

(19) 日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11) 特許番号

特許第3101153号
(P3101153)

(45) 発行日 平成12年10月23日 (2000. 10. 23)

(24) 登録日 平成12年 8 月18日 (2000. 8. 18)

(51) Int.Cl. ⁷	識別記号	F I
G 0 6 K 9/68		G 0 6 K 9/68 G
G 0 6 F 17/22	5 2 0	G 0 6 F 17/22 5 2 0 B

請求項の数 1 (全 11 頁)

(21) 出願番号 特願平6-143583

(22) 出願日 平成 6 年 6 月24日 (1994. 6. 24)

(65) 公開番号 特開平8-16711

(43) 公開日 平成 8 年 1 月19日 (1996. 1. 19)

審査請求日 平成10年 1 月30日 (1998. 1. 30)

(73) 特許権者 000005049

シャープ株式会社

大阪府大阪市阿倍野区長池町22番22号

(72) 発明者 船山 竜士

大阪府大阪市阿倍野区長池町22番22号

シャープ株式会社内

(74) 代理人 100079843

弁理士 高野 明近

審査官 脇岡 剛

(56) 参考文献 特開 平4-151761 (J P, A)

特開 昭62-251986 (J P, A)

(58) 調査した分野(Int.Cl.⁷, D B名)

G06K 9/68

G06F 17/22

(54) 【発明の名称】 日本語入力装置

(57) 【特許請求の範囲】

【請求項 1】 読みのわからない文字を入力する際に、形の似ている文字で読みのわかっている文字をキー文字入力する漢字入力手段と、オンライン手書き文字認識用辞書及び光学文字認識用辞書を利用して、似ている文字群と該文字群の各文字の類字度を集めて自動生成される類字辞書と、前記漢字入力手段により入力されたキー文字に基づいて前記類字辞書から形の似ている文字を検索し、類字度を評価する類字検索手段とを備えた日本語入力装置であって、前記漢字入力手段から見方の異なる複数のキー文字を入力し、該キー文字ごとに前記類字辞書から形の似ている文字を検索することにより目的とする文字の候補を絞り込み、キー文字と類字度の高い順に候補文字を表示させることを特徴とする日本語入力装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、日本語入力装置に関し、より詳細には、日本語入力可能な機器において、読みのわからない文字を入力する際、形の似ている文字を集めた辞書を用いて、その字に形が似ている文字をまず入力してから前記辞書を用いて目的の字を検索するようにした日本語入力装置に関する。例えば、ワープロ、パソコン、オフコン、ワークステーション、電子手帳など、辞書を検索する方式を用いた日本語入力方式を持つすべての機器に適用されるものである。

【0002】

【従来の技術】現在主流になっているかな漢字変換は、読みをかな入力し、入力されたかなを変換することによって目的の字を得ようになっているが、読みのわからない文字を入力する方法としては、コード入力や部首入

(2)

力、画数入力などがある。コード入力は、図25に示すように、字や記号と1:1に対応したコード(JIS、シフトJIS、区点コードなど)を、通常は紙に印刷された一覧表(通常、第1水準の漢字は音読みでのよみがなの50音順、第2水準は部首の画数順の配列になっている)から拾い出し、そのコードを入力して変換することによって目的の字を得るものである。一覧表が紙に印刷されておらず、電子化されている場合もあり、その場合は、コードを入力するのではなく、ディスプレイ上で紙の上と同じように目的の字を探し出し、ポインティングデバイスなどを用いてその字を指定することもある。

【0003】部首入力は、図26に示すように、目的の字がどの部首に属するかを、通常は紙に印刷された一覧表から拾い出し、そこに書いてあるコードを入力して目的の字を得るものである。一覧表が紙に印刷されておらず、電子化されている場合もあり、その場合は、コードを入力するのではなく、ディスプレイ上で紙の上と同じように目的の字を探し出し、ポインティングデバイスなどを用いてその字を指定することもある。

【0004】画数入力は、文字の総画数が同じものを集め、画数の順にその漢字群を配列した表から、目的の文字を目視によって検索するものである。紙に印刷された一覧表から、該当する文字のコードを探し出し、そのコードを入力して目的の文字を得る場合と、一覧表が電子化されており、ディスプレイ上で紙の上と同じように目的の文字を探し出し、ポインティングデバイスなどを用いてその文字を指定する場合とがある。

【0005】従来の日本語入力装置について記載した公知文献としては、例えば、特開平2-18661号公報がある。この公報のものは、音訓入力手段と部首入力手段と画数入力手段の少なくとも1つの入力手段と、字形構成要素辞書及び字形構成要素入力手段とによって、それぞれの辞書の共通をとって候補文字を表わすもので、候補対象漢字の数を減らすことができるというものである。すなわち、読みや部首、画数の情報を複合して、複数の条件を設定し、目的の文字を絞り込むというものである。

【0006】また、特開平5-12249号公報のものは、文書作成時に画面に表示された漢字パターンを、例えばカーソルで指定すると、部首解析処理手段は、部首解析用テーブルに基づいて、指定された漢字の部首名及びその部首の位置を解析し、検索出力手段は、部首解析処理手段による解析結果に基づいて、指定された漢字の部首名及びその部首の位置に対応する全ての漢字コードを部首辞書メモリ手段から検索する。これにより、部首名がわからない場合でも、同一の部首名を有する漢字を指定するだけで、目的の漢字を得ることができる。また、目的の漢字を構成する部首名の一部しかわからない場合でも、その部首名に基づいて目的の漢字を得ることができるものである。

【0007】

【発明が解決しようとする課題】パソコンや電子手帳などの日本語入力が可能な情報機器などにおいて、通常は読みをかなで入力し、入力されたかなを漢字などに変換する方式などを用いているが、それだけでは、読みのわからない字や記号を入力することができない。そのような文字を入力する場合、従来の方法では、総画数順に漢字が配列された表を、紙上あるいは電子の形で持ち、画数を手がかりに目的の字をその表から目視によって検索していた。あるいは、部首ごとに漢字を配列した表を、紙上あるいは電子の形で持ち、部首を手がかりに目的の字を、目視によって検索していた。通常は、部首の画数順に同じ部首を持った漢字群が配列され、同じ部首を持った漢字群の中では、総画数の順に字が配列されている。

【0008】しかし、総画数より検索する方式だと、同じ画数を持つ漢字は膨大な数に昇り、そこからひとつの文字を探し出すのは、非常に困難であり、また、画数のはっきりと特定できない文字も存在するため、検索はさらに困難を極める。また、部首より検索する方式だと、目的の漢字の部首がきちんとわかっていなければならないことが条件となる。漢字によっては、部首が何になるのかわからないものもあり(例えば、彗星の「彗」の部首は「彗」である)、その場合、目的の字を探し出すのは非常に困難である。また、多くの部首では、同じ部首を持つ文字の数が、非常に多くなり、同様に検索には困難が伴う。

【0009】また、コード入力の場合、一覧表は通常、第一水準ならば音読みでのよみがなの50音順に配列されており、もとより読みのわからない字を入力しようとしているわけであるから、これでは役に立たないことがわかる。第2水準ならば部首の画数の小さい順に配列されており、目的の字の画数を数えるという手間がかかる上に、部首のわからない文字を探すには相当手間がかかる。部首入力の場合は、目的の字の部首がわからない場合は探すのに大変な手間がかかるし、同じ部首を持つ字が多く存在する場合の検索も時間がかかる。総画数入力の場合は、目的の文字の画数を数えるのがまず面倒であり、多画の文字であれば、画数の数え間違いなどが発生して検索の効率が悪い。また、同じ画数の文字も多く存在するので、目的の文字を探すには手間がかかる。

【0010】前述した特開平2-18661号公報に示されている方法を用いた場合、「読み」はわからないわけであるから、部首と画数から目的の文字を検索するわけであるが、部首がわからないものに関しては、総画数入力と同じ条件になるし、部首がわかって、文字の画数を数えるという作業は、非常に根気のいる作業であり、あまり、効率のいい検索方法であるとは言えない。一覧表が電子化されている場合も同様で、ディスプレイ上で表示できる字の数は、紙の上で表示できる数より少

(3)

ない場合が多く、その場合、目的の字を探すにはさらに時間がかかることになる。

【0011】また、類字変換を実現するための類字辞書を人力で作成することにおいて、形の似ている字を人手を使って集め、それらの類字度を主観によって評価するとなると、似ている字がまだあるにもかかわらず辞書から洩れてしまったり、類字度の評価が類字グループで統一されていなかったり、また、使用者と辞書作成者で類字の評価に違いがある場合など、使用感が著しく損なわれる可能性がある。もとより、人手を使ってこの作業を行うとなると、それにかかる労力・時間（コスト）は大変なものとなる。

【0012】また、前述したように、特開平5-12249号公報のものは、同じ部首を持つ漢字から目的の字を検索するだけのものである。本発明はあくまでも、読みのわからない字を入力するために、「形の似ている字（＝類字）」をキーにして、類字辞書を引き、目的の字を検索するものである。本発明においても、共通に持つ部首をキーとして目的の文字を検索する手段を有しているが、それはあくまでも一つの手段であり、本発明では、部首が違っていても、形が似ている字（「瓦」と「互」など）も、検索の対象としている。

【0013】また、特開平5-12249号公報においては、ただ単に、部首の同じ字を辞書から拾ってくるだけであるが、本発明では、「類字度」という概念を導入し、より似ている字を候補の先に配列することができるようになっている。例えば、特開平5-12249号公報では、「苑」を入力しようとして、「苑」をキーにしても、通常では、「うかんむり」の文字がたくさん出てきてしまい、効率的に「苑」を見つけることはできない。また、部首の位置を指定すれば、もう少し候補を絞ることができるかもしれないが、この作業自体が非効率的である。そこで、本発明における辞書の自動生成方法に従って作成された辞書を用いて類字検索を行えば、特開平5-12249号公報より遥かに効率的に、読みのわからない文字を検索することができる。

【0014】本発明は、このような実情に鑑みてなされたもので、日本語入力可能な機器において、読みのわからない文字を入力する際、形の似ている文字を集めた辞書を用いて、その字に形が似ている文字をまず入力し、前記辞書を用いて目的の字を検索するようにした日本語入力装置を提供することを目的としている。

【0015】

【課題を解決するための手段】本発明は、上記目的を達成するために、読みのわからない文字を入力する際に、形の似ている文字で読みのわかっている文字をキー文字入力する漢字入力手段と、オンライン手書き文字認識用辞書及び光学文字認識用辞書を利用して、似ている文字群と該文字群の各文字の類字度を集めて自動生成される類字辞書と、前記漢字入力手段により入力されたキー文

字に基づいて前記類字辞書から形の似ている文字を検索し、類字度を評価する類字検索手段とを備えた日本語入力装置であって、前記漢字入力手段から見方の異なる複数のキー文字を入力し、該キー文字ごとに前記類字辞書から形の似ている文字を検索することにより目的とする文字の候補を絞り込み、キー文字と類字度の高い順に候補文字を表示させる日本語入力装置であることを特徴とする。

【0016】

【作用】前記構成を有する本発明の日本語入力装置は、読みのわからない文字の検索を簡単に行うことができる。すなわち、ほとんどの漢字は、それほど多くない、いくつかの部品から成ると考えられる。例えば、「腕」という字は、「月」「艹」「夕」「巳」という4つの部品から成る。また、「苑」という字は、

【0017】

【表1】

「艹」

【0018】「夕」「巳」の3つの部品から成る。「腕」と「苑」では「夕」「巳」の部品が共通であり、従って、この二つの字は、似ているということが言える。共通の部品の割合が多いほど、二つの字の形の「類字度」は高いということが言える。こういったタイプの形の似ている字の集合は、オンライン手書き文字認識のための、認識用の辞書を利用して作成することができる。また、部品が共通でなくても、例えば「互」と「瓦」のように、形が似ている文字もまたある。これは、光学文字認識用の辞書を利用して作成することができる。

【0019】このように、形の似ている字と、それがどれだけ似ているかの指標（類字度）を要素としてもつ「類字辞書」を用意し、この辞書を利用して読みのわからない文字を検索する。そのためには、まず、入力したいが読みのわからない文字に似た字を、通常のかな漢字変換方式などを用いて入力する。そして、その字をキーにして、類字辞書を引く。同様に、別の似ている字をキーにして類字辞書を引き、目的の字を絞り込んでゆく。このようにして、（複数の）キー文字から候補を得、その中から目視で目的の字を見つけ出す。従来の総画数表や部首別表を目視して探し出すことに比べ、類字辞書を引いて検索すると、少ない候補文字の中から検索すれば良いため、検索時間は圧倒的に有利になる。

【0020】

【実施例】実施例について、図面を参照して以下に説明する。図1は、本発明による日本語入力装置の一実施例を説明するための構成図で、図中、1は漢字入力手段、2は類字検索部、3は類字辞書、4はバッファ管理部、5はバッファ部、6は候補表示・選択部、7はかな漢字変換辞書等である。漢字入力手段1は、類字入力以外の

(4)

漢字入力を行うもので、かな漢字変換手段などで実現する。該漢字入力手段1によりキー文字を与えると類字検索部2により類字辞書3を検索する。類字辞書3は、日本語入力が可能な機器において、形の似ている文字群とその類字度を集めたものである。前記類字検索部2により候補文字を与えると、バッファ管理部4ですでにバッファ部5に記憶されている文字と新しい候補文字とから新しいバッファを作成し、候補表示・選択部6で候補文字を表示・選択する。なお、かな漢字変換辞書7は、前記漢字入力手段1によるキー文字を与える場合に、必要に応じて用いられるもので、かならずしも必要とするものではない。

【0021】図2は、類字辞書の構造を示す図である。類字変換を行うためには、まず類字辞書を用意する必要がある。ある字に注目し、その字に似ている字を集める。例えば、「腕」という字をキーにそれに似ている字を列挙してみると、「腕婉宛苑蛇怨」などの字が見つかる。これらはお互いに形が似ており、その一群の字の集合を「類字グループ」と呼ぶ。類字グループの集合が類字辞書である。そして、これら類字グループの要素文字に点数を付ける。この点数は、キーとなる字に対する類字度であり、点数が大きいほどより形が似ているということを意味する。類字辞書は、すべての文字をキーにしてそれに似た字を羅列し、各自に類字度を付加したものの集合である。

【0022】図2において、まず、JISコードの順に、キー文字に対する類字グループへのポイントが置かれる。辞書の先頭は、JISコードの1番目である、2121(16進数)の文字(全角スペース)に対する、類字グループへのポイントが置かれることになる。そのポイントの先には、全角スペースの類字文字とその類字度が配列される。以下、JISコード順に、類字グループへのポイントが羅列される。

【0023】例えば、「腕」の類字は、「腕」のJISコードである4F53のポイント格納場所に格納されているポイントが示す場所から、その類字度と共に配列されている。「腕」の類字として「腕」「婉」「宛」「苑」「蛇」「怨」の6字があり、類字度がそれぞれ、90、80、70、60、50、40と仮定する。そうすると、「腕」の類字グループへのポイントが示す場所から、「腕、90、婉、80、宛、70、苑、60、蛇、50、怨、40、0」と配列されることになる(最後の0は、エンドコード)。このように、すべてのJIS文字に関して、その類字を配列した辞書を類字辞書として使用する。

【0024】次に、類字辞書を検索して入力した文字に似た字を出力する動作について説明する。図3は、類字変換の手順を示すフローチャートである。以下、各ステップ(S)に従って順に説明する。読みのわからない文字を入力するために、この方式では、形の似ている字で

読みのわかっている文字をまず入力する(S1)。これは、通常のかな漢字変換を用いるなどして、入力する。そして、その文字をキーにして、類字辞書を検索する。【0025】次に、キー文字のJISコードをもとに、その文字の類字グループを指すポイントを格納しているアドレスを算出する。そのアドレスからポイントを読み出し、類字グループの先頭アドレスを得る。そのアドレスからキー文字の1つ目の類字及びその類字度を読み出し、バッファに格納する。ポイントを進め次の類字及びその類字度を読み出す。エンドコードを読み出すまでこの動作を続ける。これによって、バッファには、キー文字の全ての類字と類字度が格納されることになる(S2)。

【0026】同様に、別のキー文字がある場合(S4)、同じ検索を行い、類字及び類字度を読み出す。そして、そのキー文字から検索された類字の中に、既にバッファにある文字と同じものがあれば、その類字度を、バッファにある文字の類字度に加え、それをその文字の新たな類字度とする。つまり、別々のキー文字から同じ文字が抽出された場合、類字度を足して、候補表示の優先順位をあげるのである。新たに読み出された類字が、バッファの中に入らない場合は、バッファにその文字を加える。バッファが溢れる場合は、類字度の小さいものから削除する(S5)。

【0027】このようにして、キー文字を複数指定し、探している可能性の高い文字を優先して前の方に表示することにより、目的の文字を探し易くすることができる。類字度順に配列された候補文字の中から目的の文字を選択し、これをもって類字変換が完了する(S6)。

【0028】図4は、図3に示す類字変換の手順に従って目的の文字を選択する動作例を示すフローチャートである。以下、各ステップ(S)に従って順に説明する。ここでは、「怨」という字を入力したいのだが、その読みがわからない場合を想定して考える。まず、この字に似ている字を、キー文字として入力する必要がある。ここでは、「宛」という字をまず、キー文字として入力している。キー文字の入力は、通常のかな漢字変換などを用いて行う(S11)。そして、このキー文字を類字変換する。候補文字として、「腕婉宛苑蛇怨…」が見つかり、類字度順に候補文字として表示される(S12)。

【0029】次に、キー文字「怒」を入力する(S13)。このキー文字を類字変換すると、「怨努恋忘…」などが候補として上がるが、すでに「宛」を類字変換した候補が上がっているため、それと統合した新しい候補文字が表示されることになる(S14)。ここでは、「宛」と「怒」を類字変換し、目的の「怨」が、候補の先頭に来ていることがわかる(S15)。

【0030】次に、類字辞書の自動作成方法について説明する。この類字辞書の自動作成方法には、(1)光学文字認識用辞書を用いる場合と、(2)オンライン手書

(5)

き文字認識用辞書を用いる場合とがある。ある文字と、別の文字がどれだけ似ているかを定量的に表現するのは難しい。また、個々の文字について、それに形の似ている文字を見つけ、類字度を評価して、類字辞書を作成するのは、人手を使って行うため、大変な労力・時間を必要とする。ここでは、光学文字認識やオンライン手書き認識の手法を用い、類字辞書の自動生成を行う手法を説明する。

【0031】まず、光学文字認識用辞書を用いて類字辞書を自動生成する方法について説明する。これは通常、認識すべき各文字のドット情報をベクトルの形で持ち、これと、OCR (Optical Character Reader : スキャナなど) で入力されたドットパターンとのマッチングを行うものである。図5 (a) , (b) にその概念を示す。ここでは、8ドット×8ドットの文字を認識することを考える。

【0032】まず、8ドット×8ドットの格子に、文字

一致度の計算は、

辞書データの特徴ベクトルを $\vec{a} = (0, 0, 0, 0, 1, 0, 0, 0, \dots, 0)$

入力データの特徴ベクトルを $\vec{b} = (0, 0, 0, 0, 1, 0, 0, 0, \dots, 0)$

とすると、

$$\frac{(\vec{a} \cdot \vec{b})}{|\vec{a}| |\vec{b}|} = \text{一致度}$$

で求められる。

【0035】従って、この光学文字認識の手法を用い、辞書データとOCRからの入力データのマッチングを行うのではなく、辞書にある各文字間でマッチングを行う。マッチングの手法は全く同じで、各文字の64次元ベクトル間の距離を求め、それを類字度とする。各文字について、一定の類字度以上を持つ文字とその類字度を、目的文字の類字グループとして辞書に登録する。

【0036】以下、光学文字認識用辞書を用いて類字辞書を自動生成する方法について、さらに詳細に説明する。

類字度の評価法

光学文字認識の原理について説明すると、システムが内部に持っている文字のマトリックスパターンと入力されたマトリックスパターンとが、どれだけ似ているかということ、そのマトリックスから得られるベクトルデータの距離でもって評価するというものである。話しを簡単にするために、3ドット×3ドットのデータで説明する。

【0037】図6 (a) を辞書にあるマトリックスパターン、図6 (b) を入力パターンとする。図6 (a)、図6 (b) の3×3ドットのマトリックスパターンは、ドットがある部分を1、ない部分を0とすると、3×3=9個の0、1の羅列、すなわち、9次元のベクトルデータで表わすことができる。これら、2つの9次元ベクトルデータの「距離」は、以下の表2で表わされる。

パターンが記録されている。それを、ドットがONになっている部分を1、OFFになっている部分を0とし、左上から順にその値を並べ、64次元のベクトルとして表現する。このベクトルが、認識対象としている各文字について存在している。この各文字のベクトルの集合が、文字認識用辞書である。

【0033】ここで、OCRから文字の入力がある。1文字8ドット×8ドットになるよう正規化されたデータを用い、辞書データとのマッチングを行う。辞書作成時と同様に、入力された8ドット×8ドットのデータを64次元のベクトルにし、それと辞書に登録されている各文字のベクトルとの距離を計算する。この2つのベクトルの距離が一致度である。辞書にある全ての文字との一致度を計算し、一番、一致度が高いものを、認識結果として出力するのが、光学文字認識の基本手法である。

【0034】

【数1】

【0038】

【数2】

$$\frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} = \text{類字度}$$

【0039】従って、辞書にあるパターンと、入力パターンが全く同じであれば、類字度は1となり、2つのパターンがずれていればいる程、0に近くなる。この値を類字度とする。

【0040】マトリックスデータより類字辞書を作成する方法

図7は、マトリックスデータより類字辞書を作成するためのブロック図で、図中、11は光学文字認識用辞書、12は類字度評価部、13は光学文字認識版類字辞書である。光学文字認識用辞書11には、入力されたマトリックスパターンと比較するための、各文字のマトリックスパターンデータがストアされており、通常は、図8に示すように、文字コード順に並んでいる。図8に示すマトリックスパターンデータから、類字辞書を自動作成するために、類字度評価部において、前述した類字度の計算法により、まず、図8の辞書の最初の文字 (J I Sコード2121の文字) と、辞書の残りの文字全部との間の類字度をそれぞれ計算する。そして、ある一定の類字度 (例えば0.7) 以上の文字について、その文字を J I Sコード2121の文字の類字として類字辞書に登録する。

(6)

【0041】次に、JISコード2122の文字と、それ以外の、辞書に含まれる全部の文字との類字度を計算し、同様にして、一定の類字度以上の文字を、JISコード2122の文字の類字として、類字辞書に登録する。同様にして、すべての文字について、残りの文字との類字度の値が一定以上の文字を、その文字の類字として類字辞書を作成するのである。このようにして、光学文字認識版類字辞書13ができる。

【0042】以上が、光学文字認識用辞書を用いて類字辞書を自動作成する方法についての説明であるが、次に、オンライン手書き文字認識用辞書を用いて類字辞書を自動作成する方法について説明する。

【0043】図9は、オンライン手書き文字認識の手法を説明するための図で、図10は、オンライン手書き文字認識用の辞書構造の概念図である。オンライン手書き文字認識であるが、これは、文字を単純なストローク（基本ストローク）の組合せとして捉え、オンラインによる手書きの入力を基本ストロークごとに認識し、基本ストロークの集合がサブパターンで、そのサブパターンの集合が文字という風に、ボトムアップ的に文字を認識するものである。

【0044】従って、オンライン手書き文字認識用辞書には、文字がどのようなサブパターンから構成されているのかを示す情報があり、これを類字辞書の生成に役立てる。すなわち、サブパターンとは、認識のアルゴリズムにもよるが、漢字における、へん、つくり、あるいは、それ自身が別の漢字になるような漢字の構成部品であり、ほとんどの漢字はいくつかのサブパターンの組合せで構成されている。従って、これらサブパターンを共有する文字どうしは形が似ているということになり、共有する割合が大きいほど、類字度が高いといえる。

【0045】例えば、「腕」という字は、「月」「㇀」「夕」「巳」の4つのサブパターンで構成されている。サブパターンの一致率が高いほど、各文字間の類字度が高くなるように設定する。また、出現頻度の高いサブパターンほど、類字度へのウェイトを高くすることにより、より正確な辞書を生成することができる。そのような方法で、全ての文字に関して、他の文字との類字度を計算し、一定の類字度以上を持つ文字とその類字度を、目的文字の類字グループとして辞書に登録する。

【0046】以下、オンライン手書き文字認識用辞書を用いて類字辞書を自動生成する方法についてより詳細に説明する。

類字度の評価法

オンライン手書き文字認識用辞書は、図10に示すような構成になっている。この例では、「結」は「糸」「土」「口」の3つのサブパターンから構成され、「詰」は、「言」「土」「口」の3つのサブパターンから構成されていることを示している。具体的には、辞書は各々図11～図13に示すような構成になっている。

【0047】図11に示す基本辞書には、文字（「詰」「結」など）が、こういったサブパターンから構成されているかを示す情報が含まれている。図12に示すサブパターン辞書には、各サブパターンが、こういった基本ストロークから構成されているかを示す情報が含まれている。図13に示す基本ストローク辞書には、各基本ストロークがどのようなものかを示す情報が含まれている。

【0048】さて、類字度の評価法であるが、2つの文字の間で共通のサブパターンが多いほど、2つの文字はより似ていると言うことができる。例えば、「苑」「宛」「草」の3つの文字について考えてみる。それぞれの文字は、図14に示すようなサブパターンから構成されるとする（サブパターンの選び方は、認識システムにより違っている）。

【0049】「苑」と「宛」では、「夕」「巳」が共通であり、3つのサブパターンのうち、2つが共通と言うことになる。一方、「苑」と「草」では、

【0050】

【表2】

「㇀」

【0051】が共通であり、同じく3つのサブパターンのうち、1つが共通と言うことになる。従って、ごく単純に考えれば、「苑」と「宛」の類字度は $2/3 = 0.67$ 、「苑」と「草」の類字度は、 $1/3 = 0.33$ ということになる。

【0052】しかし、共通なサブパターンの比率のみを2つの文字の類字度とすると、類字度が同じ値のものが数多く出現してしまい、これで構成した辞書を用いると、候補が極端に多かったり、少なかったり、ということになり、本発明の本来の目的から外れてしまう。そこで、サブパターンごとに重みを設定することにより、この問題を解決することができる。いま仮に、図15に示すように重みを設定したとする。

【0053】重み付けの原則は、「多くの文字に共通に現れるものは小さく、滅多に現れないものほど大きくする」ということである。例えば、「くさかんむり」や「うかんむり」を持った文字は非常に多く存在し、従って、これらを共通に持つ文字どうしは、それほど似ているとは思わないであろう。しかし、「巳」のように、それほど多くの文字に現れないサブパターンを共通に持つ文字どうしは、より似ていると感じるはずである。

【0054】従って、サブパターンの出現頻度に合わせて重み付けを行うことにより、より似ていると感じるであろう文字どうしの類字度を高くすることができる。この重みに関しては、全ての文字のサブパターンの出現頻度を調べ、その値に比して決定することが原則となるが、利用者の主観により、似ていると感じるものは、必ずしも出現頻度に依存するとは限らない。従って、より使い易い辞書を作成するためには、重みの調整が必要と

(7)

なる。

【0055】さて、図15に示す重みにしたがうと、「苑」と「宛」、「苑」と「草」の類字度を計算すると、図16(a)~(d)に示すようになる。図16(a)に示す重みにしたがうと、共通なサブパターンは「夕」と「巳」だから、「苑」から見た「宛」の類字度は、図16(b)に示すように0.86となる。また、図16(c)に示す重みにしたがうと、共通なサブパターンは

【0056】

【表3】

「サ」

【0057】だから、「苑」から見た「草」の類字度は、図16(d)に示すように0.14となる。この式を一般化すると、文字Aから見た文字Bの類字度は、図17に示すようになる。この式を見ればわかる通り、類字度は0から1の値をとり、1に近いほど類字度が高いということになる。

【0058】オンライン手書き認識用辞書から類字辞書を作成する方法

図18は、オンライン手書き認識用辞書から類字辞書を作成するためのブロック図で、図中、21はオンライン手書き文字認識用辞書、22は類字度評価部、23はオンライン手書き文字認識版類字辞書、24は重み評価部、25は重みテーブルである。前述した図11~図13に示す辞書により、図17の類字度評価法を用いて類字辞書を自動生成する。

【0059】まず、オンライン手書き文字認識用辞書21を用いて、重み評価部24によりサブパターンの重みテーブル25を作成する。図11に示す基本辞書の全ての文字を調べ、サブパターンごとにその数をカウントする。最も多くカウントされたサブパターンの重みを1として、サブパターンの重みを正規化し、重みテーブル25を作成する。前述したように、主観による類字度の影響を考慮した重みの調整を行うことが望ましい。次に、その重みに従って、類字度評価部22により、ある文字とその残りの文字との類字度を計算する。その類字度が一定以上(例えば0.7以上)のものを、その文字の類字として、一定以上の類字度を持つものを類字辞書に登録して、類字辞書23を作成する。

【0060】以上が、オンライン手書き文字認識用辞書を用いて類字辞書を自動生成する方法についての説明である。次に、光学文字認識版類字辞書とオンライン手書き文字認識版類字辞書の統合について説明する。前述のようにして生成した光学文字認識用辞書から作成した類字辞書と、オンライン手書き文字認識用辞書から作成した類字辞書とをひとつにまとめ、図2に示すような類字変換のための類字辞書とする。

【0061】具体的には、一つのキー文字に対し、光学

文字認識版類字辞書とオンライン版類字辞書からそれぞれ類字を取り出し、それぞれに適当な重みを付けて類字度を足し合わせる。「音」の類字が図19(類字度は満点が100になるよう正規化している)に示すようになっており、光学文字認識版の重みを0.4とし、オンライン版の重みを0.6として「音」の総合類字度を計算すると、図20に示すようになる。この値を「音」の「類字度」として類字辞書に登録する。ここでも、類字度が一定以上(例えば30以上)の数値のみを辞書に登録することにより、登録数が多くなり過ぎないようにする。また、重みは仮に、0.4と0.6としたが、この値は、実際に複数の利用者が類字変換を行ってみて、最も適当と思われる数値に調整すべきである。

【0062】次に、類字辞書の検索/表示方法について説明する。類字変換の手順は、図3に示している通りである。ここでは、辞書の検索などについて追加説明する。類字辞書は、図2に示したような構造になっている。ここでは、類字度を0から100に正規化した数値を用いている。キー文字が与えられた場合、類字検索部2は類字辞書3を検索し、キー文字の類字とその類字度をバッファ部5に格納する。例えば、キー文字を「宛」とし、その類字と類字度が図21に示すようになっている場合、バッファ部5は、図22に示すようになる。

【0063】次に、別のキー文字「怒」を与えたとする。「怒」の類字と類字度は、図23に示すとおりとする。この類字と類字度をすでにバッファにあるものに追加する。もし、新しいキー文字に対する類字の中で、すでにバッファにある文字と共通なものがあれば、その類字度をバッファの中の類字度に加えると、図24に示すようになる。バッファの大きさに限りがある場合は、類字度の小さいものから削除していく。また、類字度の閾値を適応的に変化させ、一定の類字度以上の文字だけバッファに残すようにしてもよい。閾値を適応的に変化させるとは、キー文字1つ指定の場合は、類字度50以上、2つ指定の場合は70以上、3つ指定の場合は100以上ということである。このようにして、複数のキー文字を指定することにより、目的の文字を候補の前方に持ってくるのが可能となる。以上の手法により、光学文字認識用辞書やオンライン手書き文字認識辞書から、類字辞書を自動生成することが可能になる。

【0064】

【発明の効果】以上の説明から明らかなように、本発明によると、以下のような効果がある。

(1) 現在使われている、読みのわからない文字を入力する方法では、効率的に目的の文字を探し出すことは困難であるので、「類字変換」機能を用いれば、これら従来の方式を用いた検索方法より、効率的に目的の文字を探し出すことができる。検索にかかる時間の短縮と操作の簡便化が達成される。また、読みがわからなくても、部首がわからなくても、目的の文字を検索することがで

(8)

き、さらに、画数を数える必要もない。

(2) 類字辞書を、光学文字認識やオンライン手書き文字認識の手法を用いて自動生成することにより、類字辞書編成における字洩れの防止や類字度の客観的な評価、類字グループ間での類字度評価の統一が実現でき、快適な使用感を持つ類字変換のための辞書を作成することができる。また、正確な辞書を短時間かつ低労力(低コスト)で作成することができる。

(3) 漢字入力手段から見方の異なる複数のキー文字を入力し、該キー文字ごとに類字辞書から形の似ている文字を検索し、目的とする文字の候補を絞り込み、キー文字と類字度の高い順に候補文字を表示させるので、目的とする文字の発見を容易にすることができる。

【図面の簡単な説明】

【図1】本発明による日本語入力装置の類字検索の一実施例を説明するための構成図である。

【図2】本発明における類字辞書の構造を示す図である。

【図3】本発明における類字変換の手順を示すフローチャートである。

【図4】本発明における類字変換の動作例を示す図である。

【図5】本発明における光学文字認識の手法を説明するための図である。

【図6】本発明における辞書パターン及び入力パターンを示す図である。

【図7】本発明におけるマトリクスデータより類字辞書を作成するためのブロック図である。

【図8】本発明における類字度を示す図である。

【図9】本発明におけるオンライン手書き文字認識の手法を説明するための図である。

【図10】本発明におけるオンライン手書き文字認識用の辞書構造の概念図である。

【図11】本発明における基本辞書を示す図である。

【図12】本発明におけるサブパターン辞書を示す図である。

【図13】本発明における基本ストローク辞書を示す図である。

【図14】本発明における類字度の評価法を説明するための図(その1)である。

【図15】本発明における類字度の評価法を説明するための図(その2)である。

【図16】本発明における類字度の評価法を説明するための図(その3)である。

【図17】本発明における類字度の計算を示す図である。

【図18】本発明におけるオンライン手書き認識用辞書から類字辞書を作成するためのブロック図である。

【図19】本発明における光学文字認識版類字辞書とオンライン手書き文字認識版類字辞書を示す図である。

【図20】本発明におけるキー文字に対する総合類字度を計算した例を示す図である。

【図21】本発明における類字と類字度を示す図である。

【図22】本発明における類字と類字度のバッファの状態を示す図である。

【図23】本発明における類字と類字度を示す図である。

【図24】本発明における類字と類字度のバッファの状態を示す図である。

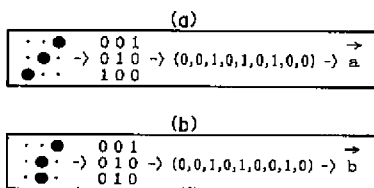
【図25】従来のコード入力用の一覧表を示す図である。

【図26】従来の部首入力用の一覧表を示す図である。

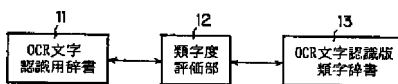
【符号の説明】

1...漢字入力手段、2...類字検索部、3...類字辞書、4...バッファ管理部、5...バッファ部、6...候補表示・選択部、7...かな漢字変換辞書等、11...光学文字認識用辞書、12...類字度評価部、13...光学文字認識版類字辞書、21...オンライン手書き文字認識用辞書、22...類字度評価部、23...オンライン手書き文字認識版類字辞書、24...重み評価部、25...重みテーブル。

【図6】



【図7】

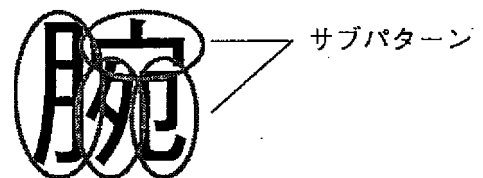


【図17】

$$\text{類字度} = \frac{\text{文字Aと文字Bで共通なサブパターンの重みの合計}}{\text{文字Aのサブパターンの重みの合計}}$$

【図9】

オンライン手書き文字認識の手法



【図14】

「兆」	→ 「㇀」	「夕」	「巳」
「宛」	→ 「㇁」	「夕」	「巳」
「章」	→ 「㇂」	「日」	「十」

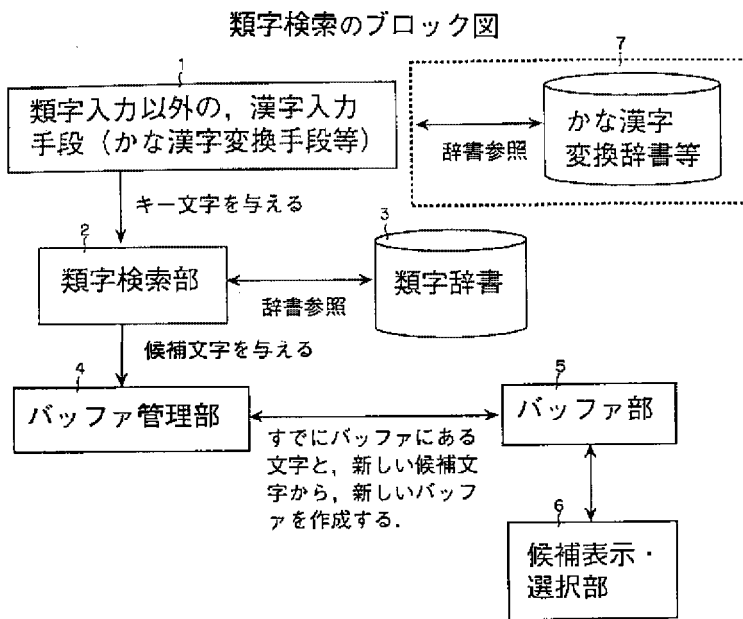
【図8】

JISコード2121の文字のマトリクスパターンデータ	苑	95
JISコード2122の文字の...	侏	88
JISコード2122の文字の...	統	88
JISコード2122の文字の...	境	70
JISコード2122の文字の...	蛇	65
JISコード2122の文字の...	怨	57

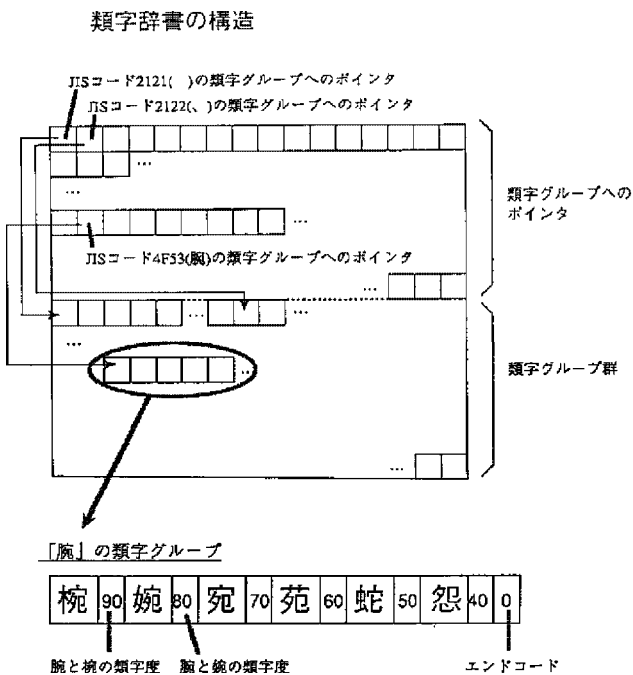
【図22】

(9)

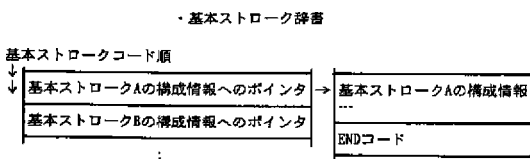
【図1】



【図2】



【図13】

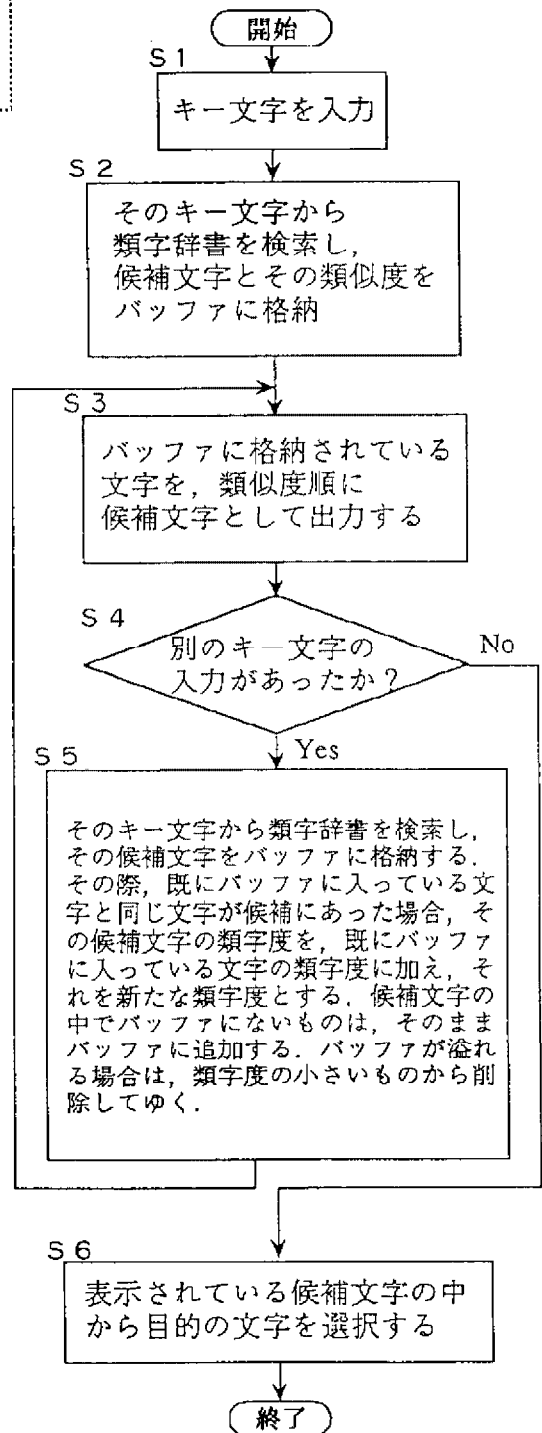


【図15】

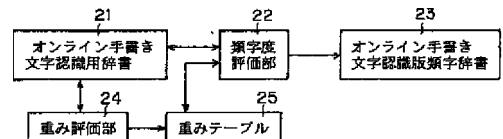
「*」	= 1
「ハ」	= 1
「日」	= 1.5
「十」	= 1.5
「夕」	= 2
「巳」	= 4

【図3】

類字変換の手順



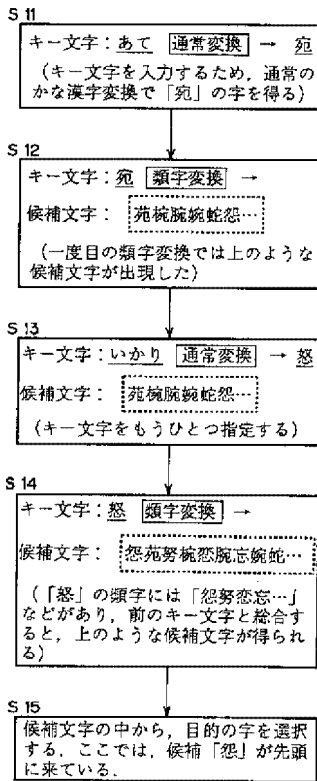
【図18】



【図4】

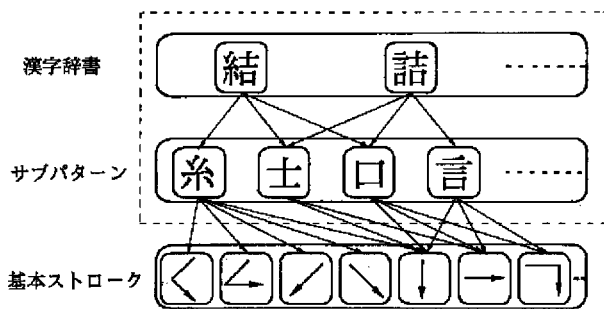
類字変換の動作例

(「怨」を入力したいのだが、読みが分からない場合の動作例)



【図10】

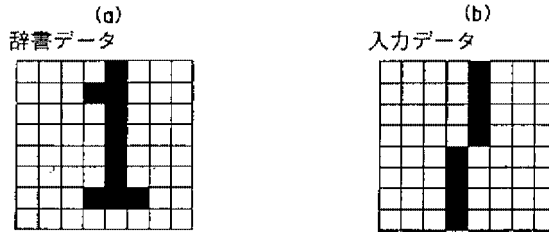
オンライン手書き文字認識用の辞書構造の概念図



【図19】

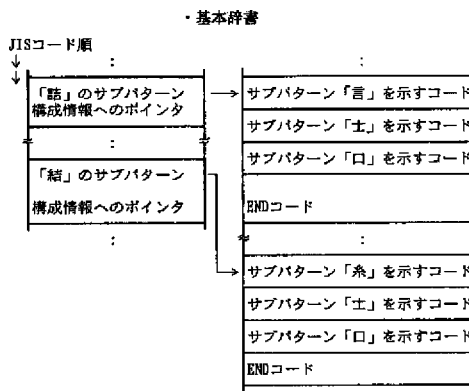
【図5】

OCR文字認識の手法



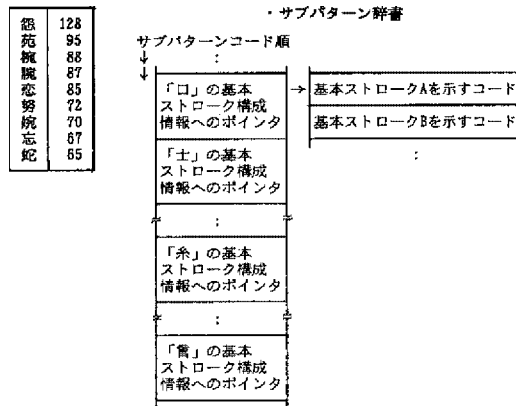
【図11】

キー文字「宛」：類字	宛	類字度	95
	宛		88
	宛		87
	宛		70
	宛		65
	怨		57



【図24】

【図12】



【図20】

【図23】

OCR版類字辞書	キー文字「音」：類字	首	類字度	72
		音		69
		音		61
		音		59
オンライン版類字辞書	キー文字「音」：	早		66
		音		66
		音		66
		日		50

	OCR版	オンライン版	
「音」類字度	$61 \times 0.4 + 66 \times 0.6$		= 64
「音」	$59 \times 0.4 + 66 \times 0.6$		= 63.2
「早」	$0 \times 0.4 + 66 \times 0.6$		= 39.6
「日」	$0 \times 0.4 + 50 \times 0.6$		= 30
「音」	$72 \times 0.4 + 0 \times 0.6$		= 28.8
「音」	$69 \times 0.4 + 0 \times 0.6$		= 27.6

キー文字「怨」：類字	怨	類字度	85
	怨		72
	怨		71
	怨		67

