

姓名の漢字仮名変換システム

中村祥次郎(日本ユニバック) 田中康仁(日本ユニバック)

・はじめに

漢字を取り扱う時に常に問題になることの1つにふりがなの問題があります。この問題はただ単に漢字と仮名のデータ処理が可能になることにより問題が解決したわけではありません。

漢字の入力と同時にふり仮名も入力しなければならぬのか? 等を悩む問題があります。ふり仮名はデータの分類や索引などで重要な意味があります。そこで漢字と仮名の関係はどのようなになっているかどのようにすれば処理できるかを調べてみました。

1. 漢字1文字に対する仮名文字の割合

毎日わかれわかれは日常生活の中に漢字を使用していますか。漢字1文字に対する仮名文字の桁数はと聞かれるとまどうものです。日本人の姓も万種名も万と千種を調べてみると次のような結果が得られました。

漢字:仮名 = 1:1.686 (姓)

漢字:仮名 = 1:1.767 (名)

この結果からみると名前の方が0.1文字だけ仮名が長いことがわかります。(この原因については別の機会に分析することにします。漢字に対する仮名の割合はこの結果からわかるように平均1.6~1.8文字です。

これにより仮名文字の桁数から漢字の桁数を割り出すと逆を行うこともできます。

2. 漢字の単語に対する仮名の割合

漢字の単語の読み方は日常文においてほぼ一定しており、その単語の読み方がある場合は必ずしもその読み方があるものを簡単に見つけることができます。たとえば青田(アノダ ヌノタ カクダ カクダ スダ スミダ カダダ)では7種類の読み方があり、村主(ムラヌシ スグリ)では2種類の読み方があります。日本人の姓も万種類、名前も万種類で漢字の読み方を調べると次のようになります。

表1 1つの漢字姓名に対する読み方

読み方件数	姓 件数	名 件数
1	42,934	30,267
2	6,748	3,007
3	1,318	543
4	347	119
5	88	30
6	25	11
7	9	8
8	3	6
9	1	1
10	1	0
11	0	1
12	0	0

この表からわかるように漢字に対する読み方は1~3種類で例外的なものでもせいぜい10種類程度です。これは仮名文字の単語に対する漢字の単語つまり同音異義語の発生に比べると発生頻度が少ないということになります。

姓や名の同音異義語の発生頻度を調べると次のようになります。

表2
同音異字件数

同音異字件数	姓	名
1	24,779	5,727
2	6,706	1,684
3	2,428	889
4	1,288	511
5	623	367
6	362	216
7	227	183
8	160	153
9	88	123
10	55	91
11	46	71
12	34	80
13	29	59
14	15	44
15	11	33
16	9	40
17	7	25
18	6	22
19	1	23
20	2	24
21	3	21
22	0	24
23	1	16
24	1	12
25	1	12
26	1	17
27	0	8
28	0	7
29	1	15
30	0	9
30以上	0	145
計	36,885	

この表からわかるように、1つの発音に対して20件~30件というものを現わします。漢字の姓名の読み方か読み方に対する音化が、ローマ字の姓名に対する漢字の発生割合は非常に多いことがわかります。

3. コンピュータによるふり仮名付けについて

コンピュータによるふり仮名を考えると全ての処理が自動的に処理されなければならぬと考えますが、何れもすべての処理を自動的に処理する必要はなく、ある部分を人間の介入により現実に近づけることが望ましくと考えています。これからその処理方法について述べてみます。

3.1 姓と名のファイルによる処理方法

日本人の姓と名前のファイル(姓10万件、名前15万件)をディスク・ファイル上に辞め漢字の姓と名前によりふり仮名を付けなければならぬファイルから取り出し、漢字ディスプレイ上に表示し、これと同時に漢字の姓漢字の名前により、姓のファイル、名前のファイルを検索し、該当する姓と名前のカナ文字を表示し、原字と照合し、キーボードまたはライトペンによる入力により正しいふり仮名付けを行います。もし該当するカナ文字の姓や名前がない場合はキーボードより入力します。

次の例によって具体的に説明します。ふり仮名を付けたいファイルはディスク上にカタログします。

← 01234567	顧客NOのキーイン
→ 東 猛	ファイルの漢字名の表示
→ 1アズマ スヒガシ 3 タケシ 4 イチノ 5 ユヨシ	姓名ファイルより該当するカナ表示
← 1,3	・1,3
→ アズマ タケシ	確認の表示
← Y =	・正しいという確認

以上のような漢字ディスプレイ上の画面処理でふり仮名付けを行います。

次にこの方式の長所、欠点について述べてみます。

- 長所 (1) 正確であり、正確である。
(2) プログラムが簡単である。

- 欠負 (1) 単語を集め続けなければならない。
 (2) ファイルの容量が大きいのでミニコン程度では処理できない。
 (3) 単語のないものに対しては人間の力に頼らなければならない。
 (4) 有限な単語から有限な単語のふり仮名付けであるためシステムの汎用性が低い。

この方法の長所、欠負を十分考えた上でシステム設計されることを望みます。

次に別の方法による漢字仮名変換の方法について述べてみます。

3.2 文字単位にふり仮名を付ける。

文字単位にふり仮名を固定して行なう方法であると早のみこみしないでいた。このような誤解を避けるために文字と姓、名について予備知識について述べておきましょう。

(1) 文字と位置

文字には位置関係があります。たとえば「仁」の名前について考えてみましょう。

康 仁 徳 元 助

「康」「徳」という文字は先頭の文字であります。「仁」は文字と文字の中間にある文字です。「元」「助」は末尾にある文字です。これらを簡単に表現するため「先頭の文字」は「頭」「中間の文字」は「中」末尾の文字は「尾」とコード付けをします。

(2) 読み方

「康」「仁」「徳」「元」「助」これらの文字には文字単位に読み方が付けられています。この読み方は辞書による読み方ばかりでなく、人名としての読み方もあります。ここではこれを含めて考えてみます。

(3) 文字の読み方と音割

読み方には音読み、訓読みなどがあります。

(4) 文字の読み方と頻度

読み方は1つの文字に3~5種類あります。しかし読み方の頻度の多いものとあまり使われない読み方もあります。また、この文字と位置の説明のように同一文字でも先頭の文字に現われる読み方と末尾に現われる文字では読み方の頻度が違ってきます。

(5) 読み方の正当性

「康仁」「徳元助」の読み方は「ヤスヒ」「トクノスケ」です。「ヤスヒ」「トクノスケ」のカナ文字を調べると次のことに気づきます。つまり「ヤ」の次に「ス」という文字は接続することもある。「ス」の次に「ヒ」と「ケ」という文字は接続している等ということ。このことにより姓名のカナ文字の連結関係を表わしたマトリックスを考えます。

後続文字

	ア	イ	ウ	エ	オ	カ	キ	ク	ケ	コ	サ	シ	ス	セ	ソ	タ	チ	ツ	テ	ト	ナ	ニ	ヌ	ネ	ノ	ハ	ヒ	フ	ヘ	ホ	マ	ミ	ム	メ	モ	ヤ	ユ	ヨ	ラ	リ	ル	レ	ロ	ワ	ヰ	ヱ	ヲ	ヾ	ン
先行文字	ア	イ	ウ	エ	オ	カ	キ	ク	ケ	コ	サ	シ	ス	セ	ソ	タ	チ	ツ	テ	ト	ナ	ニ	ヌ	ネ	ノ	ハ	ヒ	フ	ヘ	ホ	マ	ミ	ム	メ	モ	ヤ	ユ	ヨ	ラ	リ	ル	レ	ロ	ワ	ヰ	ヱ	ヲ	ヾ	ン

この表によりある姓又は名に対してのふり仮名付けをした場合の正当性の検証を行うことができます。つまり、仮名振の発音のふり仮名をした場合、この表により削除することがあります。しかしこの表だけですべての検証は行なえません。

予備知識として5つの事柄を覚えておきました。次にこれらの知識を基にして文字についての実際のデータを集めてみました。電話帳より24万5千件の名前を調査し文字別に読み方、位置別に頻度を集計しました。

表3

文字	位置	読み方	頻度	文字	位置	読み方	頻度
康	△	ヤスシ	35	徳	△	トク	14
	△	コウ	5		△	ノリ	5
	頭	ヤス	988		△	イサオ	1
	〃	コウ	111		頭	トク	1062
	〃	ミチ	2		〃	ノリ	101
	〃	ノブ	1		〃	ヨシ	4
	〃				〃	ヤス	2
	尾	ヤス	410		〃	イツ	1
	〃	コウ	7		尾	ノリ	323
					〃	トク	162
仁	△	ヒトシ	104	〃	ヨシ	1	
	△	マサシ	27	之	△	ヒサシ	1
	△	ジン	9		頭	ユキ	54
	△	シノブ	4		〃	クニ	1
	頭	ニ	146		〃	シ	1
	〃	ジン	75		〃	ノ	1
	〃	ヒト	74		中	ノ	1550
	〃	ヨシ	7		尾	ユキ	2053
	〃	キミ	3		〃	ジ	50
	〃	ジ	3		〃	ノ	22
	〃	マサ	3		〃	シ	18
	〃	マサ	3	〃	ヤ	1	
	〃	サト	1	〃	ヨシ	1	
	〃	シ	1	助	△	タズク	5
	〃	トシ	1		頭	スケ	112
	〃	ミ	1		〃	シヨ	1
	〃	ヤス	1		尾	スケ	2007
	中	ニ	33		〃	ヨシ	9
	〃				〃	スケ	1
	尾	ヒト	100				
〃	ジン	16					
〃	ジン	12					
〃	トシ	1					

この一次をまとめてみると表3のようになります。

康仁

表4

順位	読み方	ウエイト計算	カナ文字 連結性
1	ヤスヒト	988 x 100 = 98,800	○
2	ヤスジ	988 x 16 = 15,808	○
3	ヤスジン	988 x 12 = 11,856	○
4	コウヒト	111 x 100 = 11,100	○
5	コウジ	111 x 16 = 1,776	○
6	コウジン	111 x 12 = 1,332	○
7	ヤストシ	988 x 1 = 988	○
8	ミチトシ	2 x 100 = 200	○
9	ノブトシ	1 x 100 = 100	○
10	ミチジ	2 x 16 = 32	○
11	ミチジン	2 x 12 = 24	○
12	ノブシ	1 x 16 = 16	○
13	ノブジン	1 x 12 = 12	○
14	ミチトシ	2 x 1 = 2	○
15	ノブトシ	1 x 1 = 1	○

徳之助

表5

順位	読み方	ウエイト計算	カナ文字 連結性
1	トクノスケ	1062 x 2007 = 2,131,434	○
2	リノスケ	101 x 2007 = 202,707	○
3	トクノヨシ	1062 x 9 = 9,558	○
4	ヨシノスケ	4 x 2007 = 8,028	○
5	ヤスノスケ	2 x 2007 = 4,014	○
6	イツノスケ	1 x 2007 = 2,007	○
7	リノヨシ	101 x 9 = 909	○
8	ヨシノヨシ	4 x 9 = 36	○
9	ヤスノヨシ	2 x 9 = 18	○
10	イツノヨシ	1 x 9 = 9	○
	トクノズケ	1062 x 1 = 1,062	×
	リノズケ	101 x 1 = 101	×
	ヨシノズケ	4 x 1 = 4	×
	ヤスノズケ	2 x 1 = 2	×
	イツノズケ	1 x 1 = 1	×

次にこの表を基として「康仁」「徳之助」にふり仮名を付けてみます。「康仁」の場合は「康」は先頭に位置しますので「ヤスコ」ミチ「ノブ」と読むことが可能です。また「仁」は末尾に位置しますので「ヒト」ジ「ジン」トシと読みます。これらのおのおの読み方と出現頻度により読み方のウエイト計算を1順序づけを行いますと表4のようなになります。「徳之助」についても同じように読み方と出現頻度を基としてウエイト計算を行ない順序づけすると表5のようになります。

「仁」は読み方が1とだけしか読まないためウエイト計算から省きました。「徳之助」の場合「ト」「ノ」「ズ」「ケ」「リ」「ノ」「ズ」「ケ」「ヤ」「ス」「ノ」「ズ」「ケ」「ト」「ノ」「ズ」「ケ」のふり仮名が考えられます。

1か1の次にズが接続することがありません。そのためこの5つの読み方の可能性はありません。5つの読み方が発生しても妙なことに気づくことでしよう。

このように「康仁」「徳之助」の読み方は「ヤスコ」「トノズケ」が順位1位となり正しく仮名ふりが行なえました。参考までに運転免許証では「康仁」「徳之助」は「ヤスニ」「トノズキ」となっています。

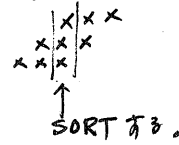
- この方式の長所は
1. ファイルの容量が小さくてよくミニコンでも十分処理できる。
 2. 単語を集めなくてよい。
 3. 有限のデータで無限の組合わせが扱え汎用性がある。
- 欠点は
1. いかなる読み方が発生する可能性がある。
 2. 特殊な読み方に対して弱い。
- があります。
- コンピュータによるふり仮名づけの2つの方式を組み合わせると互いに

良い結果が得られます。

4. コンピュータによる実験の準備。

日本人の姓名(姓10万, 名15万)のデータを解析し表3のようなものを作成しなければなりません。このためにはデータ収集のために大変な苦勞がかかるかあります。そこでこの作業を少しでも簡単に済ませる方法を考え、機械的処理と単純作業により処理することを考えました。

1. 表6は姓, 名ファイルと単純リストを作成する。
2. 表7は姓名の漢字の個数だけデータを増加し1桁づつデータを再分類した。



3. 表7を見ながら表8のデータを作る。
4. 表8のデータをファイルに収めた。この作業を姓, 名について行った。

5. コンピュータによる実験。

4. で準備した漢字から変換データに基つて表4, 表5のような計算を行ないシミュレーションを行った。

実験のデータ・テープは適当なデータがないため、女生名のテープ(女生10万, 名15万)を用いた。この姓名ファイルには姓, 名を単純に集めたもので旧来の姓, 名にウエイトを加えて(佐藤, 鈴木, 田中などは多い)でその頻度を考える)処理すべきであるがそれは行われなかった。

表6

76-04-01		姓索引表(アイツ-アイマ)				3
SEQ-NO	よみ	姓	SEQ-NO	よみ	姓	
1000193	アイツキ	相槻	1000237	アイハナ	相花	
1000194	アイツチ	合土	1000238	アイハマ	合浜	
1000195	アイツボ	合坪	1000239	アイハラ	合原	
1000196	アイツヤ	會津屋	1000241		相原	
1000197	アイツ	会津	1000244		逢原	
1000198		合津	1000245	アイバ	会場	
1000199		相津	1000246		会場	
1000200	アイツカ	相塚	1000247		合庭	

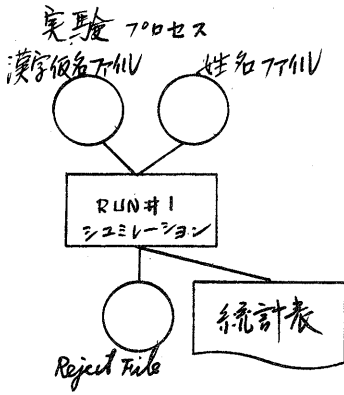
表7

NO.	文字の位置と移動	CODE	読み方	PAGE 2377	名前
85139	琢路	017746	タクジ		琢路
85140	琴路	017746	キンジ		琴路
85141	琴路	017746	キンヂ		琴路
85142	琴路	017746	コトジ		琴路
85143	琴路	017746	コトヂ		琴路
85144	環路	017746	カンジ		環路
85145	環路	017746	タマジ		環路
85146	甚路	017746	シゲジ		甚路
85147	甚路	017746	シゲヂ		甚路
85148	甚路	017746	シゲミチ		甚路
85149	甚路	017746	ジンジ		甚路
85150	甚路	017746	ヂンジ		甚路
85151	由路	017746	ユウジ		由路
85152	由路	017746	ヨシジ		由路
85153	由路	017746	ヨシミチ		由路

表8

SEQ.NO	文字	論理コード	位置	頻度	音/訓	仮名文字	旧番号
16441	鼎	22744	尾	2		カネ	16476
16442	鼓	22747	頭	1		ツツミ	16481
16443			尾	1		コ	16480
16444	齡	23006	頭	2		トシ	16484
16445			尾	3		トシ	16482
16446			尾	2		ヨ	16483
16447	龍	23024	頭	45		タツ	16488
16448			頭	29		リュウ	16487
16449			頭	1		ウタ	16491
16450			頭	1		タキ	16490
16451			頭	1		ロウ	16489
16452			中	1		リュウ	16492
16453			尾	5		タツ	16486
16454			尾	4		リュウ	16485

275 頁



姓名ファイルには漢字の姓名とふりがなが入っているため次のような処理を行った。漢字の姓名により漢字仮名ファイルの表により表4, 表5で示したふりがな付の計算を行い文字の連続性を調べコンピュータで作成したふりがなを作業順に分類し、姓名ファイルのふりがなテーブルの何個目で一致するか調べた。このテーブルと一致しなかったものはリジェクトファイルにアウトした。

表9, 表10の分析結果は次のように見る。

順位1はコンピュータ内部で発生させた作業計算の最大のものと姓名ファイルのふりがなと一致した件数である。順位2は作業計算で2位のものと姓名ファイルのふりがなと一致した件数である。

この分析結果を棒グラフにしたものが図1, 図2である。

この結果は予想以上の良い結果であった。

6. 実験結果の分析.

- 姓名の漢字仮名変換をディスプレイを使用して行う場合1つの画面に10~20ヶ対応する仮名文字を表示すれば十分であることがわかった。
- この実験の処理時間
1件の処理に大型コンピュータで200ms程度必要であった。
しかしこれは各種統計を取るために色んな計算を行っているのが大部分の処理時間と必要としているが実際のシステムでは50ms以下で十分処理できると思われる。
- 姓の実験結果より、名の実験結果がよいのは実験の途中でシステムの改良を行ったためである。
- リジェクトファイルが1割程度発生したはこの原因は次の理由による。

1. 読みがなのつかないもの
世界ヒロシ, 東海林 ショウジ
2. 変換テーブルの不備, ミス
3. 途中に1のあるもの
井上 けい子
4. <リ>返<記号>のもの
マコト 等
佐々木 ササキ

1.3. 以外のものはプログラムの改良により処理可能である。この処理を行えばリジェクトファイルはごくわずかになる予定である。

7. 次期システム.

この実験結果にそとづき図3の漢字ディスプレイを使用したふりがな付システムを日本ユニバーサル総合研究所と協力し作成する予定である。

表9

姓の分析結果

順位	件数	%	累積%
1	43,317.	49.869	49.870
2	15,240	17.545	67.415
3	6,648	7.653	75.069
4	4,020	4.628	79.697
5	2,787	3.208	82.905
6	1,935	2.227	85.133
7	1,553	1.787	86.921
8	1,138	1.310	88.231
9	971	1.117	89.349
10	771	0.887	90.236
11	715	0.823	91.059
12	480	0.552	91.612
13	491	0.565	92.177
14	396	0.455	92.633
15	381	0.438	93.072
16	281	0.323	93.395
17	294	0.338	93.734
18	234	0.269	94.003
19	258	0.291	94.272
20	186	0.214	94.486
20以上	4,764	5.484	100.000
合計	86,860.	100.000	—

リ注外件数
14,873件

漢字カナ変換結果表(姓)

図1

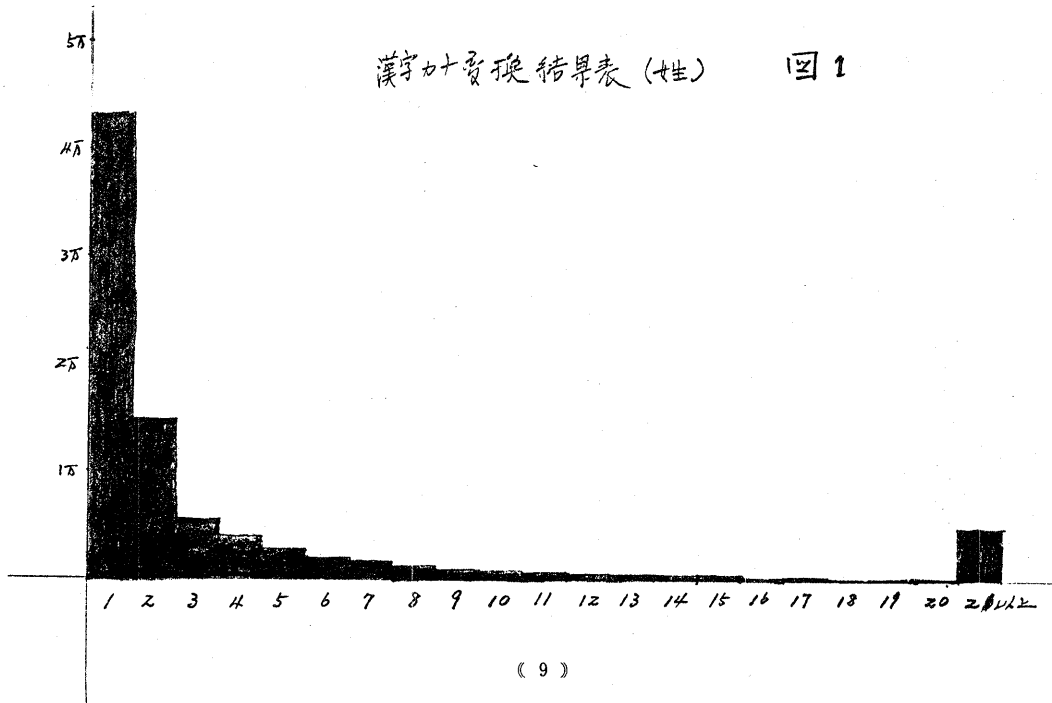
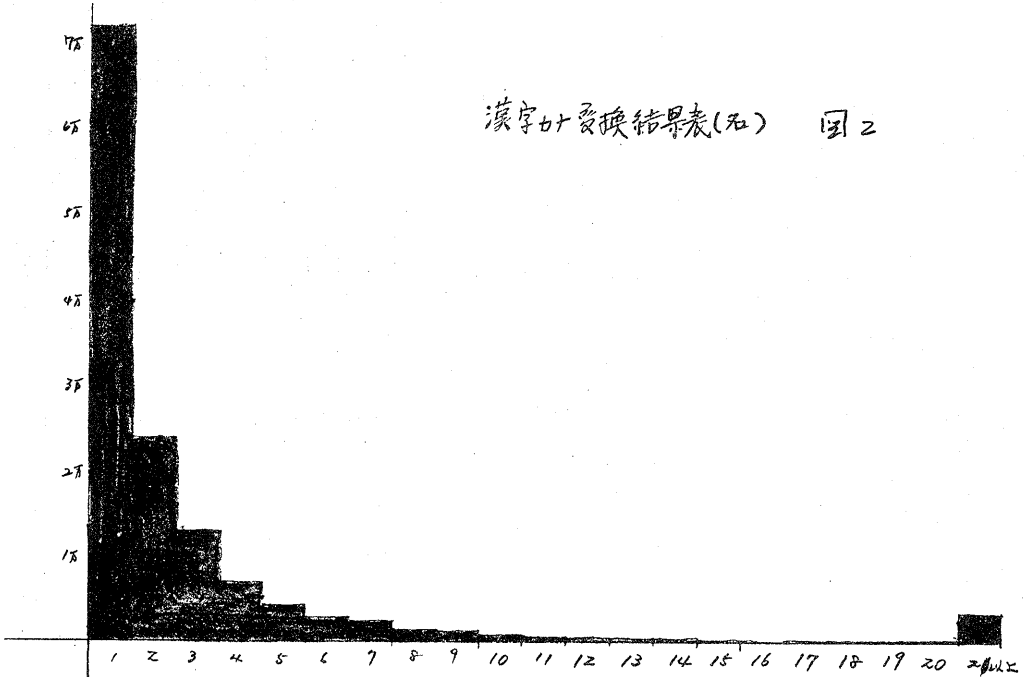


表10 名の分析結果

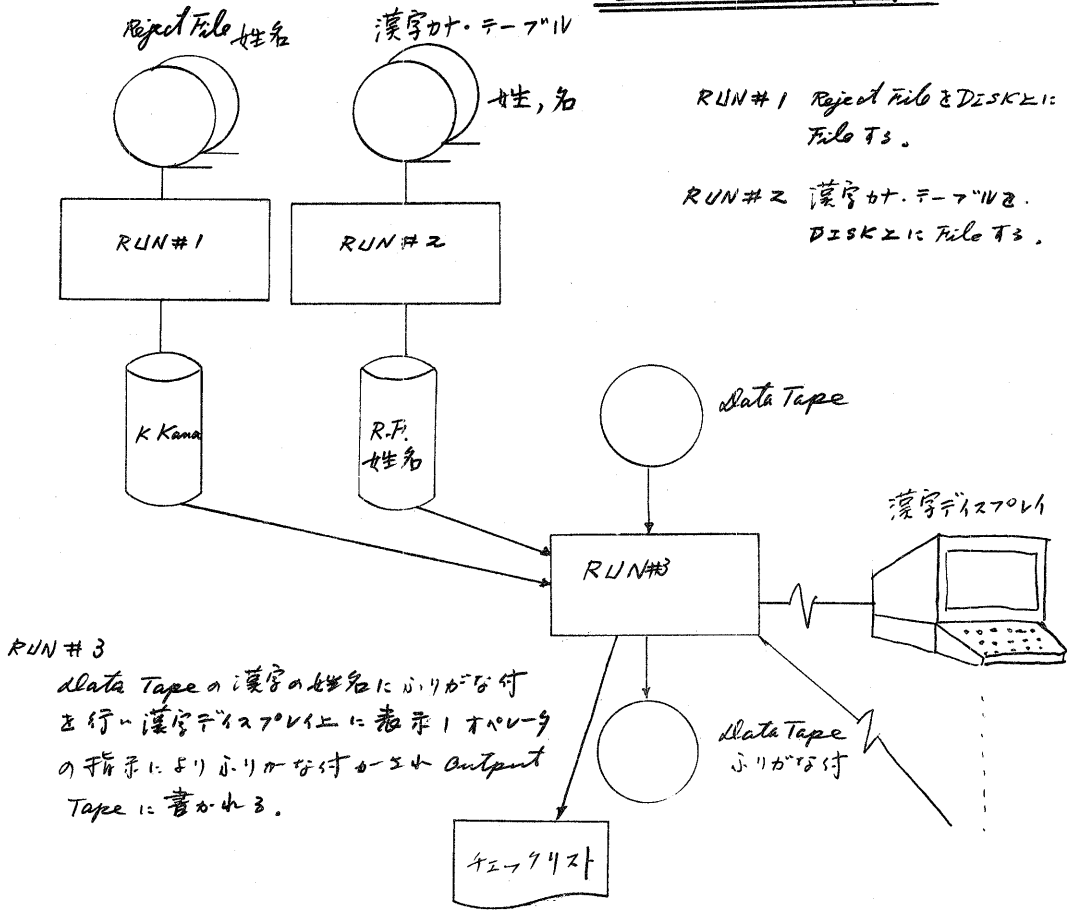
順位	件数	%	累積%
1	72,588	51.968	51.968
2	24,454	17.507	69.475
3	13,324	9.539	79.014
4	7,082	5.070	84.084
5	4,457	3.190	87.275
6	2,902	2.077	89.353
7	2,408	1.723	91.077
8	1,523	1.090	92.167
9	1,414	1.012	93.179
10	1,152	0.824	94.004
11	872	0.624	94.628
12	618	0.442	95.070
13	696	0.498	95.569
14	490	0.350	95.919
15	439	0.314	96.234
16	394	0.282	96.516
17	385	0.275	96.791
18	292	0.209	97.000
19	295	0.211	97.211
20	229	0.163	97.375
21以上	3,664	2.623	100.000
合計	150,928	100.00	—

注 外件数
11,250件



漢字の交換結果表(右) 図2

図 3. システム完成図



8. おわりに。

このシステム開発してくださった多くの方々に感謝致します。特に文字連結表の統計表を作成してくださったシステム統括第一システム推進部 安藤恵美子さん、日本ユニバーク総合研究所の方々に、データ・チェックを下さった池谷静香さんに感謝します。

中村祥次郎君にはこのシミュレーションプログラムの作成してもらいました。

9. このシステムについての問い合わせ先

日本ユニバーク株式会社 システム統括
第一システム推進部 田中康仁

(03)585-4111 (EX 3218)

東京都港区赤坂2-17-51

漢字カナ変換システム開発過程とスケジュール

0. ~ 1975年7月 思考時期
1. 1975年8月~12月 データ収集
ファイルの作成
2. 1975年1月~2月 概要説明書の作成
3. 1976年3月~4月 シミュレーションプログラムの作成
ファイルの作成(更新) データ・チェック
4. 1976年5月~6月 シミュレーション実施,
結果の分析
5. 1976年7月~10月 漢字ディスプレイ用い:システムの開発予定
6. 1976年11月~12月 評価, 次の計画立案

参考文献

1. 人名辞典 (番号案内用) 名古屋電話番号案内局
昭和46年6月1日発行
2. 計量言語学 コミカナ方式によるカナの漢字変換
田中章太(ロビ口語研究所)
- 3.