする目的だけに使用するように構築することも、コンピュータ・プログラムに従って動作する汎用コンピュータにすることも可能である。以下に説明する処理はどのコンピュータ・システムや他のエレクトロニック・デバイスにも適用可能である。具体的には、種々の汎用コンピューティング・マシンは以下に教示する説明に従って書かれたソフトウェアと共に使用することができるが、もっと特殊化されたエレクトロニック・デバイスを構築すると、必要とするオペレーションを実行できるので好都合である。

【0058】

本発明の第4側面では、切り捨て処理手法は変換のために受信されたソース文字列が、そのソース文字列が変換のためにソース文字列を格納している入力バッファの長さを越えたときでも、異なるターゲット文字コード化に正確に変換されることを保証する。切り捨て処理手法は入力バッファに置かれているソース文字列の一部を切り捨てて、切り捨てられた部分がソース文字列内の後続文字に影響されることなくターゲット文字コード化に変換されるように作用する。当然のことであるが、入力バッファに置かれていて、切り捨てられた後のソース文字列の部分（切り捨てられた部分）を最大限にすると、コード変換を効率よく実行することができる。

【0059】

本発明による切り捨て処理手法はソース文字を異なるコード化のターゲット文字に変換するコード変換システムと関連づけて以下で詳しく説明する。

【0060】

図1は本発明による代表的なコンピュータ・システム100を示すブロック図である。コンピュータ・システム100は中央処理ユニット(CPU)102を装備し、このCPUは双方向にランダムアクセス・メモリ(RAM)102と、単方向にリードオンリ・メモリ(ROM)106と結合されている。代表例として、RAM104はプログラミング命令とデータを格納し、この中には現在CPU102上で実行中のプロセスに対する他のデータや命令のほかに、以下で説明するテーブルが含まれている。代表例として、ROM106はコンピュータ・システムがその機能を実行するために使用する基本的操作命令、データおよびオブジェクトを格納している。そのほかに、ハードディスク、CD ROM、磁気光学(フロッピカル)ドライブ、テープ・ドライブなどの大量記憶装置108は双方向にCPU102と結合されている。大量記憶装置108は一般的にCPUによって活発に使用されることがない追加のプログラミング命令、データおよびテキスト・オブジェクトを収容しているが、アドレス空間は例えば仮想メモリなどのためにCPUによってアクセス可能である。上述したコンピュータの各々はさらに入出力ソース110を含んでおり、その代表例として、キーボードやポインタ・デバイス(例えば、マウスやスタイラス)などの入力メディアがある。コンピュータ・システム100はデータと命令を転送できるネットワーク接続112を含むこともできる。追加の大量記憶装置(図示せず)をネッ



トワーク接続 112 を経由して CPU 102 に接続することも可能である。コンピュータ・システム 100 はさらに、コンピュータ・システム 100 によって生成または表示されたテキストおよびイメージ（画像など）を見るための表示スクリーン 114 を含んでいる。

#### 【0061】

CPU 102 はオペレーティング・システム（図示せず）と一緒にあって、コンピュータ・コードを実行するように動作する。コンピュータ・コードは RAM 104、ROM 106、または大量記憶装置 108 に置いておくことができる。また、コンピュータ・コードはポータブル・プログラム・メディア 116 に置いておき、必要時にコンピュータ・システム 100 にロードまたはインストールすることも可能である。ポータブル・プログラム・メディア 116 の例としては、CD-ROM、PC カード・デバイス、RAM デバイス、フロッピディスク、磁気テープなどがある。

10

#### 【0062】

##### I. 定義

1. コード・ポイント：コード・ポイントとは、特定のコード化におけるビット・パターンである。通常、ビット・パターンは 1 または 2 バイト以上の長さになっている。ユニコードのコード・ポイントは常に 16 ビットまたは 2 バイトである。

#### 【0063】

2. コード化：コード化とは、文字セットとコード・ポイントのセットとを 1 対 1 で対応づけたもの（マッピング）である。例えば、ASCII コード化は a - z、A - Z、および 0 - 9 を含むセットをコード・ポイント x 00 - x 7 F に対応づけている。

20

#### 【0064】

3. テキスト要素：テキスト要素は特定のオペレーションで 1 つの単位として扱われる 1 つまたは 2 つ以上のコード・ポイントのシーケンスである。例えば、LATIN CAPITAL LETTER U とその後に続く NON-SPACING DIAERESIS は本発明によるコード変換オペレーションの対象となるテキスト要素である（例えば、この例では、2 つの隣り合う文字）。

#### 【0065】

4. 絵文字：文字を可視的に表現した表示形態。例えば、イタリックの “a” とローマ字の “a” は同じ基礎文字を表現する 2 つの異なる絵文字である。

30

#### 【0066】

5. 表現形態：表現形態はコンテキストに応じてその可視形体を変化させる絵文字である。ある種のコード化はコンテキストから独立している抽象文字だけをマッピングするのに対し、他のコード化は表現形体だけをマッピングする。例えば、“fi” のような合字は、LATIN CAPITAL LETTER F とその後に続く LATIN CAPITAL LETTER I の文字シーケンスの表現形態である。

#### 【0067】

6. フォールバック：フォールバックはソース・コードに正確に等価ではないが、オリジナルの情報の一部を残しているターゲット・コード化内の 1 つまたは 2 つ以上のコード・ポイントのシーケンスである。例えば、(C) は © を表す可能なフォールバックである。

40

#### 【0068】

7. デフォルト：デフォルトはターゲット・コード化内にソース・コード・ポイントに類似したものがないとき使用されるターゲット・コード化内の 1 つまたは 2 つ以上のコード・ポイントのシーケンスである。

50

## 【 0 0 6 9 】

## 11. ユニコード・コンバータ

本発明による一般的変換手法はソース文字を異なるコード化のターゲット文字に変換する。好ましくは、ソース文字またはターゲット文字のどちらかはユニコード文字になっている。

## 【 0 0 7 0 】

ユニコード標準は開発されたいくつかの文字コード化を汎用国際的文字コード化標準に統一化したものである。図2はユニコード文字コード化のフォーマットを示す図である。具体的には、ユニコード標準に用意されているコードは、図2に示すフォーマット200で示すように16ビット幅になっている。本明細書では、ユニコード文字は頭にuを付けた16進数で表され(例えば、u001)、他のコード化内の文字は頭にxを付けた16進数で表されている(例えば、x41は1バイト文字を表し、x8140は2バイト文字を表している)。

10

## 【 0 0 7 1 】

図3は、ソース文字列302を受信し、ターゲット文字列304を出力する本発明による基本的ユニコード・コード変換システム300を示すブロック図である。ユニコード・コード変換システム300はソース文字列302の文字を、ターゲット・ストリーム内にあって、文字コード化がソース文字列で使用されたコード化と異なっている1つまたは複数の文字に変換するように動作する。好ましくは、ユニコード・コード変換システム300はユニコードから異なるターゲット・コード化に変換するか(ユニコードから)、あるいは異なるソース・コードから変換する(ユニコードへ)。

20

## 【 0 0 7 2 】

図4は本発明によるユニコード・コード変換システム400の実施例を示すブロック図である。本発明の第4側面では、トランケータ407とバッファ405が存在し、適用されているが、第1、第2または第3側面では、これらは必ずしも存在する必要がないので、必ずしも適用されるとは限らない。

## 【 0 0 7 3 】

本発明の第1、第2および第3側面によれば、ユニコード・コード変換システム400はユニコード文字列404を受信し、ターゲット文字列406を出力するユニコードからのコンバータ402を含んでいる。ユニコードからのコンバータ402は本発明に従ってコード変換プロセスを実行する。その際に、ユニコードからのコンバータ402はスキャナ408とやりとりする。スキャナ408はスキャナ・テーブル408と一緒に使用されて、ユニコード文字列をスキャンしテキスト要素を特定する。そのあと、ユニコードからのコンバータ402はルックアップ・ハンドラ412を使用して、スキャナ408によって特定されたテキスト要素のターゲット・コード化内の1つまたは複数の文字を調べる。ルックアップ・ハンドラ412はマッピング・テーブル414を使用してテキスト要素のターゲット・コード化内の1つまたは複数の文字を取得する。さらに、ユニコードからのコンバータ402はフォールバック・ハンドラ416を使用することも可能である。フォールバック・ハンドラ416はマッピング・テーブル414と一緒に使用され、ルックアップ・ハンドラ412がテキスト要素のターゲット・コード化内の1つまたは複数の文字を特定できなかった場合に、ターゲット・コード化内にあって、テキスト要素のフォールバック・マッピングとして使用できる1つまたは複数の文字を特定するように動作する。状態管理機構(state administrator)418は変換の現在状態に関する情報を維持し、あるいはストアしている。この情報の例としては、シメトリック・スワッピングのコンテキスト、方向および状態がある。

30

40

## 【 0 0 7 4 】

本発明の第4側面によれば、ユニコード・コード変換システム400はユニコード文字列404を受信し、ターゲット文字列406を出力するユニコードからのコンバータ402を含んでいる。ユニコード文字列はバッファ405にストアされる。トランケータ407はバッファ405にストアされたユニコード文字列404の長さを切り捨てて、正確な

50

変換が行われることを保証する。

#### 【 0 0 7 5 】

ユニコードからのコンバータ 4 0 2 はコード変換の全体的オペレーションを制御する。その際に、ユニコードからのコンバータ 4 0 2 はスキャナ 4 0 8 とやりとりする。スキャナ 4 0 8 はスキャナ・テーブル 4 0 8 と一緒に使用されて、切り捨てられたユニコード文字列 4 0 4 (トランケータ 4 0 7 から与えられたもの) をスキャンしテキスト要素を特定する。そのあと、ユニコードからのコンバータ 4 0 2 はルックアップ・ハンドラ 4 1 2 を使用して、スキャナ 4 0 8 によって特定されたテキスト要素のターゲット・コード化内の 1 つまたは複数の文字を調べる。ルックアップ・ハンドラ 4 1 2 はマッピング・テーブル 4 1 4 を使用してテキスト要素のターゲット・コード化内の 1 つまたは複数の文字を取得する。さらに、ユニコードからのコンバータ 4 0 2 はフォールバック・ハンドラ 4 1 6 を使用することも可能である。フォールバック・ハンドラ 4 1 6 はマッピング・テーブル 4 1 4 と一緒に使用され、ルックアップ・ハンドラ 4 1 2 がテキスト要素のターゲット・コード化内の 1 つまたは複数の文字を特定できなかった場合に、ターゲット・コード化内にあって、テキスト要素のフォールバック・マッピングとして使用できる 1 つまたは複数の文字を特定するように動作する。状態管理機構 (state administrator) 4 1 8 は変換の現在状態に関する情報を維持し、あるいはストアしている。この情報の例としては、シメトリック・スワッピングのコンテキスト、方向および状態がある。この種の情報は、切り捨てられたユニコード文字列がブロックの終わりで終わっていない非ブロック区切り変換を行うとき必要になるものである。そのような場合には、変換の現在状態に関する情報をストアしておく、コード変換プロセスは入力文字列 4 0 4 がバッファ 4 0 5 のサイズよりも大きいときでも、正確なコード変換を行うことができる。

#### 【 0 0 7 6 】

##### A . スキャナおよびスキャナ・テーブル

スキャナ 4 0 8 はスキャナ・テーブル 4 1 0 と一緒に使用されて、ユニコード文字列 4 0 4 をスキャンし、ルックアップ・ハンドラ 4 1 2 が必要とする次のテキスト要素と追加情報を戻す。追加情報には、方向情報、コンテキスト情報、および種々の状態インジケータの 1 つまたは 2 つ以上が含まれている。以下では、スキャナ 4 0 8 の一般的オペレーションについて説明する。スキャナ 4 0 8 は入力ユニコード文字列 4 0 4 の文字をスキャンして行く。ターゲット・コード化のために方向情報が必要であれば、テキスト要素内の各文字ごとに文字方向が取得される。また、ターゲット・コード化のために文字コンテキスト情報が必要であれば、テキスト要素内の各文字ごとに文字コンテキスト情報が取得される。そのあと、スキャナ 4 0 8 が文字の各々をスキャンしていくとき、スキャナ 4 0 8 はスキャナ・テーブル 4 1 0 に入っている情報に従って文字に対するアクションをとる。スキャナ 4 0 8 がどのようなアクションをとるかは、状態と文字クラスに基づいて判断される。スキャナ 4 0 8 がとることができるアクションとしては、現在の文字にマークを付けること、シメトリック・スワッピング・ビットをセットまたはクリアすること、テキスト要素のコンテキスト形式を記録すること、テキスト要素が再配列を必要とすることを示すフラグをセットすること、テキスト要素の終わりを示すこと、などがある。シメトリック・スワッピング・ビット、コンテキストおよび方向はスキャナの状態に関する情報として状態管理機構 4 1 8 によってセーブされる。戻す前に、スキャナ 4 0 8 はテキスト要素のコンテキスト情報をセーブしておく。スキャナ 4 0 8 はテキスト要素 (入力文字列内の各テキスト要素) とその属性を返却する。属性には次のものがある。方向、クラス、優先順位、シメトリック・スワッピング状態、サブセットおよびコンテキストである。スキャナ 4 0 8 がテキスト要素を判断したあと、文字を標準形順序に再配列する必要がある場合がある。1 つの例として、テキスト要素内の文字の再配列はユニコードで定義されている標準形順序になっていない非スペース・マークがテキスト要素に含まれているときに行われる。

#### 【 0 0 7 7 】

好ましくは、スキャナ 4 0 8 はスキャナ・テーブル 4 1 0 と一緒に、並列に動作するべ

アのステート・マシンとして実現されている。第1ステート・マシンは文字方向を解明し、第2ステート・マシンは該当する場合には、テキスト要素と文字形式コンテキスト情報を計算し、シメトリック・スワッピング状態も記録しておく。別々になった2つのステート・マシンを使用すると、ユニコード・コード変換システム400の設計と保守が容易化される。第1および第2ステート・マシンは状態とクラスによってインデックスされる2次元配列（またはテーブル）として実現することができる。スキャナ408がとるべきアクションが文字方向によって決まる場合には、ステート・マシン・エントリは各方向についてスキャナ408がとるべき該当アクションを収めている別のテーブルまでのインデックスである。

#### 【0078】

スキャナ408の機能は入力ユニコード文字列404をテキスト要素に変換し、テキスト要素とその属性を戻すことである。スキャナ408はテキスト要素のある種の特性をセーブしておく必要があり、そうすれば、ターゲット・コード化で正しく変換されることになる。すなわち、特性には、方向、コンテキストおよびシメトリック・スワッピング状態がある。しかし、スキャナ408は、そのオペレーションが特定のターゲット・コード化から独立しているので、どのようなターゲット・コード化であるかを知っている必要はない。それにもかかわらず、ユニコード変換システム400は、テキスト要素の定義（つまり、チャンク行動）が、スキャナ・テーブル410を変更するだけでターゲット・コード化と共に変化できるように実現されていることが好ましい。

#### 【0079】

文字の方向性は文字を表示するために使用される。例えば、アラビアまたはヘブライ文字が表示スクリーンに表示されるとき、これらは右から左への順序になっている。大部分のユニコード文字は黙示的方向をもっている（ユニコード・バージョン1.0の407ページ（セクション4.6）および611ページ（付録A）を参照）。ユニコード標準に用意されている暗黙的方向クラスとその値には次のものがある。左から右へ（0）、右から左へ（1）、ヨーロッパ数字（2）、ヨーロッパ数字セパレータ（3）、ヨーロッパ数字ターミネータ（4）、アラビア数字（5）、共通数字セパレータ（6）、ブロック・セパレータ（7）、セグメント・セパレータ（8）、ホワイトスペース（9）およびその他の数字である。スキャナ408はテキスト要素の文字の方向クラスを調べる。次に、その方向クラスはテキスト要素の方向を解明するために使用される。また、方向性をオーバーライドまたは埋め込ませる特殊なユニコード文字もある。これらの特殊なユニコード文字はスキャナ408によって単一文字テキスト要素として扱われる。

#### 【0080】

テキスト要素を形成するときスキャナ408が従う基本的ルールがいくつかある。ベース・ルールでは、適用されるルールがなければ、テキスト要素は単一ユニコード文字とされる。別のルールでは、ベース文字に続く非スペースまたは結合マークはそのベース文字と一緒に単一テキスト要素として分類される。さらに別のルールでは、シンボル（つまり、朝鮮ハングル・ジャモス文字）、合字または表意文字に関連する文字が見つかったとき、これらはテキスト要素に結合される。さらに別のルールでは、フラクション・スラッシュの両側が1つまたは複数の10進数のシーケンスで囲まれていると、これらは数字フラクション・テキスト要素として結合される。

#### 【0081】

以下では、非スペースまたは結合マークに関するルールについて詳しく説明する。ユニコード標準によれば、非スペース・マークはベース文字のあとに置かれている。従って、ベース文字のあとに置かれた非スペース文字はベース文字を含むテキスト要素の一部となる（ユニコード標準バージョン1.0、403ページ（セクション4.3）を参照）。例えば、単一の非スペース文字の後に非スペース文字でない文字が続いているときは、その非スペース文字は直前の文字と一緒にテキスト要素として結合される。その場合、テキスト要素の長さは2であり、テキスト要素の属性はベース文字によって定義されている。前に置かれた文字がなければ、非スペース文字は単一テキスト要素として渡されるだけであ

10

20

30

40

50

る。複数の非スペース文字もこのようにして結合することができる。

#### 【 0 0 8 2 】

以下では、朝鮮ハングル・ジャモス文字について詳しく説明する。各ハングル文字は暗黙値を持ち、これは次のクラスの1つになっている。Choseong（初期）、Jungseong（中間）またはJongseong（最終）である。ユニコード標準バージョン1.1（セクション5）はこれらの文字のコードと許容される組み合わせをリストしている。スキャナ408はこれらの文字の許容される組み合わせに従って朝鮮ジャモス文字をグループ化する。入力が許容されない組み合わせのときは、スキャナ408は文字を単一テキスト要素として返却する。前述したように、ハングル分節のあとに結合マークが付いているときは、その結合マークはハングル分節のテキスト要素内に挿入される。

10

#### 【 0 0 8 3 】

次に、数字フラクシオンに関するルールについて詳しく説明する。スキャナ408は最初にフラクシオン・スラッシュ数字の各文字を、それらが単一文字テキスト要素であるものとして取り扱う。しかし、完全なフラクシオン・スラッシュ・シーケンスが見つかったときは、スキャナ408はそのシーケンスに関連する文字を単一テキスト要素になるように結合する。ディジットが結合マークと一緒に見つかったときは、そのディジットと結合マークはフラクシオン・スラッシュの一部となることができないが、ディジットと結合マークと一緒にテキスト要素を形成することができる。

#### 【 0 0 8 4 】

非スペース文字を除き、すべてのアラビア文字は単一テキスト要素としてスキャナ408を通過する。アラビア形式のシェーピング状態文字も、単一テキスト要素としてスキャナを通過する。方向フォーマット・コードは単一テキスト要素としてスキャナ408を通過する。

20

#### 【 0 0 8 5 】

B. ルックアップ・ハンドラ、マッピング・テーブルおよびフォールバック・ハンドラ

マッピング・テーブル414は1つまたは複数のユニコード文字の入力シーケンスをターゲット・コード化における1つまたは複数の出力シーケンスと突き合わせるためにルックアップ・ハンドラ412によって使用される。ユニコード・シーケンス（つまり、テキスト要素）自体のほかに、入力シーケンスに関するある種の追加情報が得られ（例えば、方向、コンテキスト、シメトリック・スワッピング状態、垂直形式要求、フォールバック要求、許容範囲、変種）、ある種のテーブルはこの情報を利用している。好ましくは、マッピング・テーブル414はフォールバック・ハンドラ416が必要とするデータもストアしているが、別のテーブルを用意してフォールバック・ハンドラ416に使用させることも可能である。

30

#### 【 0 0 8 6 】

図5はユニコード・コード変換システム400のマッピング・テーブル414の好ましい配列を示す概略図である。マッピング・テーブル414は好ましくはヘッダ部分500を含み、マッピング・テーブル414のデータのセグメントはテキスト要素内の文字数に基づいて分割されている。ヘッダ部分500の内容は以下で詳しく説明する。図5に示すマッピング・テーブル414は1からN文字までのテキスト要素のコード化をサポートしている。ルックアップ・ハンドラ412が1文字テキスト要素のターゲット・コード化を探すためにマッピング・テーブル414をサーチするとき、マッピング・テーブル414のセグメント502が使用される。同様に、テキスト要素が2文字幅であれば、セグメント504が使用され、テキスト要素がN文字幅であれば、セグメント506が使用される。図4と図5は単一マッピング・テーブル414を示しているが、ユニコード・コード変換システム400は複数の異なるマッピング・テーブル414を使用する。つまり、各ターゲット文字セットごとに1つのマッピング・テーブルを使用する。各マッピング・テーブルは複数のサブテーブルを含んでいる。

40

#### 【 0 0 8 7 】

マッピング・テーブル414はサイズと全体的変換速度要件を考慮に入れて設計されて

50

いる。マッピング・テーブル 4 1 4 はルックアップ時間を重大に低下させることなく可能な限り小さくしておくべきであり、ルックアップ時間はテーブル・サイズを大幅に大きくすることなく可能な限り高速にしておくべきである。ユニコード・コード変換システム 4 0 0 は複数のテーブル・フォーマットをサポートしているので、各サブテーブルごとに異なるフォーマットにすることが可能である。そのようにすると、速度とサイズのトレードオフを特定のテーブルに合わせて必要時に調整することができる。好ましくは、マッピング・テーブル 4 1 4 の設計は単一ユニコード文字からターゲット・コード化内の単一文字にマッピングするのを可能な限り高速化するものでなければならないが、これは最も普通の使い方であるからである。

#### 【 0 0 8 8 】

マッピング・テーブル 4 1 4 の設計は、テーブルがフォールバック・ハンドラ 4 1 6 の要求の少なくとも一部をサポートし、複数のマッピング許容範囲をサポートし、複数のターゲット文字セットの変種をサポートするようになっている。テーブル・フォーマットは 1 つまたは複数のユニコード文字シーケンスをゼロまたは 1 個以上の文字の出力シーケンスにマッピングすることもできる。マッピング・テーブル 4 1 4 は単一入力シーケンスに対して起こり得る複数の出力シーケンスを指定することが可能であり、特定の出力シーケンスは方向、コンテキストおよびシメトリック・スワッピング状態などの属性によって判断される。本発明の第 3 側面に関しては、中心となるのは出力シーケンスの 1 つをコンテキストに基づいて選択することである。マッピング・テーブルは容易に拡張することができるので、ユニコード・コード変換システム 4 0 0 のコード化の振舞を容易にカスタマイズすることができる。

#### 【 0 0 8 9 】

マッピング・テーブル 4 1 4 が必要とする情報には次のものがある。スキャナ 4 0 8 からのテキスト要素（つまり、結合マークを標準形順序にして変換される入力文字シーケンス）、垂直形式が使用可能なときに水平形式の代わりに使用するかどうか、解明されたテキスト要素の方向、テキスト要素の文字形式コンテキスト情報（初期、中間、最終、または隔離）、シメトリック・スワッピングの現在状態、どのレベルのルックアップを起動させるかの情報（許容レベル（厳格または緩和）とフォールバック（オンまたはオフ））、および特定のコード化変種の識別子である。ルックアップ・レベルに関する情報はコールまたはユニコード・コード変換システム 4 0 0 をコールするアプリケーション・プログラムから与えられる。方向およびコンテキストが重要でない言語または文字では、解明された方向とコンテキスト情報は必要でない。

#### 【 0 0 9 0 】

変種の定義、ユニコード・シーケンスからターゲット・シーケンスへの実際のマッピング、およびこれらにアクセスするために使用されるテーブル・フォーマットはマッピング・テーブル 4 1 4 の設計によって変更可能である。従って、正確性および若干であるが、パフォーマンスとサイズとのトレードオフはマッピング・テーブル 4 1 4 の設計に大きく依存している。好ましいことは、マッピング・テーブル 4 1 4 が厳格および緩和マッピング、フォールバック・マッピング、およびデフォルト・マッピングをサポートすることである。

#### 【 0 0 9 1 】

厳格マッピングはラウンドトリップ忠実度が保証されるコード変換である。ユニコードからターゲット文字セットへの厳格マッピングはその文字セットからユニコードへのマッピングとは逆である。ユニコードからターゲット文字セットへの緩和マッピングはターゲット文字セット内の文字の定義または隔離された用途の範囲に属する追加マッピングである。緩和マッピングは正しくマッピングされるように見えるが、若干のあいまいさがある。例えば、多くの文字セットでは、単一文字は、明示的定義、あいまいな定義、または確立された用法のいずれかによって複数の意味を持つことがある。例えば、Shift-JIS 文字 x 8 1 6 1 は 2 つの意味を持つように規定されている。すなわち、「2 重垂直線」と「平行」である。これらの意味の各々は異なるユニコード文字 u 2 0 1 6 「2 重垂直

10

20

30

40

50

線」とu 2 2 2 5「に平行」に対応している。Shift-JISからユニコードにマッピングするときは、コード変換システムはこれらのユニコード文字の1つ、つまり、2重垂直線を選択しなければならない。ユニコードからShift-JISにマッピングするときは、コード変換システムは両方のユニコード文字を同一Shift-JIS文字にマッピングすることができ、そのようにするのが通常である。しかし、これらのユニコードからのマッピングの1つだけはユニコードへのマッピングとは逆になっている。

#### 【0092】

厳格マッピングと緩和マッピングの比較例

- ・ ユニコードu 0 0 0 DがASCII x 0 D「キャリッジリターン」に厳格にマッピングされていれば、ユニコードu 2 0 2 9「パラグラフ・セパレータ」はASCII x 0 Dにゆるやかにマッピングすることができる。

10

- ・ ユニコードu 0 0 2 D「ハイフン・マイナス」がASCII x 2 D「ハイフン・マイナス」に厳格にマッピングされていれば、ユニコードu 2 0 1 0「ハイフン」とu 2 2 1 2「マイナス記号」はASCII x 2 Dにゆるやかにマッピングすることができる。

- ・ ユニコードu 0 0 E 0「グループ付きのラテン小文字A」がISO 8859-1 x E 0「アクサングラフ付きの小文字a」に厳格にマッピングされていれば、2文字のユニコード・シーケンスu 0 0 6 1 + u 0 3 0 0「ラテン小文字A」+「結合アクサングラフ」はISO 8859-1 x E 0にゆるやかにマッピングすることができる。

- ・ Shift-JISは半幅文字と全幅文字を区別しているので、Shift-JISの緩和マッピングもこれらを区別しておかなければならない。つまり、ユニコードu F F 4 0「全幅アクサングラフ」はShift-JIS x 8 1 4 D「アクサングラフ(全幅)」に厳格にマッピングされており、これはShift-JIS x 6 0「アクサングラフ(半幅)」と区別されている。ユニコード・シーケンスu 3 0 0 0 + u 0 3 0 0「表意文字スペース」+「結合アクサングラフ」はShift-JIS x 8 1 4 Dにゆるやかにマッピングすることができる。しかし、ユニコード・シーケンスu 0 0 2 0 + u 0 3 0 0「スペース」+「結合アクサングラフ」はShift-JIS x 8 1 4 Dにゆるやかにマッピングしてはならない。これはShift-JIS x 6 0にマッピングされるべきである。

20

#### 【0093】

本発明の第1側面によれば、ユニコードからある種の他の文字にマッピングし再び元に戻すマッピング(ラウンドトリップ・マッピング)は、他の文字への厳格なマッピングが存在しているユニコード文字だけを使用するとき可能である。

30

#### 【0094】

さらに、本発明の第1から第5までの側面によれば、フォールバック・マッピングはユニコード文字の意味または同一性を保存していないユニコードからのマッピングである。つまり、これらはユニコード文字(または文字シーケンス)をその定義または用途がユニコード文字の意味または用途を含んでいないターゲット文字セット内の文字(または文字シーケンス)にマッピングする。これに対して、フォールバック・マッピングによると、ユニコード文字(または文字シーケンス)に最も近く対応しているターゲット・コード化内の文字(または文字シーケンス)が得られる。

40

#### 【0095】

フォールバック・マッピングの例

- ・ ユニコード文字u 0 3 0 0「結合アクサングラフ」はフォールバック・マッピングとしてASCII x 6 0「アクサングラフ[スペース]」にマッピングすることができる。違いは、ユニコード文字が結合マーク(非スペース)であるのに対し、ASCII文字はスペース・マークであることである。

- ・ ユニコード文字u 0 1 C 0「ラテン文字デンタル・クリック」はフォールバック・マッピングとしてASCII x 7 C「垂直線」にマッピングすることが可能である。

- ・ ユニコード文字u 2 0 0 1「EM QUAD」はフォールバック・マッピングとしてASCII x 2 0「スペース」にマッピングすることが可能である。

50

## 【 0 0 9 6 】

従って、上記の例に示すように、フォールバック・マッピングはユニコード文字（またはシーケンス）にグラフィックが近似しているターゲット文字（またはシーケンス）を生成するために使用される。

## 【 0 0 9 7 】

パフォーマンス上の理由により（つまり、コード化をマッピング・テーブルから得るときの速度）、マッピング・テーブル 4 1 4 までのインデックスにはいくつかのフォーマットがある。可能とされるフォーマットには、セグメント・フォーマット、リスト・フォーマット、範囲フォーマット、チェーン・フォーマットがある。文字シーケンスの長さが異なるごとに、別々のインデックスを設けることが好ましい。その結果として、セグメント 5 0 2、5 0 4、5 0 6 の各々に関連するインデックスは異なるフォーマットにすることができ、各インデックスの先頭の情報はそのフォーマットを指定している。フォーマットに関係なく、各インデックスは最終的には、入力シーケンスを直接に出力シーケンスにマッピングするか、あるいは出力シーケンスが長ければ、対応する出力シーケンスのロケーションを指定しているオフセットにマッピングする。

## 【 0 0 9 8 】

マッピング・テーブル 4 1 4 までのインデックスのチェーン・フォーマットについては、以下で詳しく説明する。チェーン・フォーマットでは、セクションの先頭がチェックされ、それがチェーン・フォーマット・テーブルのチェーン・ヘッダであるか、他のフォーマットであるかが判断される。チェーン・フォーマットは複数のインデックス・テーブルのチェーンを指定し、各々は異なるフォーマットになっている場合がある。必要とするマッピングが最初のインデックス・テーブルに見つからなければ、ルックアップ・ハンドラ 4 1 2 は 2 番目のインデックス・テーブルを調べ、以下同様である。チェーン・フォーマットは、例えば、あるインデックス・フォーマット（これは空間的および/または時間的に効率的である）が入力シーケンスのすべてではないが、大部分をマッピングすることができるのに対し、別の非効率なインデックス・フォーマットは少数の残りのシーケンスを処理できるときに利用すると、便利である。チェーン・メカニズムがないと、インデックス・シーケンスのすべてに非効率のフォーマットを使用しなければならないことになる。チェーン・フォーマットは変種と許容レベルが異なるたびに、異なるサブテーブルが必要になるときにも便利である。チェーン内の各サブテーブルは、現在取り扱われているマッピング許容範囲と変種に基づいてそれを含めるか、除外させるビット・フラグをもっている。本発明の第 2、第 3 および第 4 側面によれば、ルックアップ・ハンドラ 4 1 2 がターゲット・コード化のマッピング・テーブル 4 1 4 をサーチするときは、含められたサブテーブルだけが考慮される。

## 【 0 0 9 9 】

さらに、本発明の第 1 から第 4 までの側面によれば、各サブテーブルに関連するこれらのビット・フラグはサブテーブル・マスクを形成している。また、呼び出し側要求（例えば、コード化の変種と許容範囲）と判断された属性（例えば、解明された方向とコンテキスト）は選択マスクを形成している。サブテーブル・マスクと選択マスク内のビット割当ては同じである。従って、特定のサブテーブルを含めるかどうかの判断はサブテーブル・マスクと選択マスクのビットワイズ AND をとり、そのあと結果をサブテーブル・マスクと比較することによって行われる。ビットワイズ AND の結果がサブテーブルのサブテーブル・マスクと同じであれば、そのサブテーブルが含まれる。そうでなければ、それは含まれない。

## 【 0 1 0 0 】

マッピング・テーブル 4 1 4 のヘッダ 5 0 0 は好ましくは次のものを収めている。

- ・ 一般的識別情報 - フォーマット、長さ、チェックサムおよびバージョン。
- ・ 最小ターゲット文字サイズ（バイト数）（文字サイズともいう）。
- ・ 一般的フラグ（例えば、そのルックアップ・テーブルが方向またはコンテキスト・データを必要としているかどうか）。



- ・ そのテーブルによって処理される最大入力シーケンス長、および 1 からこの最大長までの入力シーケンス長を処理するテーブルを指定しているオフセット / 長さのペアのリスト。
- ・ そのユニコードからのマッピングのデフォルト・フォールバック文字または文字シーケンス。
- ・ そのテーブルによってサポートされる変種のカウントとリスト。各変種ごとに、1 つまたは 2 つ以上の関連ビット・マスクが指定される。単一の変種に複数のビット・マスクがある場合は、属性情報（方向、コンテキスト、および垂直形式の要求）はどのビット・マスクが使用されるかを判断するために使用される。ビット・マスクで「1」にセットされたビットは異なる変種をサポートするために種々のサブテーブルをオンにするために使用される。
- ・ 可能とされる 4 つの許容範囲設定値（厳格 / 緩和、フォールバック・オン / オフ）の各々に関連する追加のビット・マスクのセット。該当する許容範囲レベル・マスクは変種マスクと OR がとられて、サブテーブルを使用可能または使用禁止にするために使用されるビット・マスクを形成する。

10

【0101】

### C. コード変換処理

以下では、ユニコード・コード変換システム 400 の好適実施例によって実行される処理について詳しく説明する。

【0102】

20

図 6 A はユニコード・コード変換システム 400 を利用するアプリケーション（つまり、呼び出し側アプリケーション・プロセスまたはプログラム）によって実行される処理 600 を示すフローチャートである。具体的には、図 6 A はユニコードからの処理に関係しているが、当然に理解されるように、類似のオペレーションを実行して他方の方向に変換することも可能である（ユニコードへの処理）ユニコードからのコンバータ 402 は全体的変換プロセスを制御する。

【0103】

最初に、処理 600 は新しい状態インスタンスと変換のための制御情報を作成し、初期化する（602）。処理 600 はインスタンスを設定するので、複数のスキャン・オペレーションを進行中にし、そのインスタンスによって区別することが可能である。次に、切り捨てが必要であるかどうかの判断 604 が行われる。切り捨てが必要な場合には、トラシケータ 407 が起動される（606）。切り捨てが必要になるのは、入力データ・ストリームが、変換のためのデータを収めている受入れバッファ 405 の容量を越えたときである。切り捨てが必要でなければ、あるいは切り捨てが必要な場合には、ブロック 606 に続いて、ユニコードからのコンバータ 402 がコールされ（608）ユニコード文字列 404 を変換する。ユニコードからのコンバータ 402 の機能は、テキスト要素を取得し、そのためのターゲット・マッピングを調べることである。これについては、以下で詳しく説明する。変換機能から戻ると、処理 600 はユニコードからのコンバータ 402 からターゲット文字列 406 を受け取る。

30

【0104】

40

次に、処理 600 は変換が失敗したかどうかを判断する。変換が失敗していれば、エラーが処理される（614）。他方、変換が正常に行われていれば、変換が完了したかどうかの判断 616 が行われる。ユニコード文字列 404 の文字がターゲット・コード化に変換されたとき変換は完了する。変換が完了していれば、処理 600 は完了し、ターゲット文字列 406 はコード変換を要求したプロセスまたはアプリケーションに利用可能にされる。さらに、処理 600 は状態インスタンスと変換のための制御情報を破棄する（618）。他方、変換がまだ完了していないと判断 616 されたときは、処理 600 は変換が完了するか、エラーが起こるまでブロック 604 ~ 616 を繰り返す。

【0105】

図 6 B は本発明の第 4 側面による切り捨て処理 620 を示すフローチャートである。切

50

り捨て処理 6 2 0 は図 6 A のブロック 6 0 6 で起動され、トランケート 4 0 7 によって実行される。

【 0 1 0 6 】

切り捨て処理 6 2 0 は出力長をゼロ ( 0 ) に初期化する ( 6 2 2 )。出力長はバッファの実効長、つまり、バッファの切り捨てられた長さに一致している。次に、次のテキスト要素が取得される ( 6 2 4 )。次のテキスト要素の取得 ( 6 2 4 ) に関連する処理はスキャナ 4 0 8 によって行われるが、これについては図 9 A ~ 図 9 C を参照して以下で詳しく説明する。次に、テキスト要素がバッファの物理的長さ ( バッファ長 ) を越える可能性があるかどうかの判断 6 2 6 が行われる。テキスト要素がバッファ長を越える可能性がなければ、出力長はそのテキスト要素を含むように更新される ( 6 2 8 )。ここでは、実効長はテキスト要素がバッファ長を越える可能性がない限り、テキスト要素単位で大きくなる。ブロック 6 2 8 に続いて、考慮すべき追加テキストがバッファに残っているかどうかの判断 6 3 0 が行われる。考慮すべき追加テキストがバッファに残っていれば、処理はループしてブロック 6 2 4 に戻り、切り捨て処理 6 2 0 を繰り返す。他方、考慮すべき追加テキストがバッファに残っていないか、あるいはテキスト要素がバッファ長を越える可能性があるとしてブロック 6 2 6 で判断されていれば、出力長が戻される ( 6 3 2 )。戻された出力長 ( 6 3 2 ) はバッファの実効長 ( つまり切り捨てられた長さ ) であるので、バッファは実効的にテキスト要素の終わり ( つまり、バッファ内のテキストの最後のテキスト要素 ) で終了する。

【 0 1 0 7 】

以上のようにして、切り捨て処理はバッファの実効長 ( つまり、切り捨てられた長さ ) を判断する。これはバッファの切り捨てられた部分とも呼ばれる。切り捨ての後バッファに残っている余分のテキストは残余部分と呼ばれる。残余部分は入力ソース文字列の次のバッファ部分にキャリーオーバーされ、その部分と一緒に考慮される。この処理はテキスト要素が正しく判断されることを保証する。

【 0 1 0 8 】

以下は切り捨て処理の使い方の例である。

例 : “ . . . A B C D ` E F G . . . ”

バッファに置かれている部分が “ D ” の直後で終わっていれば、切り捨て処理はバッファに置かれている部分を切り捨てて、切り捨てられた長さが “ C ” のあとで終わるようにする。バッファに置かれた部分を切り捨てる必要があるのは、そのようにしないと、テキスト要素 “ D ” がそのあとに置かれた結合マーク “ ` ” から切り離されることになるからである。切り離されると、テキストはターゲット・コード化に正しく変換されないことになる。残余部分は “ D ` ” であり、次の部分にキャリーオーバーされ、“ D ` E F G . . . ” となる。

【 0 1 0 9 】

図 7 はユニコード・コンバータ処理 7 0 0 を示すフローチャートである。ユニコード・コンバータ処理 7 0 0 は図 6 A のブロック 6 0 8 で実行されるオペレーションと関連づけられている。

【 0 1 1 0 】

ユニコード・コンバータ処理 7 0 0 は判断 7 0 2 から始まる。この判断 7 0 2 は変換すべきテキストがあるかどうかを判断する。変換すべきテキストがなければ、ユニコード・コンバータ処理 7 0 0 は戻るだけである ( または完了する )。他方、変換すべきテキストがあれば ( つまり、ユニコード文字列 4 0 4 が完全に処理されていなければ )、処理 7 0 0 が続けられる。まず、オフセットがオフセット配列に対して更新される ( 7 0 4 )。オフセット配列入力文字列に関連するオフセット ( ポインタ ) の配列であり、フォント変更、ライン中断、言語変更などのある種の変更が、呼び出し側アプリケーションが有意と判断している入力文字列 4 0 4 内のどこで行われたかを示している。オフセット配列の更新 7 0 4 は、異なる文字長に合わせてオフセット ( ポインタ ) を調節することにより行われ

る。例えば、入力ユニコード文字列 4 0 4 のユニコード文字は 2 バイト長であるのに対し、ASCII のターゲット・コード化に関連する文字のサイズは 1 バイト長である。ここでは、オフセット配列の更新 7 0 4 はオフセットがターゲット・コード化内の対応する文字を指すようにオフセットを更新する。実際には、入力文字列内のオフセットはコード化が異なっている出力文字列にマッピングされる。次に、次のテキスト要素が取得される (7 0 6)。スキャナ・テーブル 4 1 0 を用いるスキャナ 4 0 8 はユニコード文字列 4 0 4 からのテキスト要素を判断する。次のテキスト要素の取得 7 0 6 は以下で詳しく説明する。次に、取得されたテキスト要素はマッピング・テーブル 4 1 4 で調べられ (7 0 8)、ターゲット・コード化内のテキスト要素の変換コードが得られる。このルックアップはマッピング・テーブル 4 1 4 を使用してルックアップ・ハンドラ 4 1 2 によって行われる。変換コードのルックアップ 7 0 8 についても以下で詳しく説明する。

10

**【 0 1 1 1 】**

次に、テキスト要素の変換コードが見つかったかどうかの判断 7 1 0 が行われる。変換コードが見つかったときは、ユニコード文字列 4 0 4 とターゲット文字列 4 0 6 の入力位置ポインタと出力位置ポインタがそれぞれ更新される (7 1 2)。入力位置ポインタは入力文字列 4 0 4 がどれだけ変換されたかを示している。出力位置ポインタはターゲット文字列 4 0 6 の長さを示している。ブロック 7 1 2 に続いて、処理 7 0 0 はユニコード・コンバータ処理 7 0 0 の先頭に戻り、変換すべきユニコード文字列 4 0 4 の次のテキスト要素 (存在する場合) が処理できるようにする。

**【 0 1 1 2 】**

20

しかるに、変換コードがマッピング・テーブル 4 1 4 に見つからないと判断 7 1 0 が判断したときは、呼び出し側 (例えば、呼び出し側アプリケーション) がフォールバック処理を要求していたかどうかの判断 7 1 4 が行われる。呼び出し側がフォールバック処理を要求していれば、フォールバック処理が実行される (7 1 6)。フォールバック処理はフォールバック・ハンドラ 4 1 6 によって実行される。これは以下で詳しく説明する。他方、呼び出し側がフォールバック処理を要求していなければ、テキスト要素がルックアップ・ハンドラ 4 1 2 によってターゲット・コード化に変換できなかったのでエラーが通知される (7 1 8)。ブロック 7 1 6 と 7 1 8 に続いて、ユニコード文字列 4 0 4 とターゲット文字列 4 0 6 の入力位置ポインタと出力位置ポインタがそれぞれ更新され (7 0 2)、そのあと処理はユニコード・コンバータ処理 7 0 0 の先頭に戻り、変換すべきユニコード文字列 4 0 4 の次のテキスト要素 (存在する場合) が処理できるようにする。

30

**【 0 1 1 3 】**

図 8 はオフセット更新処理 8 0 0 を示すフローチャートである。オフセット更新処理 8 0 0 はオフセット配列が更新される図 7 のブロック 7 0 4 と関連している。

**【 0 1 1 4 】**

オフセット更新処理 8 0 0 は現在入力位置がオフセット配列にあるかどうかを判断 8 0 2 することから始まる。現在入力位置がオフセット配列にあれば、オフセット配列の内容がブロック 8 0 4 に続いて現在出力位置長さに従って更新される (8 0 4)、オフセット更新処理 8 0 0 から戻る。他方、現在入力位置がオフセット配列にないときは、更新すべきオフセットがないのでオフセット更新処理 8 0 0 から戻る判断 8 0 2 が行われるだけである。

40

**【 0 1 1 5 】**

図 9 A および図 9 B は本発明の第 1、第 2 および第 4 側面による次のテキスト要素処理 9 0 0 を示すフローチャートである。次のテキスト要素処理 9 0 0 は次のテキスト要素を取得するとき図 7 のブロック 7 0 6 で実行されるオペレーションを具体化したものである。好ましくは、次のテキスト要素処理 9 0 0 はスキャナ・テーブル 4 1 0 を使用してスキャナ 4 0 8 によって行われる。

**【 0 1 1 6 】**

次のテキスト要素処理 9 0 0 は状態と再配列フラグを初期化 (9 0 2) することから始まる。次に、マッピング・テーブル 4 1 4 が方向情報を必要としているかどうかの判断 9

50

04が行われる。マッピング・テーブル414が方向情報を必要としていれば、ユニコード文字列404の方向が解明される(906)。本発明の第2側面では、マッピング・テーブル414が方向情報を必要としていれば、ユニコード文字列404の次の入力文字の方向が解明される(906)。方向が解明されると(906)、あるいは方向が必要でない場合には判断904の後で、文字のコンテキストに基づく判断908が行われる。ユニコード文字列404内の文字のコンテキストがコード変換(マッピング)に影響する可能性があるときは、コンテキストが解明される(910)。ブロック910に続いてまたは判断908に続いて、コンテキストが重要でないときは、ユニコード文字がユニコード文字列404から取得される(912)。ここでは、ユニコード文字列404内の次の文字が取得される。次に、取得されたユニコード文字の属性が調べられる(914)。属性ルックアップ914は図11~図13を参照して以下で詳しく説明する。次に、次のテキスト要素処理900のアクションが判断される(916)。アクションの判断916は図14A、図14Bおよび図15を参照して以下で詳しく説明する。

#### 【0117】

次に、アクションが“END”であるかどうかの判断918が行われる。アクションが“END”でなければ、アクションは少なくとも“ADVANCE”である。アクションが“ADVANCE”であるときは、アクションが“MARK”でもあるかどうかの判断920が行われる。アクションが“MARK”でもあるときは、文字がテキスト要素に挿入される(922)。さらに、“MARK”アクションがとられると双方向状態がセーブされる(924)。双方向状態は方向性埋め込みスタックと双方向ステート・マシンの現在の状態を含んでいる。ブロック924に続いてまたはアクションが“MARK”でもないときは判断ブロック920に続いて、アクション修飾子に基づいてスイッチ・オペレーションが実行される。アクション修飾子はアクションの一部であり、“S”、“ISS”、“ASS”などの修飾子を含んでいる。アクションはどのアクション修飾子ももたないこともできる。アクション修飾子が“S”であるときは、再配列フラグがセットされる(928)。再配列フラグ(セットされたときは)はテキスト要素内の文字を再配列する必要がないことを示している。アクション修飾子が“ISS”(つまり、シメトリック・スワッピング禁止)であるときは、スワップ・フラグがオフにセットされる(930)。アクション修飾子が“ASS”(つまり、スワッピング禁止活動化)であるときは、スワップ・フラグがオンにセットされる(932)。スワップ・フラグはシメトリック・スワッピングが必要かどうかを示している。スイッチ926は拡張エリア934を使用して追加のアクション修飾子を含むように容易に適応させることができる。拡張エリア934を使用すると、ユーザはスキナ408の振舞を変更することができる。また、アクション修飾子があれば、次のテキスト要素処理900はアクション修飾子に関連するどのオペレーションも実行しない。アクション修飾子オペレーションに続いて、現在文字インデックスが更新される(936)。現在文字インデックスは次のテキスト要素処理900を実行するときソース文字列をスキャンしていくために使用されるソース文字列を指すポイントである。次に、次のテキスト要素処理900はブロック904から始まるオペレーションを繰り返す。処理はアクションが“END”であると判断918が判断するまでループしてブロック904-936を繰り返す。アクションが“END”であると、判断918は判断ブロック938を実行させる。判断ブロック938は再配列フラグがセットされているかどうかを判断する。再配列フラグがセットされていれば(ブロック928)、テキスト要素内の文字は再配列される(940)。再配列は異種の文字クラスの加重値を提供する優先順位属性を使用して行うことが好ましい。再配列フラグがセットされていない場合にはブロック938に続いて、再配列フラグがセットされている場合にはブロック940に続いて、次のテキスト要素処理900は完了し、ユニコード・コンバータ処理700に戻る。

#### 【0118】

図10A、図10Bおよび図10Cは本発明の第3側面による次のテキスト要素処理900を示すフローチャートである。次のテキスト要素処理900は次のテキスト要素を取

10

20

30

40

50

得するとき図7のブロック706によって実行されるオペレーションを具体化したものである。好ましくは、次のテキスト要素処理900はスキャナ・テーブル410を使用してスキャナ408によって実行される。

#### 【0119】

次のテキスト要素処理900は状態と再配列フラグを初期化(952)することから始まる。次に、マッピング・テーブル414が方向情報を必要としているかどうかの判断954が行われる。マッピング・テーブル414が方向情報を必要としていれば、ユニコード文字列404の次の入力文字の方向が解明される(956)。次に、ユニコード文字がユニコード文字列414から取得される(958)。ここでは、ユニコード文字列404内の次の文字が取得される(958)。次に、取得されたユニコード文字の属性が調べられる(960)。属性ルックアップ960は図11～図13を参照して以下で詳しく説明する。次に、次のテキスト要素処理900のアクションと次の状態が判断される(962)。アクションと次の状態の判断は図14A、図14B、図16Aおよび図16Bを参照して以下で詳しく説明する。

#### 【0120】

次に、アクションが“END”であるかどうかの判断964が行われる。アクションが“END”でなければ、アクションは少なくとも“ADVANCE”である。アクションが“ADVANCE”であるときは、アクションが“MARK”でもあるかどうかの判断966が行われる。アクションが“MARK”でもあるときは、文字がテキスト要素に挿入される(968)。さらに、“MARK”アクションがとられると双方向状態がセーブされる(970)。双方向状態は方向性埋め込みスタックと双方向ステート・マシンの現在状態を含んでいる。ブロック970に続いてまたはアクションが“MARK”でもないときは判断ブロック966に続いて、アクション修飾子に基づいてスイッチ・オペレーション972が実行される。アクション修飾子はアクションの一部であり、“S”、“ISS”、“ASS”などの修飾子を含んでいる。アクションはどのアクション修飾子ももたないこともできる。アクション修飾子が“S”であるときは、再配列フラグがセットされる(974)。再配列フラグ(セットされたときは)はテキスト要素内の文字を再配列する必要がないことを示している。アクション修飾子が“ISS”(つまり、シメトリック・スワッピング禁止)であるときは、スワップ・フラグがオフにセットされる(976)。アクション修飾子が“ASS”(つまり、スワッピング禁止活動化)であるときは、スワップ・フラグがオンにセットされる(980)。スワップ・フラグはシメトリック・スワッピングが必要かどうかを示している。スイッチ972は拡張エリア980を使用して追加のアクション修飾子を含むように容易に適応させることができる。拡張エリア980を使用すると、ユーザはスキャナ408の振舞を変更することができる。また、アクション修飾子がなければ、次のテキスト要素処理900はアクション修飾子に関連するどのオペレーションも実行しない。アクション修飾子オペレーションに続いて、現在文字インデックスが更新される(982)。現在文字インデックスは次のテキスト要素処理900を実行するときソース文字列をスキャンしていくために使用されるソース文字列を指すポインタである。次に、次のテキスト要素処理900はブロック954から始まるオペレーションを繰り返し替える。処理はアクションが“END”であると判断964が判断するまでループしてブロック954～982を繰り返す。アクションが“END”であると、判断964はコンテキスト処理984を実行させる。コンテキスト処理984のあと、判断ブロック986は再配列フラグがセットされているかどうかを判断する。再配列フラグがセットされていれば(ブロック974を参照)、テキスト要素内の文字は再配列される(988)。再配列は異種の文字クラスの加重値を提供する優先順位属性を使用して行うことが好ましい。再配列フラグがセットされていない場合にはブロック986に続いて、再配列フラグがセットされている場合にはブロック988に続いて、次のテキスト要素処理900は完了し、ユニコード・コンバータ処理700に戻る。

#### 【0121】

次に、図10Cを参照してコンテキスト処理984について説明する。コンテキスト処

10

20

30

40

50

理 9 8 4 はここで説明している実施例では、スキャナ・テーブル 4 1 0 を使用するスキャナ 4 0 8 によって実現されている。コンテキスト処理 9 8 4 は判断 9 8 5 から始まり、アクションが End Output X n (つまり、単独コンテキスト)であるかどうか判断される。そうであれば、コンテキスト・マスクがセットされ (9 8 6)、コンテキストが単独であることを示し、コンテキスト処理 9 8 4 は完了し、戻る。他方、アクションが End Output X n でなければ、アクションが End Output X 1 (つまり、初期コンテキスト)であるかどうか判断 9 8 7 で判断される。そうであれば、コンテキスト・マスクがセットされ (9 8 8)、コンテキストが初期であることを示し、コンテキスト処理 9 8 4 は完了し、戻る。他方、アクションが End Output X 1 でなければ、アクションが End Output X r (つまり、終了コンテキスト)であるかどうか判断 9 8 9 される。そうであれば、コンテキスト・マスクがセットされ (9 9 0)、コンテキストが初期であることを示し、コンテキスト処理 9 8 4 は完了し、戻る。他方、アクションが End Output X r でなければ、アクションが end Output X m (つまり、中間コンテキスト)であるかどうかの判断 9 9 1 が行われる。そうであれば、コンテキスト・マスクがセットされ (9 9 2)、コンテキストが初期であることを示し、コンテキスト処理 9 8 4 は完了し、戻る。アクションが End Output X m でないと判断 9 9 1 されたときは、テキスト要素は関連するコンテキストをもっていないので、コンテキスト・マスクは「無視」にセットされる (9 9 3)。

#### 【 0 1 2 2 】

以上のように、コンテキスト処理 9 8 4 は上記実施例では、テキスト要素とコンテキストと一緒に判断されるように実現されている。どちらの判断も、スキャナ 4 0 8 とスキャナ・テーブル 4 1 0 によって先読みスキャンニングを利用してテキスト要素が完成し、入力テキスト・ストリーム内のテキスト要素の内容が分かるようにしている。コンテキスト・マスクに入っているコンテキスト情報は、マッピング・テーブル 4 1 4 とやりとりして、判断されたコンテキストをもつテキスト要素の正しいターゲット・コード化を探すときにルックアップ・ハンドラ 4 1 2 によって使用される。コンテキスト処理のオペレーションの詳細は図 1 6 A と図 1 6 B を参照して以下で説明する。

#### 【 0 1 2 3 】

図 1 1 はスキャナ 4 0 8 を示すブロック図である。特に、スキャナ 4 0 8 は属性ハンドラ 1 0 0 0 とテキスト要素ハンドラ 1 0 0 2 を含んでいる。テキスト要素ハンドラ 1 0 0 2 は図 9 A と図 9 B および図 1 0 A ~ 図 1 0 C を参照して上述した次のテキスト要素処理 9 0 0 を実行する。属性ハンドラ 1 0 0 0 は属性テーブル 1 0 0 4 とやりとりして、次のテキスト要素処理 9 0 0 が必要とするユニコード文字の属性を取得する (図 9 A のブロック 9 1 4 と図 1 0 A のブロック 9 6 0 を参照)。属性には次のものがある。つまり、方向、クラス、優先順位、シメトリック・スワッピング、サブセットおよびコンテキストである。方向属性は方向を解明するとき使用される (図 9 A のブロック 9 1 4 および図 1 0 A と図 1 7 A のブロック 9 5 6 を参照)。クラス属性はアクション (例えば、ADVANCE、END) を判断するためにスキャナ 4 0 8 によって使用される。優先順位属性はテキスト要素内の文字を記録するために使用される (図 9 B のブロック 9 4 0 と図 1 0 A のブロック 9 8 8 を参照)。シメトリック・スワッピング属性はシメトリック・スワッピングが必要であるかどうかを判断するために使用される。コンテキスト属性はコンテキストを解明するとき使用される (図 9 A のブロック 9 1 0 を参照)。

#### 【 0 1 2 4 】

図 1 2 は図 1 1 の属性テーブル 1 0 0 4 の好ましいフォーマットを示す概略図である。属性テーブル 1 0 0 4 はヘッダ部分 1 1 0 0、範囲テーブル部分 1 1 0 2、および属性テーブル部分 1 1 0 4 を含んでいる。ヘッダ部分 1 1 0 0 は次のものに関する情報を含んでいる。つまり、各テーブルの総テーブル・サイズ、チェックサム値、バージョン、オフセットおよび要素の数である。範囲テーブル部分 1 1 0 2 内の要素は属性が共通にグループ化されている範囲を示し、各グループごとに、属性テーブル部分 1 1 0 4 の該当部分を指すポインタをもっている。属性テーブル 1 0 0 4 の範囲テーブル部分 1 1 0 2 のフォーマ

ットは、例えば、各エントリごとに、範囲値の開始、範囲値の終了、および範囲に関連するデータ・ワードを含んでいる。属性テーブル 1004 の構成は属性情報のコンパクトなストアを容易にしている。別のストア構成も可能であるが、データ・ストアのコンパクトさから見たとき非効率的である。

#### 【0125】

図 13 は属性ルックアップ処理 1200 を示すフローチャートである。属性ルックアップ処理 1200 は属性ハンドラ 1000 によって実行され、図 9A のブロック 914 または図 10A のブロック 960 で次のテキスト要素処理 900 によって開始される。

#### 【0126】

属性ルックアップ処理 1200 は属性テーブル 1004 の範囲テーブル部分 1102 内の範囲を使用してバイナリ・サーチから開始される。該当の範囲がバイナリ・サーチで見つかり、その範囲に関連するデータ・ワードが範囲テーブル部分 1102 から取得される (1204)。好ましくは、データ・ワードの最初のビットは間接ビットである。次に、データ・ワードの間接ビットがセットされているかどうかの判断 1206 が行われる。データ・ワードの間接ビットがセットされていない場合は、データ・ワード自体は現在の文字の属性を収めているので、データ・ワードは属性として戻される (1208)。他方、データ・ワードの間接ビットがセットされていれば、属性は、範囲テーブル部分 1102 から取得されたデータ・ワードを属性テーブル部分 1104 までのインデックスまたはオフセットとして使用して属性テーブル部分 1104 から取得される。従って、この場合のデータ・ワードは範囲内の各文字の属性を収めている配列までのインデックスまたはオフセットである。ブロック 1208 または 1210 に続いて、属性ルックアップ処理 1200 は完了し、戻る。

#### 【0127】

図 14A と図 14B は次のアクションを判断するためにスキャナ 408 によって使用されるスキャナ・テーブル 1300 (410) に関連する概略図である。次のアクションの判断は次のテキスト要素処理 900 (図 9A) のブロック 916 によって起動されるか、ブロック 962 (図 10A) によって起動される。図 14A は本発明で使用されるスキャナ・テーブル 1300 の好ましいフォーマットを示す図である。スキャナ・テーブル 1300 は「現在状態」を 1 つのインデックスとして、「クラス」を別のインデックスとして持つ 2 次元配列である。「クラス」はクラス属性を指している。これらのインデックスはスキャナ・テーブル 1300 内の代表的な要素 1302 を選択する。図 14B はスキャナ・テーブル 1300 の代用的な要素 1302 を示す図である。代表的な要素 1302 はスキャナ 408 の次の状態を収めている次の状態部分 1304 を含み、アクション部分 1306 はスキャナ 408 のアクションを示している。実際には、スキャナ 408 はスキャナ・テーブル 1300 と一緒になって、スキャナをどのように動作させるかを判断するステート・マシンを実現している。

#### 【0128】

スキャナ・テーブル 1300 は異種の文字コード化に関する異なる次の状態とアクションを収めている。

#### 【0129】

本発明の第 1、第 2 および第 4 側面によれば、図 15 は好ましいレイアウトとスキャナ・テーブル 1300 にストアされる情報の両方を示しているテーブル 1400 である。このテーブルに関連する表記は次の通りである。

#### 【0130】

文字クラス：

- CC - 制御文字
- OS - 他のスペース
- NS - 非スペース
- LD - ラテン・ディジット
- FS - フラクション・スラッシュ

J L ( f ) - ジャモス先頭子音 ( フィラー )  
 J V ( f ) - ジャモス母音 ( フィラー )  
 J T - ジャモス・トレーラ  
 N U - 有効なユニコード文字でない  
 I S S - シメトリック・スワッピング禁止  
 A S S - シメトリック・スワッピング活動化

## 【 0 1 3 1 】

次の状態：

状態 0 - 終了。テキスト要素を戻すべきかどうかを、2 重および半分音符号の状態に基づいて判断する。

10

## 【 0 1 3 2 】

状態 1 - 開始状態  
 状態 2 - 非スペース ( 分音符号 ) の追加  
 状態 3 - 数値フラクシヨンの有無チェック  
 状態 7 - 朝鮮語ジャモス

## 【 0 1 3 3 】

アクション：

A d v - [ A D V A N C E ] 次の文字へ進む ( 現在文字は現在テキスト要素 ( T E ) に含まれている場合と含まれていない場合がある )  
 A d v M a r k - [ A D V A N C E + M A R K ] 現在文字に最終文字のマークを付け次の文字へ進む  
 A d v M a r k S - [ A D V A N C E + M A R K + S ] 現在文字に最終文字のマークを付け次の文字へ進み、再配列フラグをセットする  
 A d v M a r k A S S - [ A D V A N C E + M A R K + A S S ] 現在文字に最終文字のマークを付け次の文字へ進み、シメトリック・マッピングを活動化する  
 A d v M a r k I S S - [ A D V A N C E + M A R K + I S S ] 現在文字に最終文字のマークを付け次の文字へ進み、シメトリック・マッピングを禁止する  
 E n d - テキスト要素を最終のマークを付けた文字で終了する

20

30

## 【 0 1 3 4 】

注意：すべての機能は再配列フラグがセットされているかをチェックして確かめ、非スペース文字を開始ポイントから始めて再配列する。再配列する必要のある非スペース・マークのセットは複数存在することがあるので、文字列全体をチェックするのが最良の方法である。もちろん、文字を再配列すると、再配列フラグはクリアされる。

## 【 0 1 3 5 】

本発明の第 3 側面によれば、図 1 6 A および図 1 6 B は好ましいレイアウトと情報の両方を示すテーブル 1 4 0 0 であり、情報はソース文字列内のテキスト要素とそのコンテキストの両方の判断を容易にするためにスキャナ・テーブル 4 1 0 にストアされるものである。このテーブルに関連する表記は次のとおりである。

40

## 【 0 1 3 6 】

文字クラス：

C C - 制御文字  
 O S - 他のスペース  
 N S - 非スペース  
 L D - ラテン・ディジット  
 F S - フラクシオン・スラッシュ  
 D D - 二重分音符号

50



H D	- 半分音符号
C H	- ジャモス先頭子音（フィラー）
J O	- ジャモス母音（フィラー）
J V	- ジャモス子音トレーラ
N U	- 有効なユニコード文字でない
I S S	- シメトリック・マッピング禁止
A S S	- シメトリック・マッピング活動化
H H	- 高ハーフゾーン
L H	- 低ハーフゾーン
V	- V i r a m a
Z W N J	- ゼロ幅の非ジョインダ
R	- 右リンク
D	- 2重リンク
C	- リンクを引き起こす
T	- 透過

10

## 【 0 1 3 7 】

次の状態：

状態 0 - 終了。テキスト要素を戻すべきかどうかを 2 重分音符号と半分音符号に基づいて判断する

状態 1	- 開始状態
状態 2	- 非スペース（分音符号）追加
状態 3	- ラテン・ディジット
状態 4	- ラテン・ディジット・シーケンス
状態 5	- ラテン・ディジット・シーケンスとそのあとに続くフラクション・スラッシュ

20

シュ

状態 6 - ラテン・ディジット・シーケンス、フラクション・スラッシュ、ラテン・ディジット・シーケンス

状態 7	- C h o s e o n g シーケンス
状態 8	- C h o s e o n g シーケンスとそのあとに続く J u n g s e o n g シーケンス

30

状態 9 - C h o s e o n g シーケンス、J u n g s e o n g シーケンス、J o n g s e o n g

状態 10	- 高半文字
状態 11	- 高半文字とそのあとに続く低半文字
状態 12	- V i r a m a とそのあとに続くゼロ幅ジョインダまたはゼロ幅非ジョインダ

ダ

状態 13	- コンテキストの特殊開始
状態 14	- 右リンク文字（特殊コンテキスト状態）
状態 15	- 2重リンク文字（特殊コンテキスト状態）
状態 16	- リンクを引き起こす文字（通常または特殊状態）
状態 17	- 右リンク文字（通常状態）
状態 18	- 2重リンク文字（通常状態）

40

## 【 0 1 3 8 】

次の開始状態：[ 実際は次の状態のサブセット ]

開始 1	- 開始状態
開始 13	- コンテキストの特殊開始状態

## 【 0 1 3 9 】

アクション：

A d v	- [ A D V A N C E ] 次の文字へ進む（現在文字は現在テキスト要素（T E）に含まれている場合と含まれていない
-------	---

50

	い場合がある)	
AdvMark	- [ADVANCE + MARK] 現在文字に最終文字の マークを付け次の文字へ進む	
AdvMarkS	- [ADVANCE + MARK + S] 現在文字に最終文字 のマークを付け次の文字へ進み、再配列フラグをセット する	
AdvMarkASS	- [ADVANCE + MARK + ASS] 現在文字に 最終文字のマークを付け次の文字へ進み、シメトリック・ マッピングを活動化する	
AdvMarkISS	- [ADVANCE + MARK + ISS] 現在文字に 最終文字のマークを付け次の文字へ進み、シメトリック・ マッピングを禁止する	10
End	- テキスト要素を最終のマークを付けた文字で終了する	
EndOutputXn	- [END + 単独コンテキスト] テキスト要素を終 了し、単独コンテキストであることを示す	
EndOutputX1	- [END + 初期コンテキスト] テキスト要素を終 了し、初期コンテキストであることを示す	
EndOutputXr	- [END + 終了コンテキスト] テキスト要素を終 了し、終了コンテキストであることを示す	
EndOutputXm	- [END + 中間コンテキスト] テキスト要素を終 了し、中間コンテキストであることを示す	20

## 【0140】

注意：すべての機能は再配列フラグがセットされているかをチェックして確かめ、非スペース文字を開始ポインタから始めて再配列する。再配列する必要のある非スペース・マークのセットは複数存在することがあるので、文字列全体をチェックするのが最良の方法である。もちろん、文字を再配列すると、再配列フラグはクリアされる。

## 【0141】

以下では、スキャナ・テーブル1400の使用例を3つ示して説明する。最初の2つの例は本発明の第1、第2、第3および第4側面に関係し、最後の例は第3側面に関係するものである。

## 【0142】

例1：入力文字列“ A A B ”

文字クラスは3文字ともOSである。最初の文字“ A ”が取得される。開始状態（状態1）から始まり、最初のアクションはAdvMark、次の状態は状態2である。これにより、最初の文字“ A ”が現在テキスト要素内に挿入され、次の文字（2番目の文字“ A ”）が取得される。次に、状態2では、アクションはEnd、次の状態は状態0である。従って、テキスト要素は最初のテキスト要素だけを含んでいる。同じシーケンスはこの特定入力文字列の2番目と3番目の文字について繰り返される。このようにして、入力文字列の文字の各々は分離しているが、隣接しているテキスト要素に割り当てられる。

## 【0143】

例2：入力文字列“ A ` B ”

文字クラスは入力文字列の最初と最後の文字ではOSである。2番目の文字の文字クラスがNSであるのは、これが結合マークであるためである。最初の文字“ A ”が取得される。開始状態（状態1）から始まり、最初のアクションはAdvMark、次の状態は状態2である。これにより、最初の文字“ A ”が現在テキスト要素内に挿入され、次の文字（2番目の文字“ ` ”）が取得される。次に、状態2では、アクションはAdvMarkS、次の状態は状態2である。これにより、2番目の文字“ ` ”が現在テキスト要素に挿入される。次に、3番目の文字が取得される。この時点の状態2のアクションはEnd、次の状態は状態0である。従って、テキスト要素は入力文字列の1番目と2番目の文字を含んでいる。3番目の文字は例1の場合と同じように、自身のテキスト要素に置かれる。

10

20

30

40

50

## 【 0 1 4 4 】

例 3：入力文字列 “ O S R D O S O S ” [ 文字間にスペースがなく、各文字はその文字クラスで表されている。R と D 文字クラスはコンテキスト・ベースの表示形態をもつ文字を含んでいるが、O S 文字クラスは含んでいない。]

## 【 0 1 4 5 】

最初の文字は文字クラス O S である。開始状態（状態 1）から始まり、最初のアクションは A d v M a r k、次の状態は 2、次の開始状態は 1 である。これにより、最初の文字が現在テキスト要素に挿入され、2 番目の文字が取得される。2 番目の文字は文字クラス R である。次に、状態 2 では、アクションは E n d、次の状態は 0、次の開始状態は 1 である。従って、最初のテキスト要素は最初の文字だけを含み、コンテキスト・フラグは無視にセットされる。

10

## 【 0 1 4 6 】

次のテキスト要素を判断するときは、処理は 2 番目の文字から始まり、新しい開始状態は 1 になっている。この時点で、スキャナ・テーブルから得られるアクションは A d v M a r k であり、次の状態は 1 7、次の開始状態は 1 3 である。これにより、2 番目の文字が現在テキスト要素に挿入され、3 番目の文字が取得される。3 番目の文字は文字クラス R である。次に、状態 1 7 では、アクションは E n d O u t p u t X n、次の状態は 0、次の開始状態は 1 である。従って、2 番目のテキスト要素は 2 番目の文字だけを含み、コンテキストは X n（単独）である。

20

## 【 0 1 4 7 】

次のテキスト要素を判断するときは、処理は 3 番目の文字から始まり、新しい開始状態は 1 3 になっている。この新しい開始状態は最後の文字についてテーブルで示された次の開始状態である。ここでは、スキャナ・テーブルから得られるアクションは A d v M a r k であり、次の状態は 1 5、次の開始状態は 1 3 である。3 番目の文字は現在テキスト要素の一部であるので、4 番目の文字が取得される。4 番目の文字は文字クラス O S である。次に、状態 1 5 では、アクションは E n d O u t p u t X r、次の状態は 0、次の開始状態は 1 である。従って、3 番目のテキスト要素は 3 番目の文字だけを含み、コンテキストは X r（終了）である。

## 【 0 1 4 8 】

次のテキスト要素を判断するときは、処理は 4 番目の文字から始まり、新しい開始状態は 1 3 になっている。スキャナ・テーブルから得られるアクションは A d v M a r k であり、次の状態は 2、次の開始状態は 1 である。これにより、4 番目の文字が現在テキスト要素に挿入され、5 番目の文字が取得される。5 番目の文字は文字クラス O S である。次に、状態 2 では、アクションは E n d、次の状態は 0、次の開始状態は 1 である。従って、4 番目のテキスト要素は 4 番目の文字だけを含み、コンテキスト・フラグは無視にセットされる。

30

## 【 0 1 4 9 】

図 1 7 A は本発明の第 2 側面による方向解明処理 1 5 0 0 を示すフローチャートである。方向解明処理 1 5 0 0 は次のテキスト要素処理 9 0 0（図 9 A）内のブロック 9 0 6 で実行される処理である。方向解明処理 1 5 0 0 は、スキャナ 4 0 8 のために解明された方向の状態が初期状態にあるかどうかを判断する判断 1 5 0 2 から開始される。初期状態では、スキャナ 4 0 8 はテキスト要素内の文字の方向を知らない。スキャナ 4 0 8 のために解明された方向の状態が初期状態にあれば、入力文字列の初期方向が判断される（1 5 0 4）。初期状態 1 5 0 4 の判断 1 5 0 4 に関連する処理は図 1 8 を参照して以下で詳しく説明する。他方、スキャナ 4 0 8 のために解明された方向の状態が初期状態にないときは、初期方向を判断する必要はない。どちらの場合も、ブロック 1 5 0 2 またはブロック 1 5 0 4 に続いて、ユニコード文字が入力ストリームから取得される（1 5 0 6）。次に、ユニコード文字に関連する属性が調べられる（1 5 0 8）。方向属性は方向解明処理 1 5 0 0 を行うとき重要な属性である。方向解明処理 1 5 0 0 のブロック 1 5 0 6 と 1 5 0 8 が実行するオペレーションは次のテキスト要素処理 9 0 0 のブロック 9 1 2 と 9 1 4 のそ

40

50

れと同じであるので、これ以上詳しく説明することは省略する。

#### 【 0 1 5 0 】

次に、方向、次の状態およびアクションが判断される ( 1 5 1 0 )。方向、次の状態およびアクションは方向属性と現在状態を使用して判断され、テーブル属性ルックアップ・プロセスを使用して方向テーブルから取得される。方向テーブルは方向属性と現在状態によってインデックスされる 2 次元配列である。インデックスが指している方向テーブル内の要素は文字の方向、次の状態、および方向解明処理 1 5 0 0 がとるべきアクションを含んでいる。方向は次のものの 1 つである。すなわち、左、右、グローバル、および NO\_\_OUTPUT である。取り得るアクションは、NO\_\_ACTION、PUSH RO、PUSH RE、PUSH LE、PUSH LO、POP、および RESET である。方向は明示的オーバーライド文字により入力文字列の中で変化することがあるので、以前の方向は方向スタック上にストアされている。従って、「プッシュ」および「ポップ」の使用はこの分野で周知であるスタック操作コマンドを意味する。“RO”は右から左へのオーバーライドを意味し、“LO”は左から右へのオーバーライドを意味し、“RE”は右から左への埋め込みを意味し、“LE”は左から右への埋め込みを意味する。好ましくは、方向テーブルと方向解明処理 1 5 0 0 はステート・マシンとして動作する。ステート・マシンは基本的にユニコード標準バージョン 1 . 0 6 1 1 - 6 2 1 ページ ( 付録 A ) に記載されている双方向アルゴリズムに従って動作するが、1 回のパスだけで結果が得られるのに対し、ユニコード標準に記載されているアルゴリズムは複数のパスを必要とする。

#### 【 0 1 5 1 】

図 1 7 B ~ 図 1 7 D は本発明の第 2 側面による双方向アルゴリズムの好ましいインプリメンテーションによる双方向状態テーブル 1 5 1 1 を示している。双方向状態テーブル 1 5 1 1 はテーブル駆動型ステート・マシンを実現している。各カラムは単一状態である。状態はそこに記録される情報を示唆する名前が付けられている。各行は図 1 7 B ~ 図 1 7 D にまたがっており、以下に示す文字クラス名の 1 つが付けられている。

#### 【 0 1 5 2 】

LR	左から右へが主流
RL	右から左へが主流
AL	アラビア文字 ( 右から左へが主流 )
LRE	左から右への埋め込みマーク
RLE	右から左への埋め込みマーク
LRO	左から右へのオーバーライド・マーク
RLO	右から左へのオーバーライド・マーク
PDF	ポップ方向フォーマット・マーク
AN	アラビア数字
EN	ヨーロッパ数字
ET	ヨーロッパ数字ターミネータ
ES	ヨーロッパ数字セパレータ
CS	共通数字セパレータ
ON	その他の中立文字
BS	ブロック・セパレータ

#### 【 0 1 5 3 】

双方向状態テーブル 1 5 1 1 の各セルは関連のアクションと出力を持つ遷移を表している。これらのセルは新しい状態の名前、とられるオプションのアクション、および出力 ( もしあれば ) を収めている。新しい状態は現在のグローバル方向に左右される場合がある。これは新しい状態名に ( G ) を入れることで示される。グローバル方向が未知のときこれらの遷移の 1 つを行うとエラーになる。取り得るアクションは次のとおりである。

#### 【 0 1 5 4 】

プッシュ	新しい埋め込み状態を埋め込みスタック上にプッシュする。プッシュされる実際の値は現れた実際の埋め込み制御によって決ま
------	---

る。このインプリメンテーションでは、これ进行处理するアクション動詞は4つある。

#### 【0155】

ポップ 現在の埋め込み状態をスタックからポップし、スタックの新しいトップを現在の埋め込み状態にする。新しい埋め込みがオーバーライドであれば、セルに入っているターゲット状態ではなくOR状態に遷移が行われる。

#### 【0156】

リセット 埋め込み状態をクリアし、文字を消費することなくSTARTに移る。すなわち、エプシロン遷移を行う。リセットには出力はない。

10

#### 【0157】

エラー 即時エラーを引き起こす。

#### 【0158】

出力はそれぞれ左から右へと右から左へを意味するLまたはRのどちらか、現在のグローバル方法の出力を意味するG、または出力なしを意味する\*である。マシンには、小文字の's'で始まらない名前を持つ状態で入ることができる。START状態は新しいテキスト・ブロックを目的としている。sDIR状態はグローバル方向を判断するときのエントリ・ポイントとして使用される。この計算はメイン・スキャンと同時に行うことが可能であるが、別々にすると単純化される。

20

#### 【0159】

図15Aに戻って説明すると、ブロック1510に続いて、アクションが“RESET”であるかどうかの判断1512が行われる。アクションが“RESET”であれば、処理は方向解明処理1500の先頭に戻る。そうでなければ、方向解明処理1500が続行される。つまり、判断1512に続いて、アクションが実行される(1514)。

#### 【0160】

次に、方向(ブロック1510で判断されたもの)が“NO\_\_OUTPUT”であるかどうかの判断1516が行われる。方向が“NO\_\_OUTPUT”に等しくなければ、方向がコンテキスト(各インスタンスごとに)にセットされ(1518)、状態が保存される(1520)。ブロック1520に続いて、方向解明処理1500は完了し、戻る。

30

#### 【0161】

しかるに、出力が“NO\_\_OUTPUT”であると判断1516で判断されたときは、入力文字列全体が処理されたかどうかの判断1522が行われる。入力文字列全体が処理されていないければ、処理は方向解明処理1500の先頭に戻り追加のユニコード文字を取得し、処理できるようにするが、これはテキスト要素の方向がまだ判断されていないためである。他方、入力文字列全体が処理されたら判断1522されたときは、処理がテキストの終わりまで達したかどうか判断1524される。つまり、ターゲット・コード化に変換すべきテキストが方向解明処理1500によって完全に処理されたかどうか判断される。変換すべき追加テキストがあれば、入力文字列の現在の文字では方向が計算できなかったのでエラーになる(1528)。しかし、追加テキストがなければ、ブロック・セパレータの方向がブロックの全体方向から判断される(1526)(1504を参照)。ブロック1526に続いて、ブロック1518と1520が前述したように実行され、そのあと方向解明処理1500は完了し、戻る。

40

#### 【0162】

図18は本発明の第2側面による初期方向判断処理1600を示すフローチャートである。初期方向判断処理1600は図17Aのブロック1504で実行されるオペレーションと関連している。

#### 【0163】

初期方向判断処理1600は制御フラグに基づくスイッチ・オペレーションから開始される。制御フラグは次の1つを示している。すなわち、

50

NO\_\_OUTPUT、L - t o - RまたはR - t o - Lである。これらの制御フラグはユニコード・コード変換システム400を起動するアプリケーションによってセットされる(つまり、制御フラグはコンバータへの入力である)。制御フラグがR - t o - Lを示しているときは、グローバル方向はR - t o - Lにセットされる(1604)。制御フラグがL - t o - Rにセットされているときは、グローバル方向はL - t o - Rにセットされる(1606)。制御フラグがNO\_\_OUTPUTにセットされているときは、スモール・ループが開始され、方向が判断できるまで入力文字列のユニコード文字をスキャンしていく。ループは状態を“START STATE”にセット(1608)することから開始される。次に、ユニコード文字が入力文字列から取得される(1610)。ユニコード文字の属性が次に調べられる(1612)。属性は図9Aのブロック914で使用され、図12に詳しく説明されているものと同じ方法を用いて調べられる(1612)。次にユニコード文字の方向が属性(つまり、方向属性)を使用して判断される(1614)。方向が“NO\_\_OUTPUT”に等しいかどうかの判断1616が行われる。方向が“NO\_\_OUTPUT”であれば、入力文字列の終わりまで達したかどうかの判断1618が行われる。入力文字列の終わりまで達していなければ、処理は戻り、ブロック1610~1618を繰り返す。文字列の終わりまで達していると、判断1618は方向を見つけることなく特殊方向スキャン・ループを終わらせる(例えば、NO\_\_OUTPUT)。そうでなければ、スキャン・ループは判断1616が方向を判断したとき終了する。  
【0164】

いずれの場合も、スイッチ・オペレーション1602に関連する方向処理に続いて、スイッチ・オペレーション1620は判断に基づいて起動される。判断された方向が2NO\_\_OUTPUTであるときは、現在レベルは埋め込みボトムにセットされる(1622)。他方、方向がL - t o - Rであれば、現在レベルはゼロにセットされ、前のレベルはボトムにセットされ、オーバーライド状況は中立にセットされる(1624)。方向がR - t o - Lであるときは、現在レベルは1にセットされ、前のレベルはボトムにセットされ、オーバーライド状況は中立にセットされる(1626)。レベルとは、ユニコード標準に記載されている埋め込みレベルのことである。初期方向判断処理1600の特殊方向スキャン・ループ(例えば、1610~1618)は図18に別ループとして示されているが、総システム処理効率は特殊方向スキャン・ループをメイン方向スキャン・ループ内に組み入れると向上することができる。

【0165】

図19はテキスト要素ルックアップ処理1630を示すフローチャートである。テキスト要素ルックアップ処理1630はルックアップ・ハンドラ412によって実行され、ユニコード・コンバータ処理700(図7)内のブロック708によって起動される。

【0166】

テキスト要素ルックアップ処理1630は変種リストをサーチし、実際の属性と要求された変種に一致するエントリを見つけることから開始される。変種のサーチ1632は図24を参照して以下で詳しく説明する。次に、一致するものが見つかったかどうかの判断1634が行われる。一致するものが見つからなければ、エラーになる(1636)。他方、一致するものが見つければ、対応するビット・マスクが変種リストから取得される(1638)。好ましくは、変種リストは3つのフィールド、つまり、変種識別子、属性のセット、およびビット・マスクをもっている。変種リスト内の変種識別子と属性セットが実際の属性および要求された変種と一致していれば、対応するビット・マスクが変種リストから選択される。ブロック1638に続いて、ビット・マスクは許容範囲ビット・マスクと結合され(1640)選択フラグが得られる。選択フラグは上述したようにマッピング・テーブル414のサブテーブルを選択するとき使用される選択マスクを形成する。好ましくは、本発明の第2側面では、結合1640はビットワイズ・オペレーションである。好ましくは、本発明の第1、第3および第4側面では、結合1640はビットワイズORオペレーションである。次に、現在テキスト要素の長さのためのテーブルがマッピング・テーブル414内にあるかどうかの判断1642が行われる。なければ、エラーになる

( 1 6 4 4 )。他方、現在テキスト要素の長さのためにテーブルがあれば、ルックアップ・テーブルとそのフォーマットが現在テキスト要素の長さのために取得される ( 1 6 4 6 )。次に、スイッチ・オペレーション 1 6 4 8 がフォーマットに基づいて実行される。このインプリメンテーションで利用できるフォーマットはリスト、セグメント配列、範囲およびチェーンである。フォーマットがリスト・フォーマットであれば、リスト・フォーマット処理 1 6 5 0 が実行される。フォーマットがセグメント配列であれば、セグメント配列フォーマット処理 1 6 5 2 が実行される。フォーマットが範囲であれば、範囲フォーマット処理 1 6 5 4 が実行される。フォーマットがチェーンであれば、チェーン・フォーマット処理 1 6 5 6 が実行される。ブロック 1 6 5 0 - 1 6 5 6 に続いて、結果が戻され ( 1 6 5 8 )、これによりテキスト要素ルックアップ処理 1 6 3 0 は完了する。

10

#### 【 0 1 6 7 】

本発明の第 1 側面によれば、図 2 0 はチェーン・フォーマット処理 1 6 6 0 を示すフローチャートである。チェーン・フォーマット処理 1 6 6 0 は図 1 9 に示すチェーン・フォーマット処理 1 6 3 5 によって実行される処理である。

#### 【 0 1 6 8 】

チェーン・フォーマット処理 1 6 6 0 はチェーン内のテーブル数のチェーン・カウントを取得する ( 1 6 6 2 )。そのあと、現在のカウンタはゼロにセットされる ( 1 6 6 4 )。次に、現在のカウンタがチェーン・カウントより大であるか等しいかどうかの判断 1 6 6 6 が行われる。現在のカウンタがチェーン・カウントより大または等しければ、チェーン・フォーマット処理 1 6 6 0 は戻り ( 1 6 6 8 )、結果が見つからなかったのでエラーを通知する。他方、現在のカウンタがチェーン・カウントより大でも等しくもなければ、現在のルックアップ・テーブルとそのフォーマットが取得される ( 1 6 7 0 )。このインプリメンテーションで使用されるフォーマットは図 1 9 で使用されたものと同じである。そのあとスイッチ・オペレーション 1 6 7 2 がフォーマットに基づいて実行される。フォーマットがリストであれば、リスト・フォーマット処理 1 6 7 4 が実行される。フォーマットがセグメント配列であれば、セグメント配列フォーマット処理 1 6 7 6 が実行される。フォーマットが範囲であれば、範囲フォーマット処理 1 6 8 0 が実行される。フォーマットがチェーンであれば、チェーン・フォーマット処理 1 6 8 2 が実行される。ブロック 1 6 1 4 ~ 1 6 2 2 に続いて、現在のカウンタがインクリメントされる ( 1 6 8 4 )。次に、結果が見つかったかどうかの判断 1 6 8 6 が行われる。結果が見つからなければ、チェーン・フォーマット処理 1 6 6 0 はループして判断ブロック 1 6 6 6 に戻り、テーブル・チェーン内の次のルックアップ・テーブルを調べていく。しかるに、結果が見つかったと判断 1 6 8 6 されたときは、結果が戻され ( 1 6 8 8 )、これによりチェーン・フォーマット処理 1 6 6 0 は完了する。

20

30

#### 【 0 1 6 9 】

本発明の第 1 側面によれば、図 2 1 は範囲フォーマット処理 1 7 0 0 を示すフローチャートである。範囲フォーマット処理 1 7 0 0 は図 1 9 のブロック 1 6 5 4 と図 2 0 のブロック 1 6 7 8 によって実行される処理である。範囲フォーマットはデルタ値が各フィールドと関連づけられている文字の範囲リストである。

#### 【 0 1 7 0 】

範囲フォーマット処理 1 7 0 0 はそのサブテーブルのサブセット・フラグが選択フラグに一致しているかどうかを判断 ( 1 7 0 2 ) することから始まる。選択フラグは選択マスクに対するもので、サブセット・フラグはサブテーブル・マスクに対するものである。一致していなければ、範囲フォーマット処理は戻り ( 1 7 0 4 )、結果が見つからなかったことをエラー・コードで通知する。他方、そのサブテーブルのサブセット・フラグが選択フラグに一致していることを判断 1 7 0 2 が示していれば、テキスト要素の長さが 1 より大であるかどうかの判断 1 7 0 6 が行われる。テキスト要素の長さが 1 より大であれば、このフォーマットは誤って選択されたことになる。この特定インプリメンテーションでのマッピング・テーブルの構成は範囲フォーマットが長さ 1 のテキスト要素だけを目的とするようになっているためである。従って、テキスト要素が 1 より大であれば、ブロック 1

40

50

704が実行されて戻り、結果が見つからなかったことを通知する。他方、テキスト要素の長さが1より大でなければ、範囲フォーマット処理1700は続行される。範囲フォーマットをもつサブテーブルは範囲配列をもち、各範囲はそれぞれに関連づけられたデルタ値をもっている。次に、範囲配列がサーチされ(1708)、変換されるユニコード文字の該当する範囲が見つけれれる。範囲が見つかったかどうかの判断1710が行われる。範囲が見つからなければ、範囲フォーマット処理1700は戻り(1704)、結果が見つからなかったことをエラー・コードで通知する。しかるに、結果が見つかったときは、その範囲の対応するデルタ値が取得される(1712)。次に、このデルタ値はユニコード値に加えられる(1714)。そのあと、その結果は出力シーケンスにマッピングされる(1716)。結果を出力シーケンスにマッピングする処理は図22A~図22Cを参照して以下で詳しく説明する。ブロック1716に続いて、範囲フォーマット処理は完了し、戻る。

#### 【0171】

本発明の第1側面によれば、図22はリスト・フォーマット処理1800を示すフローチャートである。リスト・フォーマット処理1800は図19のブロック1650と図20のブロック1674によって実行される処理である。リスト・フォーマットはテキスト要素が配列されたリストであり、この配列リストまでのインデックス*i*は対応するルックアップ・ターゲット・リストまでのインデックスである。

#### 【0172】

リスト・フォーマット処理1800はそのサブテーブルのサブセット・フラグが選択フラグに一致しているかどうかの判断1802から始まる。一致していなければ、リスト・フォーマット処理1800は戻り(1804)、マッピングが見つからなかったことをエラー・コードで通知する。他方、そのサブテーブルのサブセット・フラグが選択フラグに一致していれば、リスト内のテキスト要素に基づいて最適化バイナリ・サーチが行われる(1806)。次に、サーチがリスト内のテキスト要素を見つけたかどうかの判断1808が行われる。見つからなければ、再びリスト処理は戻り(1804)、マッピングが見つからなかったことを通知する。しかし、テキスト要素が見つければ、テキスト要素が見つかった個所のインデックス*i*が取得される(1810)。このインデックス*i*はルックアップ・ターゲットを取得する(1812)ために使用される。次に、ルックアップ・ターゲットは出力シーケンスにマッピングされる(1814)。これでリスト・フォーマット処理は完了し、戻る。

#### 【0173】

本発明の第1側面によれば、図23Aと図23Bはセグメント配列フォーマット処理1900を示している。セグメント配列フォーマット処理1900は図19のブロック1652と図20のブロック1676によって実行される処理である。セグメント配列フォーマットは第1テキスト要素配列、最終テキスト要素配列、および*n*個のオフセット配列を含んでいる。オフセット配列内のオフセットは種々のルックアップ・ターゲット・リストを指している。

#### 【0174】

セグメント配列フォーマット処理1900はそのサブテーブルのサブセット・フラグが選択フラグに一致しているかどうかの判断1901から始まる。サブセット・フラグが選択フラグに一致していれば、最適化サーチが行われる(1902)。最適化サーチは最終テキスト要素配列内であって、探索していたテキスト要素よりも大であるか等しい最小エントリを見つけ、そのエントリのインデックス*i*を取得する。次に、第1テキスト要素配列内の*i*番目のエントリが探索していたテキスト要素より小であるか等しいかどうかの判断1904が行われる。そうでなければ、セグメント配列フォーマット処理1900は戻り(1906)、マッピングが見つからなかったことをエラー・コードで知らせる。また、判断1901が失敗した場合には、ブロック1906も実行される。他方、最初のテキスト要素配列内の*i*番目のエントリが探索していたテキスト要素より小であるか等しいことを判断1904が示しているときは、*i*番目のエントリは探索していたテキスト要素に



対応することが分かる。そのあと、ルックアップ・ターゲット・リストがオフセット配列内の  $i$  番目のエントリを通して取得される (1908)。すなわち、オフセット配列内の  $i$  番目のエントリに入っているオフセットはルックアップ・ターゲット・リストを示している。ルックアップ・ターゲット・リストまでのインデックス  $j$  が次に判断される。インデックス  $j$  は探索していたテキスト要素の値から、取得 (1908) したルックアップ・ターゲット・リスト (または範囲) 内の第 1 テキスト要素を引いた値が与えられる。インデックス  $j$  をもつルックアップ・ターゲットが次にルックアップ・ターゲット・リストから取得される (1912)。次に、ルックアップ・ターゲットがゼロに等しいかどうかの判断 1914 が行われる。ルックアップ・ターゲットがゼロに等しければ、セグメント配列フォーマット処理は戻り、マッピングが見つからなかったことをエラー・コードで通知する。他方、ルックアップ・ターゲットがゼロに等しくなければ、ルックアップ・ターゲットは出力シーケンスにマッピングされる (1918)。ブロック 1918 に続いて、セグメント配列フォーマット処理 1900 は完了し、戻る。

#### 【0175】

図 24 はサーチ変種リスト処理 2000 を示すフローチャートである。サーチ変種リスト処理 2000 は図 19 のブロック 1632 によって実行される処理である。言い換えれば、サーチ変種リスト処理 2000 はマッピング・テーブル 414 を使用してルックアップ・ハンドラ 412 によって実行されるルックアップ・テキスト要素処理 1630 の一部である。

#### 【0176】

サーチ変種リスト処理 2000 は変種リスト内の要素の総カウントを取得する (2002)。そのあと、現在のカウントがゼロに初期化される (2004)。現在のカウントが総カウントより大であるか等しいかどうかの判断 2006 が行われる。現在のカウントが総カウントより大または等しければ、サーチ変種リスト処理 2000 は戻り (2008)、変種が見つからなかったことをエラー・コードで通知する。他方、現在のカウントが総カウントより大または等しくなければ、変種リスト内にあって、現在のカウントに関連するエントリが実際の属性および要求された変種に一致しているかどうかの判断 2010 が行われる。一致していれば、変種リスト内のエントリからの変種フラグが戻される (2012)。一致していなければ、現在のカウントはインクリメントされ (2014)、そのあと処理はブロック 2006 に戻り、変種の 1 つが一致するか、あるいは変種のすべてが考慮されるまでループを続けて変種リスト内の使用可能な変種を見つける。

#### 【0177】

図 25 A および図 25 B は変種リスト 2100 を示す概略図である。図 25 A に示すように、変種リスト 2100 は変種領域 2102、所望属性領域 2104 および変種フラグ領域 2106 を含んでいる。図 25 B は好ましいインプリメンテーションによる実際属性ビット・マスクを示す図である。実際属性ビット・マスク 2108 は 32 ビット変数であり、シメトリック・スワッピング状態を示す第 1 部分 2110 (ビット 0 と 1)、垂直または水平形式を示す第 2 部分 (ビット 2 と 3)、解明方向を示す第 3 部分 (ビット 8 と 9)、およびコンテキストを示す第 4 部分 (ビット 16 ~ 19) をもっている。本発明の第 3 側面によれば、第 4 部分 2116 のビットはコンテキスト処理 984 (図 10 C) で判断されたコンテキスト・マスクによってセットされる。さらに、本発明の第 1 ないし第 4 側面によれば、部分内の各ビットはフラグを表している。ビットはフラグが未知であるか、あるいはセットされていないければ、値が “0” であり、セットされているときは値は “1” になっている。呼び出し側は第 2 部分 2110 をセットし、スキャナ 408 は第 1、第 2 および第 4 部分 2110、2114 および 2116 をセットする。

#### 【0178】

所望属性ビット・マスクは実際属性ビット・マスクと同じフォーマットであるが、ビットは属性のどれが特定テーブルおよび変種の正しいマッピングを得るために重要であるかに応じてセットされる (これはマッピング・テーブル 414 の設計によって決まる)。所望属性ビット・マスク内のビットはマッピング判断の際に考慮される各属性には “1” が

10

20

30

40

50

セットされ、ある部分のすべてのビットが“ 1 ”にセットされていれば、属性はマッピング期間には無視される。例えば、ビット 0 が“ 1 ”で、ビット 1 が“ 0 ”ならば、シメトリック・スワッピングはオンであり、マッピングが行われるとき考慮される。他方、ビット 0 と 1 が共に“ 1 ”であれば、シメトリック・スワッピングは完全に無視される。所望属性ビット・マスクの残りの未使用ビットは“ 1 ”にセットされ、必要ならばあとで値を割り当てることができる。以下、所望属性ビット・マスクの例をいくつか示して説明する。

#### 【 0 1 7 9 】

方向が左から右へであり、他の属性はいずれも重要でないとする。そうすると所望属性ビット・マスクは x F F F F F D F F となる。これに対して、方向が右から左へであり、シメトリック・スワッピングがオンであるとする、所望属性ビット・マスクは x F F F F F E F D となる。方向が右から左へで、シメトリック・スワッピングがオフであれば、所望属性ビット・マスクは x F F F F F E F E となる。上記の異種所望属性ビット・マスクの各々によると、異なる変換コードを選択することができる。例えば、ユニコード文字 u 0 0 2 8 を M a c A r a b i c にマッピングすると、所望属性ビット・マスク x F F F F F D F F では x 2 8 が、所望属性ビット・マスク x F F F F F E F D では x A 8 が、所望属性ビット・マスク x F F F F F E F E では x A 9 が得られる。

#### 【 0 1 8 0 】

本発明の第 1 側面によれば、図 2 6 A、図 2 6 B および図 2 6 C はルックアップ・ターゲットの出力シーケンスへのマッピング処理 2 2 0 0 を示すフローチャートである。ルックアップ・ターゲットの出力シーケンスへのマッピング処理 2 2 0 0 は図 2 2 のブロック 1 8 1 4 および図 2 3 A のブロック 1 9 1 8 に関連している。

#### 【 0 1 8 1 】

ルックアップ・ターゲットの出力シーケンスへのマッピング処理 2 2 0 0 は間接が許されるかどうかの判断 2 2 0 2 から始まる。間接が許されなければ、ルックアップ・ターゲットは c h a r s i z e の断片で出力シーケンスにコピーされる ( 2 2 0 4 )。c h a r s i z e とは、ターゲット・コード化内の文字の最小サイズのことであり、マッピング・テーブル 4 1 4 のヘッダ 5 0 0 に指定されている。ブロック 2 2 0 4 に続いて、処理 2 2 0 0 は完了し、結果が見つかったとの通知を出して戻る ( 2 2 0 6 )。他方、間接が許されると判断 2 2 0 2 が判断したときは、ルックアップ・ターゲットの上位ビットが 1 に等しいかどうかの判断 2 2 0 8 が行われる。上位ビットが 1 に等しくなければ、第 1 バイトがヌル出力シーケンス (つまり、長さが 0 の出力シーケンス) を示しているかどうかの判断 2 2 1 0 が行われる。第 1 バイトが x 7 F を示していれば、ヌル出力シーケンスが示される ( 2 2 1 2 ) (一致するものが見つかったが、文字は出力シーケンスに追加されない)。他方、第 1 バイトがヌル出力シーケンスを示していないときは、c h a r s i z e の断片が出力シーケンスにコピーされる ( 2 2 1 4 )。その場合には、第 1 バイトは出力シーケンスの長さを示していることが好ましい。ブロック 2 2 1 2 と 2 2 1 4 に続いて、ブロック 2 2 0 6 が実行され、これにより処理は完了し、結果が見つかったことを通知して戻る。

#### 【 0 1 8 2 】

しかるに、ルックアップ・ターゲットの上位ビットが 1 に等しいと判断 2 2 0 8 が判断した場合には、ルックアップ・ターゲットは所望の出力シーケンスを間接的に参照しているので追加の処理が必要になる。具体的には、間接シーケンスまでのオフセットがルックアップ・ターゲットの残余部分を使用して指定される ( 2 2 1 6 )。ルックアップ・ターゲットの残余部分の上位ビットが 1 に等しいかどうかの判断 2 2 1 8 が行われる。等しくなければ、c h a r s i z e の断片が出力シーケンスにコピーされ ( 2 2 2 0 )、そのあと処理 2 2 0 0 はマッピングが見つかったとの通知を戻して ( 2 2 0 0 ) 完了する。他方、ルックアップ・ターゲットの残余部分の上位ビットが 1 に等しいと判断されていれば ( 2 2 1 8 )、順次チェーン内のシーケンスのカウントが取得され ( 2 2 2 4 )、リニア・サーチがシーケンスを通るように行われ、マッピングが特定される。ブロック 2 2 2 6 に

10

20

30

40

50

続いて、前述したブロック 2 2 2 0 と 2 2 2 2 が実行される。

【 0 1 8 3 】

図 2 6 C は図 2 6 B のブロック 2 2 2 6 によって実行されるオペレーションの詳細図である。すでに説明したように、図 2 6 C のブロック 2 2 2 6 はシーケンスを通るリニア・サーチを実行する。最初に、現在のカウン트가ゼロにセットされる ( 2 2 2 8 )。現在のカウン트가シーケンス・カウンタより大または等しいかどうかの判断 2 2 2 8 が行われる。現在のカウン트가シーケンス・カウンタより大または等しければ、ルックアップ・ターゲットの出力シーケンスへのマッピング処理 2 2 0 0 はマッピングが見つからなかったとの通知と共に戻る ( 2 2 3 2 )。他方、現在のカウン트가シーケンス・カウンタより大または等しくなければ、次の出力シーケンスとそのシーケンス・マスクが取得される ( 2 2 3 4 )。シーケンス・マスクは複数の出力シーケンスの 1 つを選択するために使用されるマスクである。次に、サブセット・フラグが使用中であるかどうかの判断 2 2 3 6 が行われる。サブセット・フラグが使用中でなければ、シーケンス・マスクと論理 AND がとられた実際属性が実際属性にビットワイズに等しいかどうかの判断 2 2 3 8 が行われる。そうであれば、マッピングが見つかったので処理は上述したブロック 2 2 2 0 と 2 2 2 2 を実行する。他方、サブセット・フラグが使用中であると判断 ( 2 2 3 6 ) されたときは、シーケンス・マスクと論理 AND がとられた選択フラグがビットワイズにシーケンス・マスクと等しいかどうかの判断 2 2 4 0 が行われる。そうであれば、マッピングが見つかったので前述したブロック 2 2 2 0 と 2 2 2 2 が実行される。従って、ブロック 2 2 3 8 と 2 2 4 0 は異なった方法で正しい出力シーケンスを取得する。他方、どちらかの判断ブロック 2 2 3 8 または 2 2 4 0 が現在のシーケンスが正しいマッピングでないと示していれば、現在のカウン트는インクリメントされ ( 2 2 4 2 )、処理はブロック 2 2 3 0 に戻ってチェーン内の次のシーケンスのオペレーションを続ける。

【 0 1 8 4 】

図 2 7 は本発明によるフォールバック・ハンドリング処理 2 3 0 0 を示すフローチャートである。フォールバック・ハンドリング処理 2 3 0 0 はフォールバック・ハンドラ 4 1 6 によって実行され、図 7 に示すユニコード・コンバータ処理 7 0 0 のブロック 7 1 6 によって起動される処理である。

【 0 1 8 5 】

フォールバック・ハンドリング処理 2 3 0 0 はフォールバック・オプションを使用してテキスト要素を調べる ( 2 3 0 2 )。ルックアップ 2 3 0 2 は図 1 9 を参照して上述したテキスト要素ルックアップ処理 1 6 3 0 に類似している。唯一の重要な違いは、変化した選択フラグを通して追加サブセットが考慮されるようにフォールバック・オプションがセットされていることである。次に、テキスト要素の変換コードが見つかったかどうかの判断 2 3 0 4 が行われる。変換コードが見つかったら、変換またはマッピングを取得するときにどのフォールバックが使用されたかを示すようにエラー・コードがセットされる ( 2 3 0 6 )。ブロック 2 3 0 6 に続いて、フォールバック・ハンドリング処理 2 3 0 0 は完了し、戻る。

【 0 1 8 6 】

他方、判断 2 3 0 4 が変換コードが見つからなかったと示しているときは、フォールバック・オプションに基づいてスイッチ・オペレーション 2 3 0 8 が実行される。フォールバック・オプションには、次のものがある。すなわち、デフォルト、呼び出し側定義、呼び出し側定義を伴うデフォルト、またはデフォルトを伴う呼び出し側定義である。フォールバック・オプションがデフォルトであれば、スイッチ・オペレーション 2 3 0 8 はデフォルト処理 2 3 1 0 を実行させる。フォールバック・オプションが呼び出し側定義であれば、スイッチ・オペレーション 2 3 0 8 は呼び出し側定義のオペレーション 2 3 1 2 を実行させる。フォールバック・オプションが呼び出し側定義を伴うデフォルトであれば、スイッチ・オペレーション 2 3 0 8 はデフォルト処理を実行させ ( 2 3 1 4 )、続いて判断 2 3 1 6 と呼び出し側定義の処理を実行させる ( 2 3 1 8 )。判断 2 3 1 6 はデフォルト処理が正常に行われると、呼び出し側定義処理 2 3 1 8 をバイパスするように働く。フォ

ールバック・オプションがデフォルト処理を伴う呼び出し側定義であれば、スイッチ・オペレーション 2308 は呼び出し側定義処理を実行させ (2320)、続いて判断 2322 とデフォルト処理を実行させる (2324)。判断 2322 は、呼び出し側定義処理 2320 が正常に行われたときは、デフォルト処理 2324 をバイパスするように働く。スイッチ・オペレーション 2308 に関連する処理に続いて、フォールバック処理 2300 がマッピングまたは変換コードを正常に特定していたかどうかの判断 2326 が行われる。フォールバック処理 2300 が失敗していれば、文字セットのデフォルト・フォールバック文字シーケンスが取得される (2328)。デフォルト・フォールバック文字シーケンスとは、フォールバック・ルックアップ 2302 が変換コードを特定するのに失敗したとき使用される変換コードである。好ましくは、デフォルト・フォールバック文字シーケンスはマッピング・テーブル 414 のヘッダに収められている。例えば、ASCII では、デフォルト・フォールバック文字シーケンスは“?”であるのが普通である。そのあと、ブロック 2328 に続いて、またはフォールバック処理がマッピングまたは変換コードを取得するのに成功したときはブロック 2326 に続いて、フォールバック・オプションが使用されたことを示すエラー・コードがセットされ (2306)、フォールバック・ハンドリング処理 2300 は完了し、戻る。

【0187】

図 28 はデフォルト処理 2400 を示すフローチャートである。デフォルト処理は図 27 のブロック 2310、2314 および 2324 で実行される処理と関連している。

【0188】

デフォルト処理 2400 は最初に現在のカウンタをゼロにセットする (2402)。次に、現在のカウンタがテキスト要素の長さより大であるか等しいかの判断 2404 が行われる。そうであれば、デフォルト処理 2400 は完了し、戻る。そうでなければ、フォールバック・フラグがセットされた単一ユニコード文字に対してルックアップ処理が実行される (2406)。ここでは、ルックアップはテキスト要素の個別文字に対するものであるが、以前 (ブロック 3202) では、ルックアップはテキスト要素全体に対するものであった。そのあと、単一ユニコード文字の変換コードが見つかったかどうかの判断 2408 が行われる。見つからなければ、ユニコード文字で利用できる個別マッピングがなかったことをエラー・コードで通知してデフォルト処理 2400 は戻る (2410)。他方、変換コードが見つかったら、現在のカウンタがインクリメントされ (2412)、処理はテキスト要素内の次のユニコード文字の処理を行うためにブロック 2404 に戻る。

【0189】

上述した図 4 ~ 図 28 はユニコードからターゲット・コード化への変換 (ユニコードから) に関するものであるが、上述したように、ユニコード・コード変換システム 300 は異なるソース・コード化からユニコードに変換する (ユニコードへ) 機能も等しく備えている。ユニコードへはユニコードからの処理と類似しているが、実質的には複雑度が軽減されている。ユニコードへの処理はターゲット・コード化を判断するときテキスト要素を探すためにスキャンしたり、複数の文字シーケンスを調べたりする必要がないのが通常である。ユニコードへの処理はソース文字列を個々の文字に分割し、そのあとでユニコードの対応するコード・ポイントを見つけるだけである。しかし、文字のマッピングがそのあとに続く文字に影響されるようなまれな場合には (例えば、デーヴァナーガリーのようなインド・スクリプト)、上述したようにスキャンが行われる場合がある。

【0190】

図 29 は本発明によるユニコード・コード変換システム 2500 の実施例を示すブロック図である。ユニコード・コード変換システム 2500 はユニコードへの変換を行う (つまり、ユニコードへの処理)。ユニコード・コード変換システム 2500 はユニコードへのコンバータ 2502 を含み、このコンバータはソース文字列 2504 を受信し、ユニコード文字列 2506 を出力する。ユニコードへのコンバータ 2502 はスキャナ 2508 とやりとりするユニコードへのコンバータを通してコード変換プロセスを実行する。スキ

ャナ 2 5 0 8 はスキャナ・テーブル 2 5 1 0 を使用して、ソース文字列 2 5 0 4 をスキャンしソース文字列 2 5 0 4 を文字に分断する。ここでは、ユニコードからの場合と異なり、ソース文字列は個々の文字に分割されるだけである。次に、ユニコードへのコンバータ 2 5 0 2 はルックアップ・ハンドラ 2 5 1 2 を使用して個別文字を調べそのユニコード・コード化を取得する。ルックアップ・ハンドラ 2 5 1 2 はマッピング・テーブル 2 5 1 4 を使用してユニコードの文字を取得する。さらに、ユニコードへのコンバータ 2 5 0 2 はフォールバック・ハンドラ 2 5 1 6 を使用することもできる。フォールバック・ハンドラ 2 5 1 6 はマッピング・テーブル 2 5 1 4 と一緒に動作して、ルックアップ・ハンドラ 2 5 1 2 がユニコード文字を特定できなかった場合に、テキスト要素のフォールバック・マッピングとして使用できるターゲット・コード化内の 1 つまたは複数の文字を特定する。

10

#### 【 0 1 9 1 】

スキャナ 2 5 0 8、スキャナ・テーブル 2 5 1 0、ルックアップ・テーブル 2 5 1 2、マッピング・テーブル 2 5 1 4、フォールバック・ハンドラ 2 5 1 6、状態管理機構 2 5 1 8 は図 4 の対応するデバイスと類似しているが、実質的には複雑度が軽減されている。従って、これらのデバイスはユニコードへとユニコードからの両方の処理を実現するように設計することができる。さらに、与えられたコンピュータ・システムや他のエレクトロニック・デバイスがユニコードへとユニコードからの両方の処理を実行する機能を備えているときは、そのコンピュータ・システムや他のエレクトロニック・デバイスはハブとして動作することができる。このハブはユニコードがサポートする種々の各国語文字間で変換を行うように動作することができる。例えば、ソースの各国語文字セットはまずユニコードに変換され、次にターゲットの各国語文字セットに変換される。

20

#### 【 0 1 9 2 】

本発明の第 1 側面によれば、上述したテーブルは、好ましくは、上述したフォーマットを使用してテーブルにストアすべきデータの圧縮を行う。さらに、本発明の第 1 側面によれば、圧縮を重視している。例えば、圧縮すると、属性テーブルは 5 0 倍 ( 5 0 対 1 ) に小さくなる。しかし、本発明の第 1 側面によれば、この分野の精通者ならば理解されるように、テーブルのフォーマットは多数の形体をとることができるので、あるテーブルは他のテーブルよりも実現が容易になる。最も容易な実現はどの圧縮も使用しないテーブルであるが、そのためには最大のデータ記憶装置が必要になる。さらに、本発明の第 1 側面によれば、テーブルは上述したインプリメンテーションで使用されているが、これらのテーブルで引き起こされる振舞を直接にコーディングすることも可能である。しかし、直接コーディングはコード変換システムのオペレーションの変更を困難にするのに対し、テーブルを使用すると、そのような変更を取り入れるためにテーブルだけを変更するだけで済むのが普通である。

30

#### 【 0 1 9 3 】

本発明の多数の特徴および利点は上述してきた説明で明らかにした通りであり、本発明のかかる特徴と利点はすべて請求の範囲に記載されている。さらに、本発明はこの分野の精通者が容易に理解されるように、種々態様に変更することが可能であるので、本発明は上述してきた正確な構造とオペレーションに限定されるものではない。従って、すべての適当な変更および等価技術は本発明の範囲に属するものである。

40

#### 【図面の簡単な説明】

#### 【 0 1 9 4 】

【図 1】本発明による代表的なコンピュータ・システムを示すブロック図である。

【図 2】ユニコード文字コード化のフォーマットを示す図である。

【図 3】ソース文字列を受信し、ターゲット文字列を出力する本発明による基本的ユニコード・コード変換システムを示すブロック図である。

【図 4】本発明の実施例によるユニコードからのコード変換システムの実施例を示すブロック図である。

【図 5】ユニコード・コード変換システムのマッピング・テーブルの好ましい配列を示す概略図である。

50

【図 6 A】本発明の実施例によるユニコード・コード変換システムを利用するアプリケーション・プログラムによって実行される処理を示すフローチャートである。

【図 6 B】本発明の実施例による切り捨て処理を示すフローチャートである。

【図 7】本発明の実施例によるユニコード・コンバータ処理を示すフローチャートである。

【図 8】本発明の実施例による更新オフセット処理を示すフローチャートである。

【図 9 A】本発明の実施例による次のテキスト要素処理を示すフローチャートである。

【図 9 B】本発明の実施例による次のテキスト要素処理を示すフローチャートである。

【図 10 A】本発明の実施例による次のテキスト要素処理を示すフローチャートである。

【図 10 B】本発明の実施例による次のテキスト要素処理を示すフローチャートである。

【図 10 C】本発明の実施例による次のテキスト要素処理を示すフローチャートである。

【図 11】本発明の実施例によるスキャナを示すブロック図である。

【図 12】図 10 に示す属性テーブルの好ましいフォーマットを示す概略図である。

【図 13】本発明の実施例による属性ルックアップ処理を示すフローチャートである。

【図 14 A】本発明の好適実施例による次のアクションを判断するためにスキャナによって利用されるスキャナ・テーブルを示す概略図である。

【図 14 B】本発明の好適実施例による次のアクションを判断するためにスキャナによって利用されるスキャナ・テーブルを示す概略図である。

【図 15】本発明の実施例による好ましいレイアウトと、スキャナ・テーブルにストアされる情報の両方を表しているテーブルである。

【図 16 A】本発明の好適実施例による好ましいレイアウトと、スキャナ・テーブルにストアされる情報の両方を表しているテーブルである。

【図 16 B】本発明の好適実施例による好ましいレイアウトと、スキャナ・テーブルにストアされる情報の両方を表しているテーブルである。

【図 17 A】本発明の実施例による方向解明処理を示すフローチャートである。

【図 17 B】本発明による双方向状態テーブルの好ましいレイアウトを表しているテーブルである。

【図 17 D】本発明による双方向状態テーブルの好ましいレイアウトを表しているテーブルである。

【図 18】本発明の実施例による初期方向処理を判断するためのフローチャートである。

【図 19】本発明の実施例によるルックアップ・テキスト要素処理を示すフローチャートである。

【図 20】本発明の実施例によるチェーン・フォーマット処理を示すフローチャートである。

【図 21】本発明の実施例による範囲フォーマット処理を示すフローチャートである。

【図 22】本発明の実施例によるリスト・フォーマット処理を示すフローチャートである。

【図 23 A】本発明の実施例によるセグメント配列フォーマット処理を示す図である。

【図 23 B】本発明の実施例によるセグメント配列フォーマット処理を示す図である。

【図 24】本発明の実施例による変形リスト・サーチ処理を示すフローチャートである。

【図 25 A】本発明の実施例による変形リスト処理に関連するテーブルとデータ・フォーマットを示す概略図である。

【図 25 B】本発明の実施例による変形リスト処理に関連するテーブルとデータ・フォーマットを示す概略図である。

【図 26 A】本発明の実施例による出力シーケンスへのルックアップ・ターゲットのマッピング処理を示すフローチャートである。

【図 26 B】本発明の実施例による出力シーケンスへのルックアップ・ターゲットのマッピング処理を示すフローチャートである。

【図 26 C】本発明の実施例による出力シーケンスへのルックアップ・ターゲットのマッピング処理を示すフローチャートである。

10

20

30

40

50

【図 27】本発明の実施例による本発明の処理に従うフォールバック・ハンドリング処理を示すフローチャートである。

【図 28】本発明の実施例によるデフォルト処理を示すフローチャートである。

【図 29】本発明の実施例によるユニコードへのコード変換の実施例を示すブロック図である。

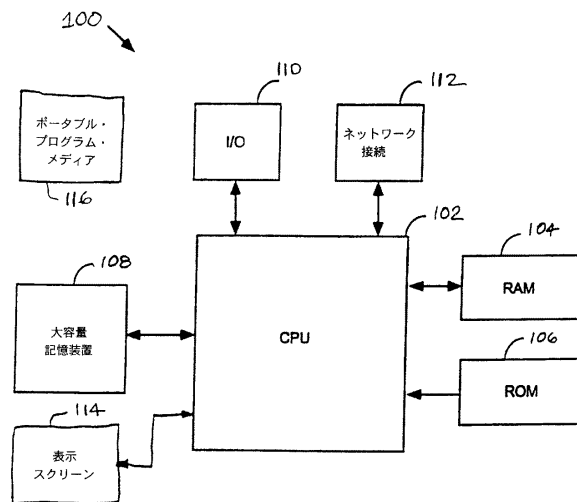
【符号の説明】

【0195】

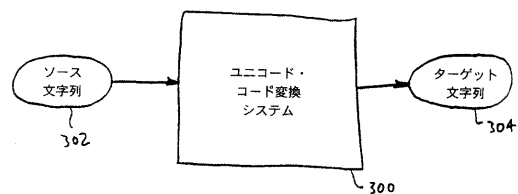
- 100 コンピュータ・システム
- 200 フォーマット
- 400 ユニコード・コード変換システム
- 414 マッピング・テーブル
- 1004 属性テーブル
- 1300 スキャナ・テーブル
- 1302 要素
- 1400 テーブル
- 1511 双方向状態テーブル
- 2100 変種リスト
- 2108 実際属性ビット・マスク
- 2110, 2112, 2114, 2116 部分

10

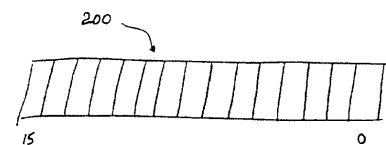
【図 1】



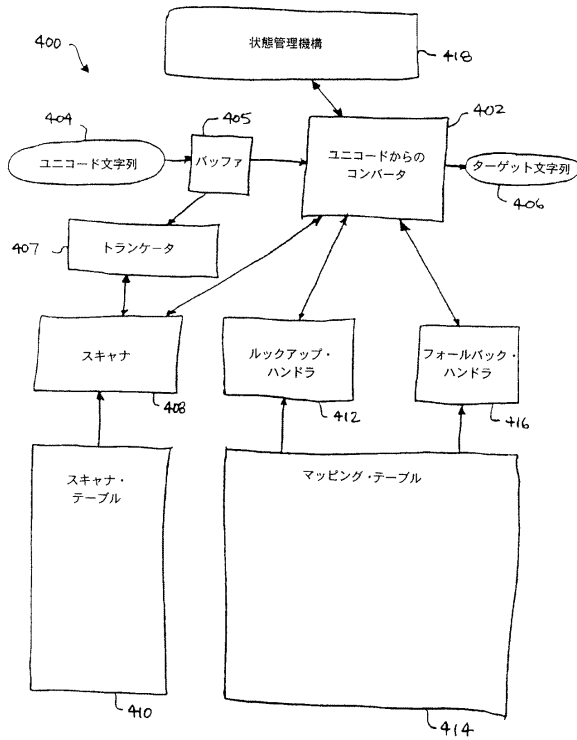
【図 3】



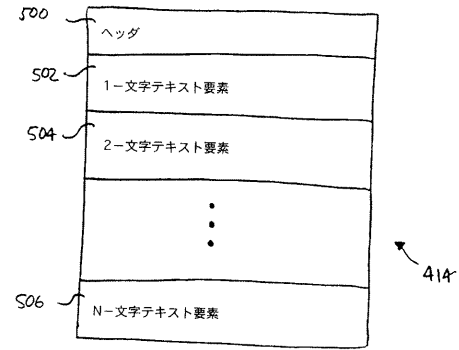
【図 2】



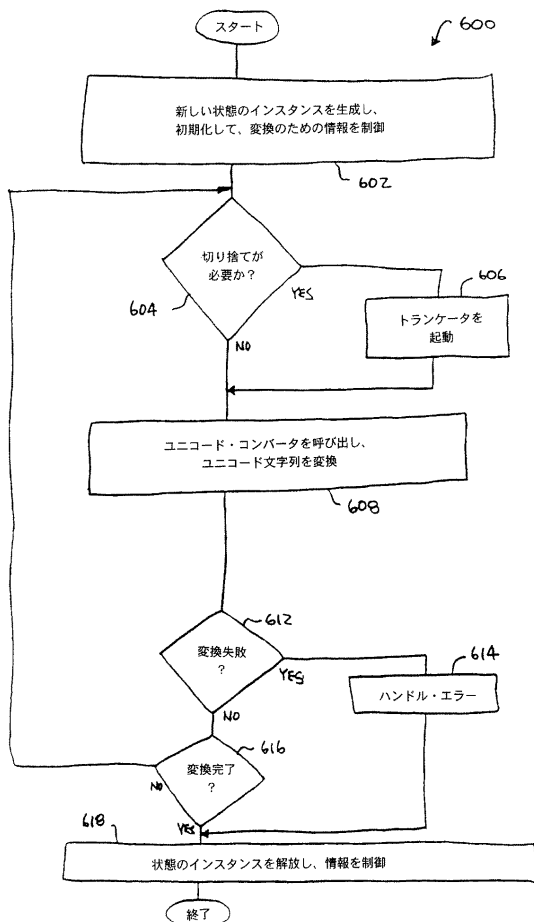
【図 4】



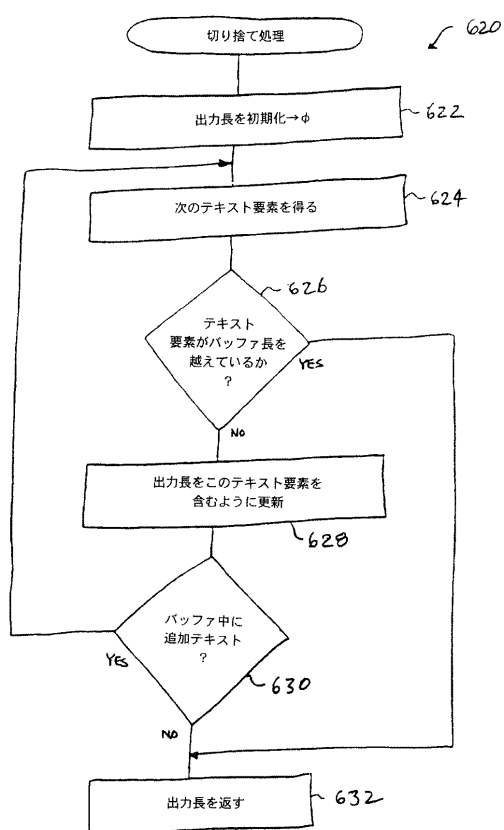
【図 5】



【図 6 A】

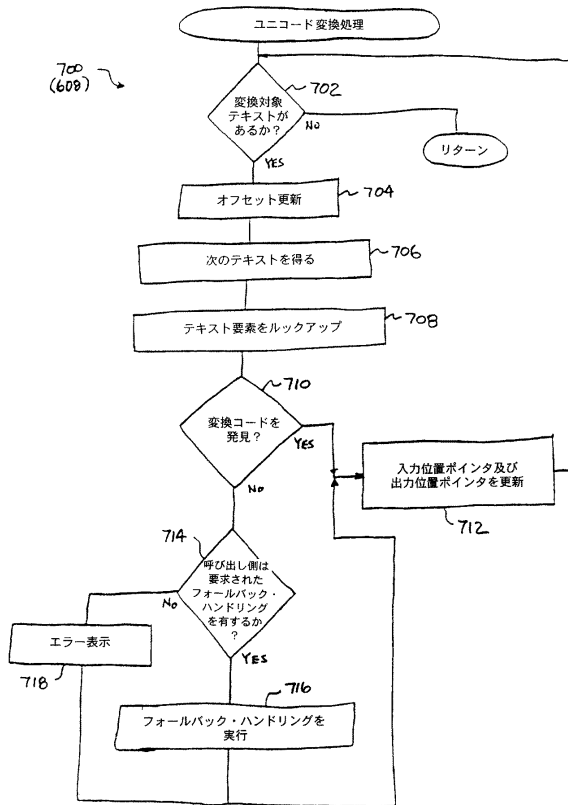


【図 6 B】

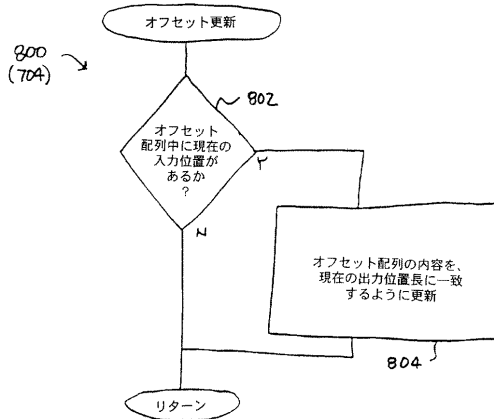




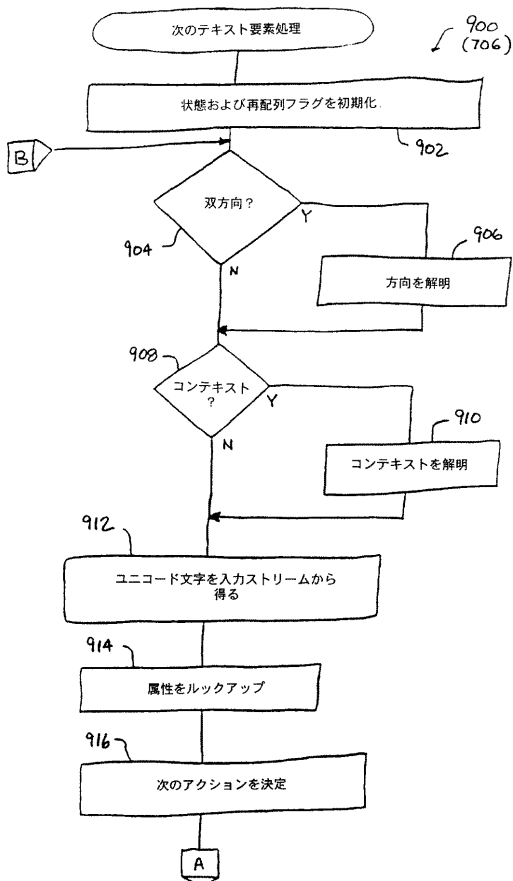
【図 7】



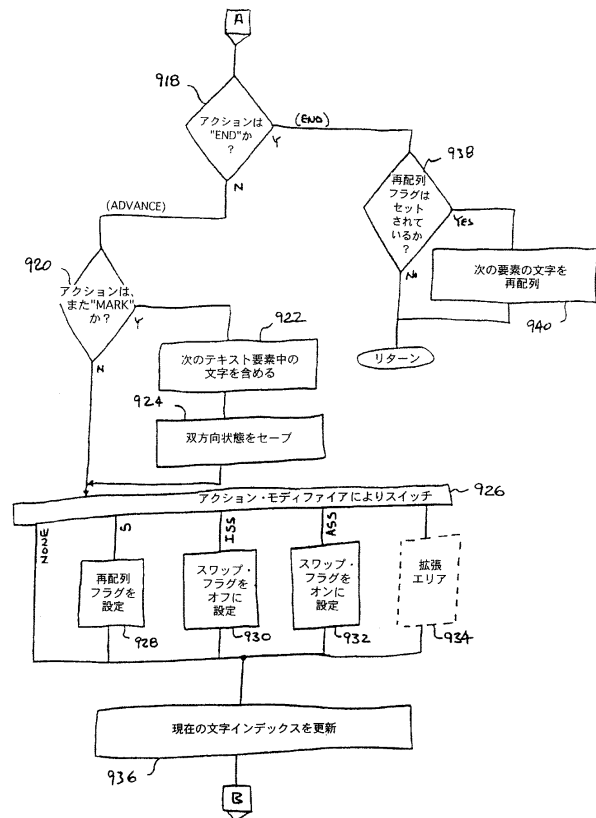
【図 8】



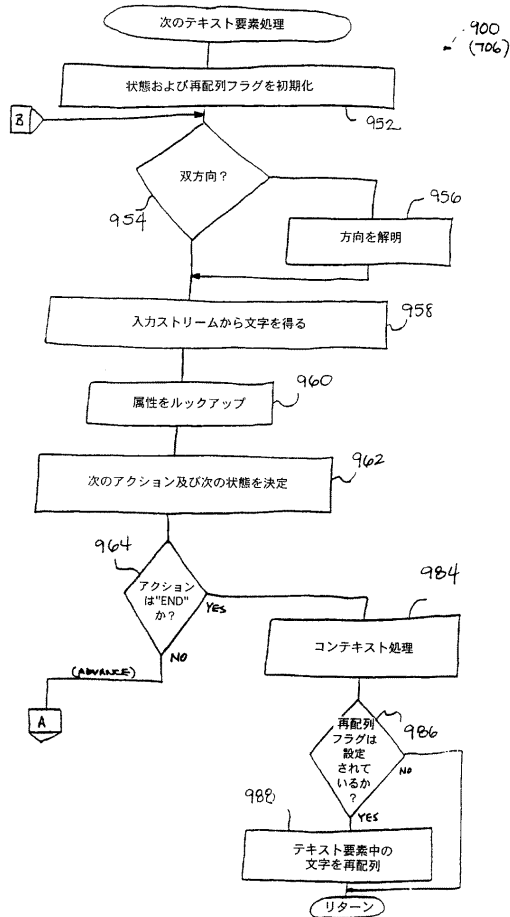
【図 9 A】



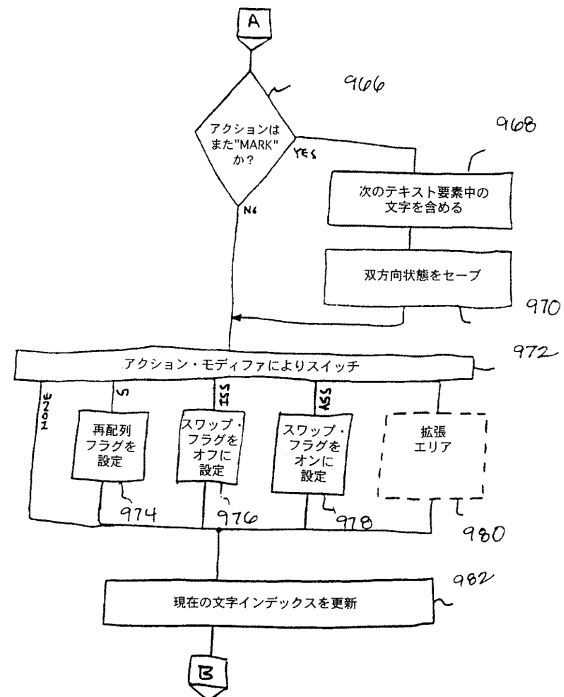
【図 9 B】



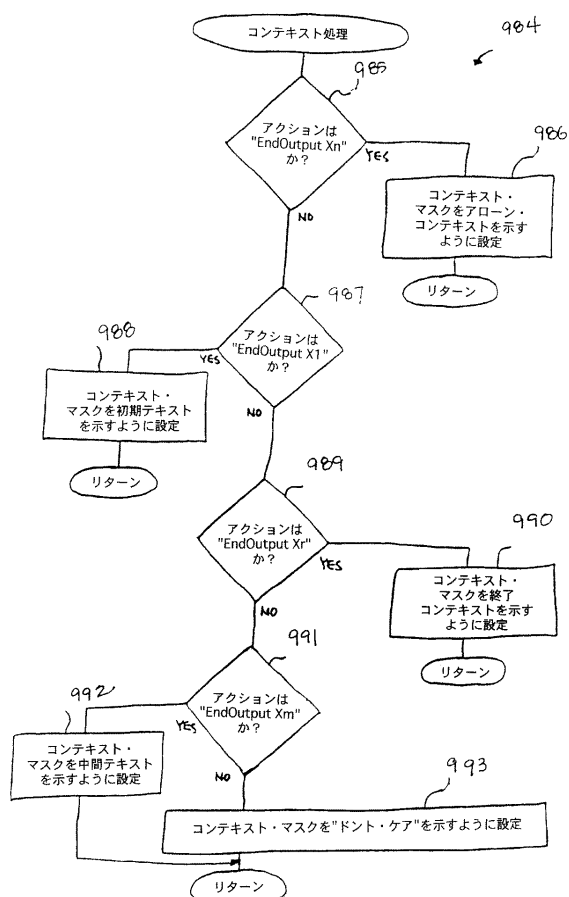
【図10A】



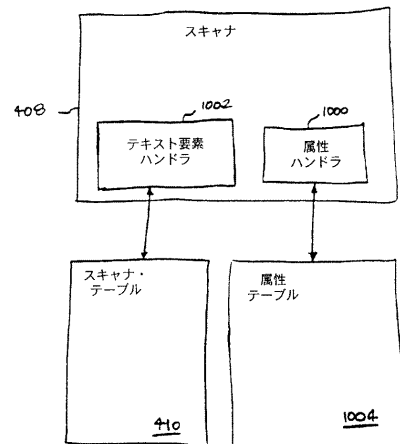
【図10B】



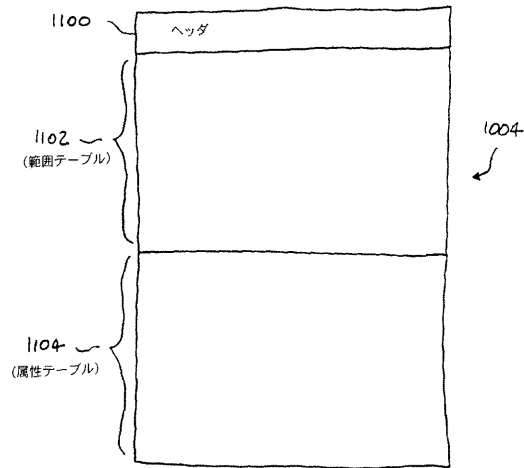
【図10C】



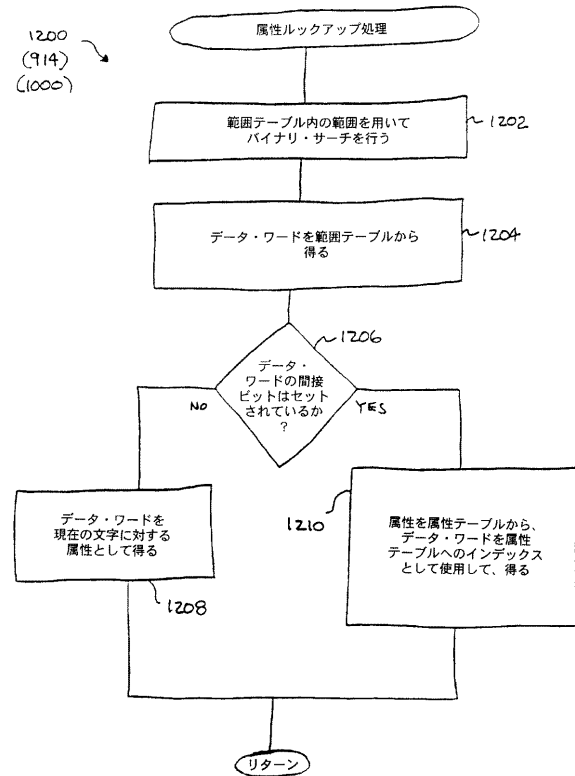
【図11】



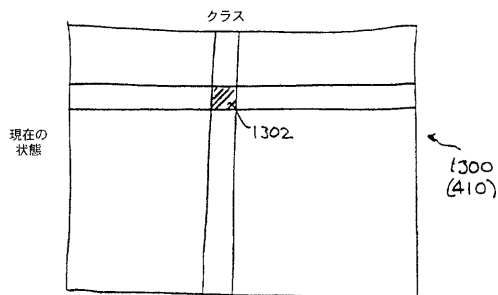
【 図 1 2 】



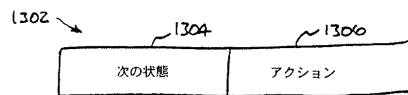
【 図 1 3 】



【 図 1 4 A 】



【 図 1 4 B 】



【 図 1 5 】

[illegible]

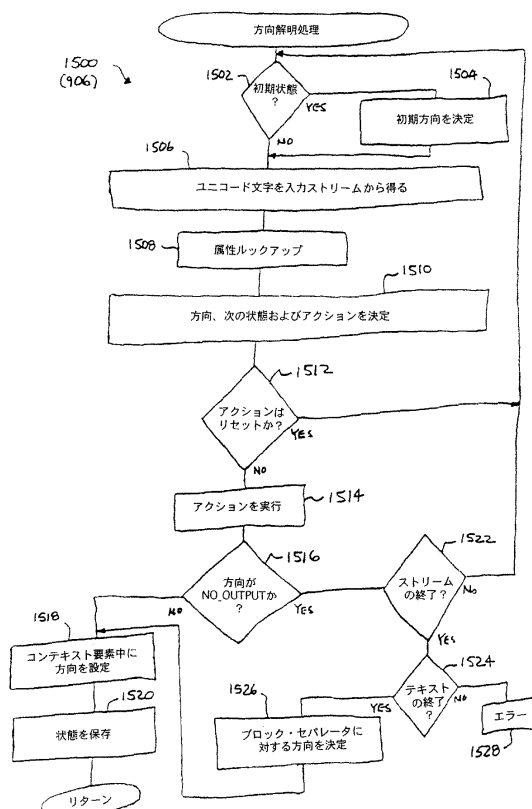
【 図 1 6 A 】

[illegible]

【 図 1 6 B 】

[illegible]

【 図 1 7 A 】



【 図 1 7 B 】

	START	sDIR	OR	SL	SR	AL	L+AN	L+EN
LR	START,L	START,L	OR,G	S,L	S,R	AL	S,L	L+EN
RL	S,R	START,R	OR,G	S,R	S,L	S,R	S,L	S,L
AL	ALL	START,R	OR,G	AL,R	AL,R	AL,R	AL,R	S,R
LRE	SL,push,L	START,L	SL,push,L	SL,push,L	SL,push,L	SL,push,L	SL,push,L	SL,push,L
RLR	SR,push,R	START,R	SR,push,R	SR,push,R	SR,push,R	SR,push,R	SR,push,R	SR,push,R
LRO	OR,push,L	START,L	OR,push,L	OR,push,L	OR,push,L	OR,push,L	OR,push,L	OR,push,L
RLO	OR,push,R	START,R	OR,push,R	OR,push,R	OR,push,R	OR,push,R	OR,push,R	OR,push,R
PDF	S(G),pop,G	error	S(G),pop,G	S(G),pop,G	S(G),pop,G	S(G),pop,G	S(G),pop,G	S(G),pop,G
AN	(G)+AN,L	sDIR,*	OR,G	L+AN,L	R+AN,L	AL+Num,L	S(G),pop,G	AL+G
EN	(G)+EN,L	sDIR,*	OR,G	L+EN,L	R+EN,L	AL+Num,L	L+EN,L	L+EN,L
ET	s#+ET,*	sDIR,*	OR,G	SL+*N	SR+ET,*	SL+*N	SL+*N	L+EN,L
ES	s(G)+N,*	sDIR,*	OR,G	SL+*N	SR+*N,*	SA+*N	SL+*N	SL+*N
CS	s(G)+N,*	sDIR,*	OR,G	SL+*N	SR+*N,*	SA+*N	SL+*N	SL+*N
ON	s(G)+N,*	sDIR,*	OR,G	SL+*N	SR+*N,*	SA+*N	SL+*N	SL+*N
BS	s(G),G	sDIR,*	reset	reset	reset	reset	reset	reset
	scan for base differentiation		override	strong left	strong right	arabic letter	strong left+	strong left+

17B	17C	17D
-----	-----	-----

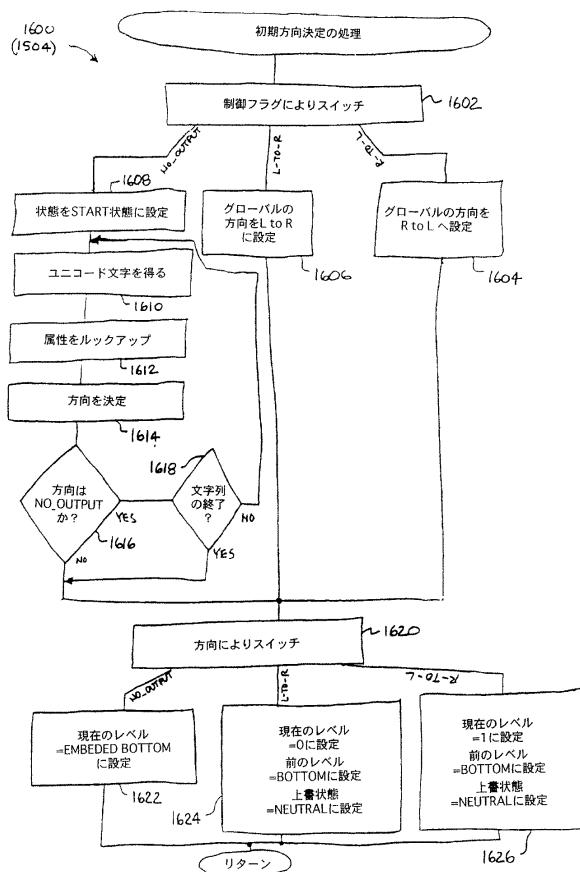
115

【 図 1 7 D 】

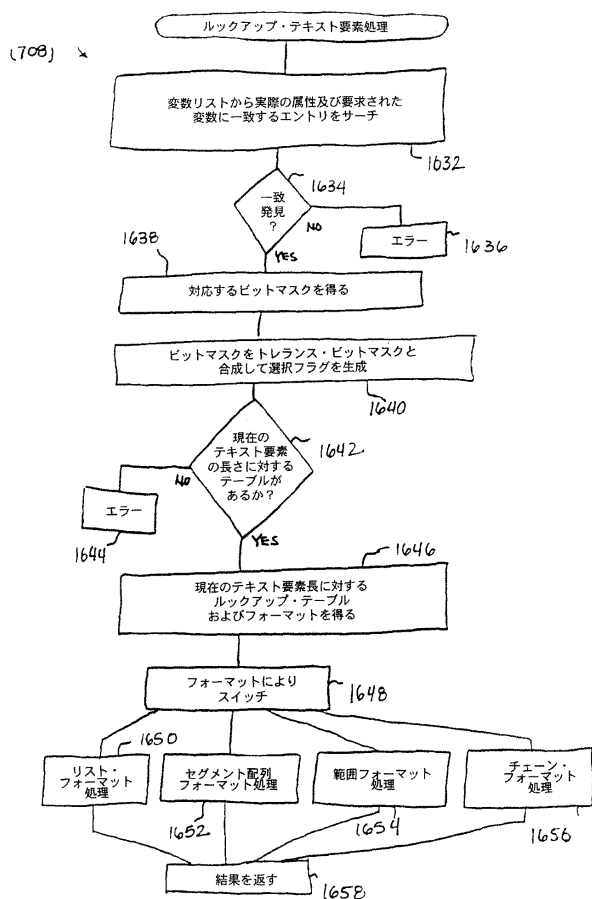
[illegible]

115

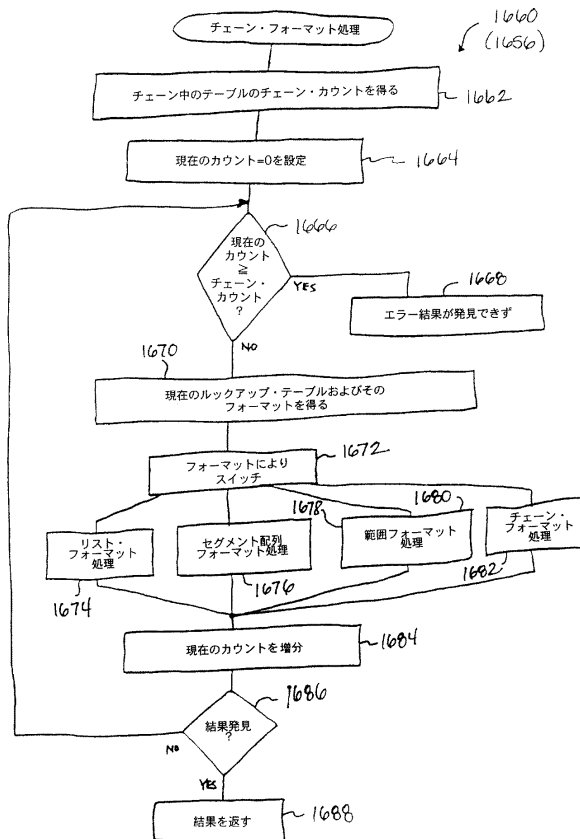
【 ㄨ 1 8 】



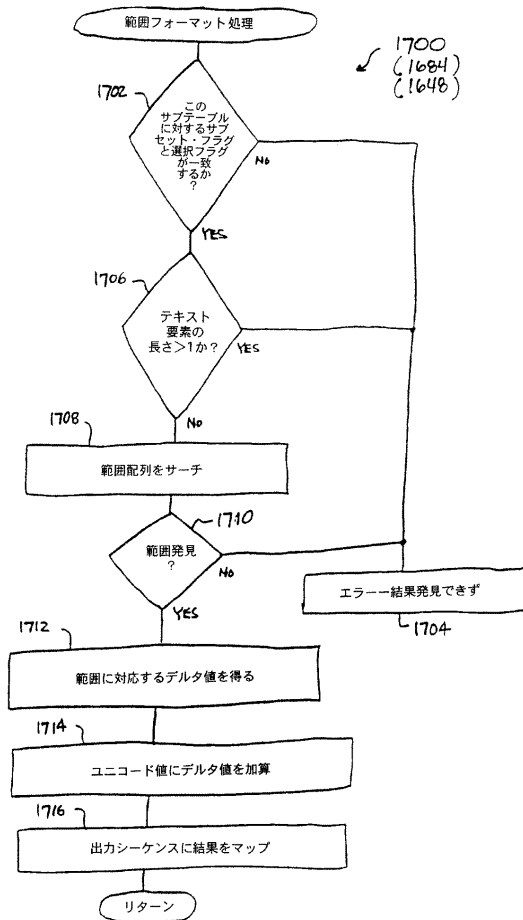
【 図 1 9 】



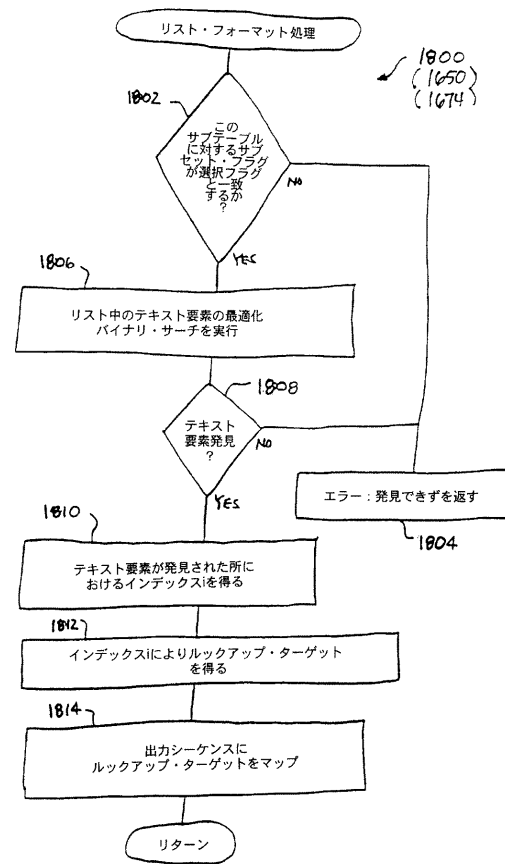
【 図 2 0 】



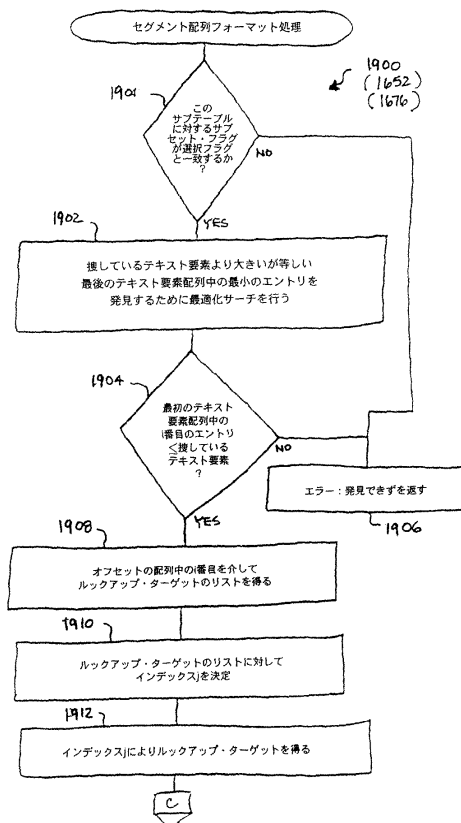
【図 2 1】



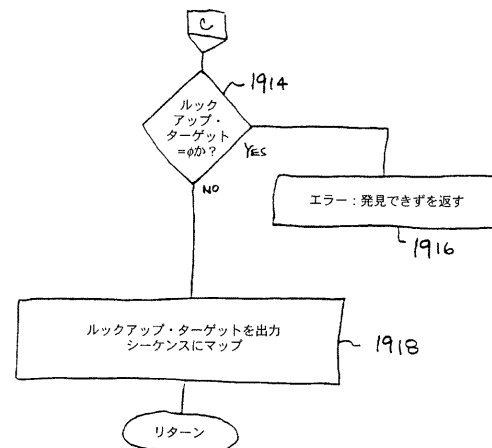
【図 2 2】



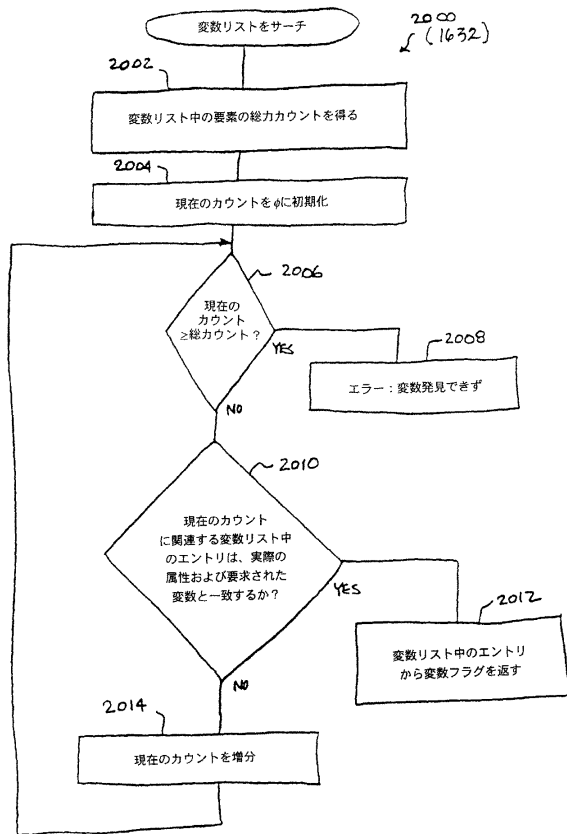
【図 2 3 A】



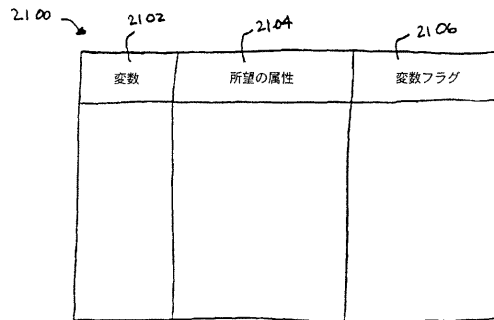
【図 2 3 B】



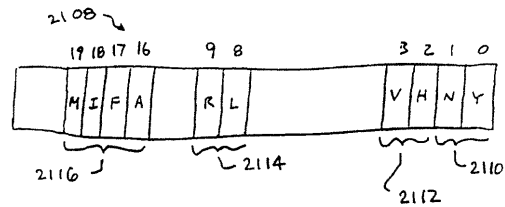
【図 24】



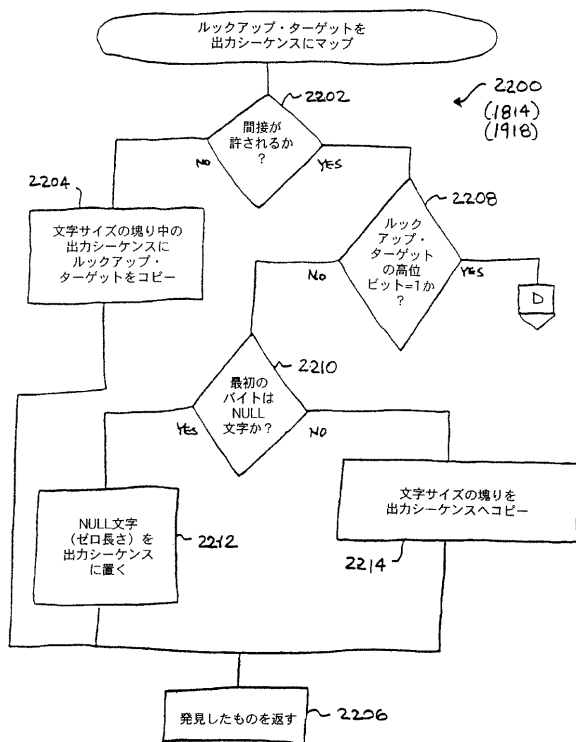
【図 25A】



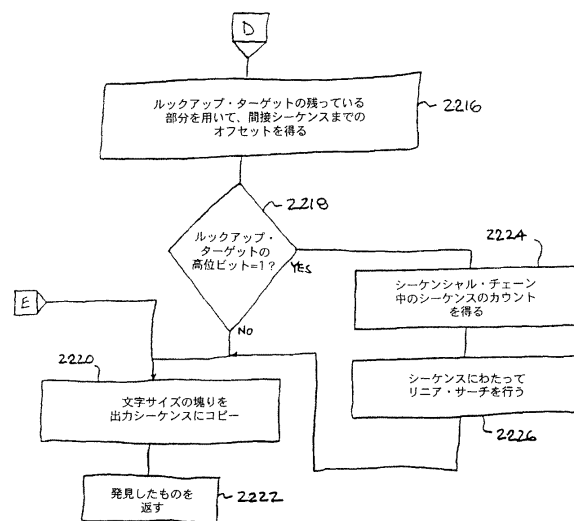
【図 25B】



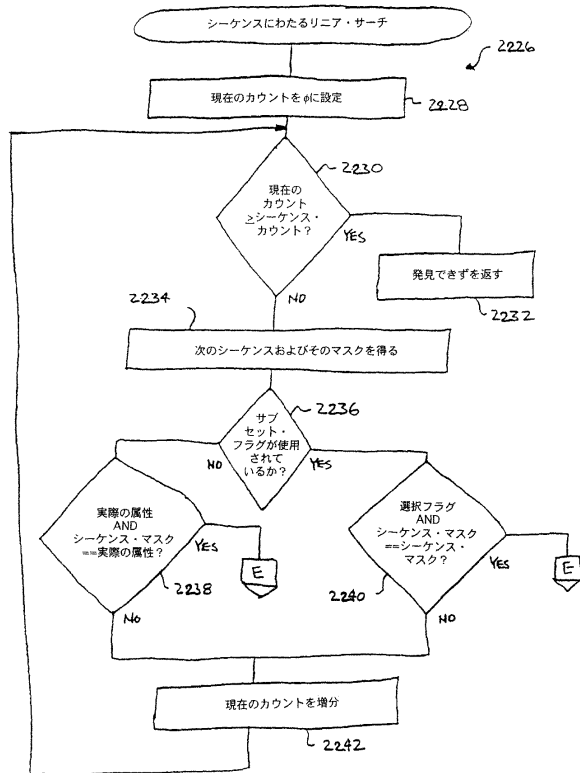
【図 26A】



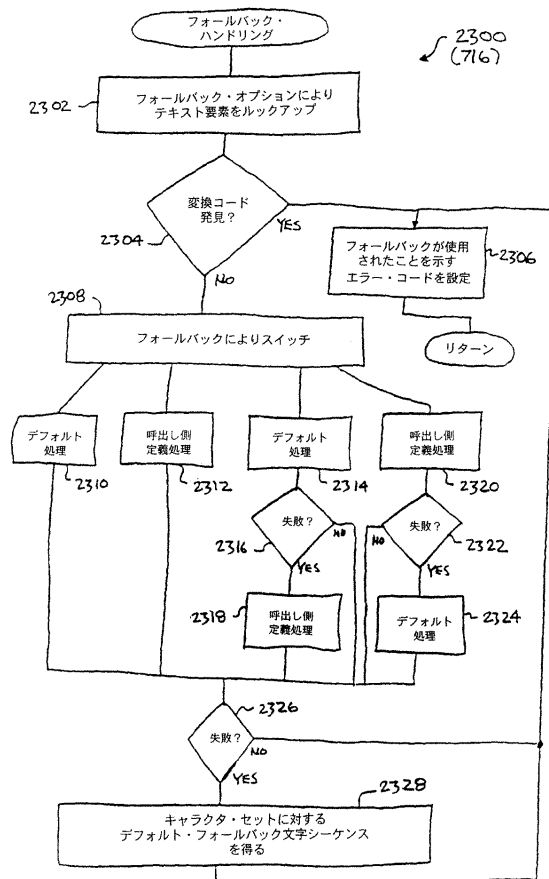
【図 26B】



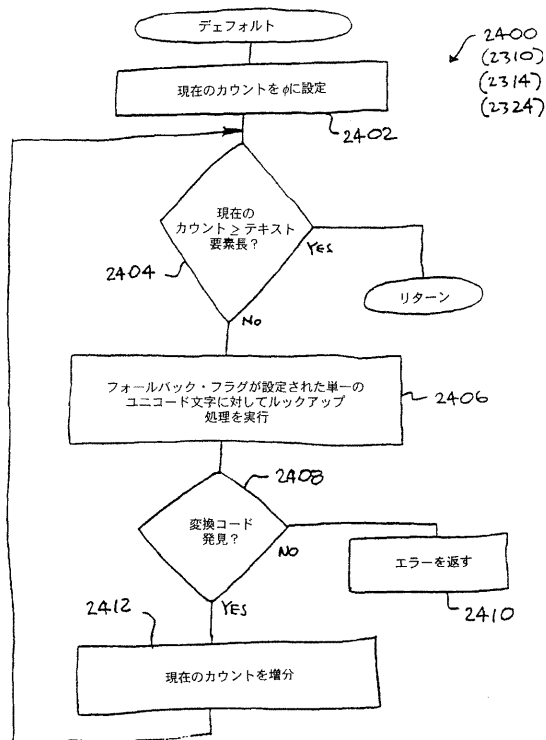
【図 26C】



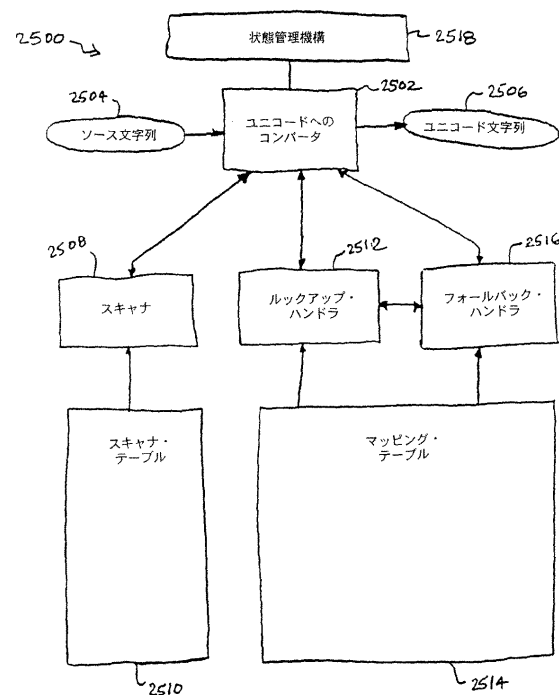
【図 27】



【図 28】



【図 29】





## フロントページの続き

- (31)優先権主張番号 08/527,837  
(32)優先日 平成7年9月13日(1995.9.13)  
(33)優先権主張国 米国(US)  
(31)優先権主張番号 08/527,831  
(32)優先日 平成7年9月13日(1995.9.13)  
(33)優先権主張国 米国(US)

## 早期審理対象出願

- (72)発明者 マコンネル, ジョン, アイ .  
アメリカ合衆国 9 4 0 2 5 カリフォルニア州 メンロ パーク マッケンドリー ドライブ  
2 2 2  
(72)発明者 タン, ユン - フォン, フランク  
アメリカ合衆国 9 4 0 8 2 カリフォルニア州 サニーヴェイル イースト ホームステッド  
ロード 1 7 5 アpartment ナンバー 8  
(72)発明者 ダニエルズ, アンドリュー, エム  
アメリカ合衆国 9 4 0 2 5 カリフォルニア州 メンロ パーク ハイト ストリート 1 2 3

## 合議体

審判長 田口 英雄  
審判官 飯田 清司  
審判官 真木 健彦

- (56)参考文献 特開平 6 - 8 3 5 7 1 ( J P , A )  
特開平 7 - 9 3 3 7 4 ( J P , A )  
特開平 4 - 2 3 8 5 3 1 ( J P , A )  
特開平 7 - 1 2 1 5 1 1 ( J P , A )  
特開平 8 - 1 6 3 6 0 ( J P , A )  
特開平 3 - 2 2 3 9 2 3 ( J P , A )