

姓名のカナ漢字変換システム

田中康仁

(日本ユニバフ(株)金融5部システム推進G)

1. まえがき

最近漢字を処理する研究が活発になってきました。漢字の *Hardware, Software* の問題は種々ありますが技術の進歩とともに解決の方向を見いだしています。

しかし漢字インプットについては新しい方式がなく相変わらず人手によって1文字1文字漢字入力装置(例えば漢字ペンコーダー、漢字ディスプレイ、ペンコーダー、漢テレ等)によって入力するのが普通です。このような入力方法ですと種々の問題があります。入力機器が多量に必要であるとか、いままでカナ文字で出来た種々のマスターの移行に長期間を要する、漢字処理システムの普及に大きな障害をおよぼす等教えあげればきりがありません。

そこでこの問題を解決する方法としてカナ文字から漢字への移行を考えてみました。特に姓名のカナ漢字変換システムをお話し申し上げ、御批判を頂きたいと思っております。

2. インプット・システムに必要な条件

漢字インプットに必要な条件は

1) 正確に

2) 早く

3) 安く

の3つの項目が上げられます。

従来使われている漢字入力機器での入力についてこの3つの条件を考えてみます。漢字入力機器での入力はオペレーターが1文字1文字、3000文字から4000文字の文字盤から望むだけなければならない。この操作を正確に行うには熟練した技術が必要である。又正確であるという保障を得るためには校正を行わなければならない。

早くという条件について考えてみると漢字入力機器の入力では平均50~60文字/分、1日の作業量は約1万5千文字であります。このくらの入力では200万件とか300万件のマスターを移行するには20台とか30台の入力機器で約1年~2年かかります。この1~2年間は約1割~2割のデータが変更されてしまいます。又新しいデータも約1割ぐらい発生します。マスターの移行が終わると入力に必要な機器は必要なくなります。数社~数十社が同時に漢字システムの導入を行うとすると単純な算術での計算だけではすまされない問題が発生します。

安くという条件を考えます。漢字の入力費用は校正の度合い、原票の状態、コーディング化されているか否か、などにより異なりますが1字2円~2円50銭程度といわれています。この価格を維持するためには数年間定期的に一定量の入力データが有るということが保障されていなければならないと思います。マスター移行のような場合を考えてみますと一時的に多量のデータが発生し、その後データ量が極端に少なくなります。

このように考えると漢字入力機器での入力は決して安い方法ではありません。

3. このシステムの開発過程.

3.1 アイデアの段階.

昭和47年3月頃新しい入力によって何か変わった方法はなにか考えている。した、その当時東邦生命の高圧部長から「学習研究社でカタ文字漢字互換が可能である」と話して、「？」ことを聞き知した。これが可能であれば大変なことなのでこっそり学習研究社の後藤氏に会いその内容を聞いてみました。その原理は簡単で「タ+カ」に対応して漢字の「田中」「田中」「多中」と対応させ、人間の介入によって「田中」を選択する方法であった。我々は全て自動的に互換することを考えていたため、このような簡単な発想が出てこなかった。この方法は人間の介入があるから従来の漢字のインポートと違うのではないかという意見もあるが、数千の文字盤の中からの2~3文字を選択のと数個の中から一つを選択するのでは大きな差があります。

そこでこの発想を理論的に発展させてみました。

3.2 カタ文字の姓名を分析する。

この大変興味深い方法について研究を行ってみた。原理はわかっても全口に氏名がいくつあるか、それがどのような頻度分布をしているか、などの基本的事項は全くわかっていない。柳田国男先生が約8~10万の姓があると書いているのはあるがこれを統計的な研究を行って書いたものではなく、一つの推測にしかすぎない。全国を網羅し重複して登録してないファイルがないのかと探してみたが適当なものが見つかった。そこでそれに近いものとして東邦生命、東邦生命のファイルと借り統計的処理を行ってみた。生命保険のファイルはあくまで全口を網羅しており、しかも氏名が独特であるから保険に入らないという性質のもので

はないので統計を取って十分意味があると判断したからです。

またカナ文字の姓と名を区別するために1桁のスペースが入っているから姓と名を区別することが容易にできるという利点もある。

その結果を表にしてみました。

表1 姓の頻度調査 (第百生命)

	種類	総頻度に対する割合	件数	総件数に対する割合
アイウエオ	1,111	4.33%	27,823	3.88%
	1,724	6.73	54,293	7.57
	823	3.21	14,957	2.09
	287	1.12	6,435	0.89
ア行計	1,535	5.88	49,428	6.90
カクケコ	1,839	7.18	42,124	5.88
	717	2.80	17,250	2.40
	827	3.23	16,530	2.30
	95	0.37	415	0.05
カ行計	1,082	4.22	30,151	4.21
サシスセソ	859	3.35	40,555	5.65
	1,442	5.63	23,829	3.33
	438	1.71	22,862	3.19
	317	1.23	5,799	0.81
サ行計	229	0.89	2,585	0.36
タチツクテト	1,387	5.41	51,913	7.25
	299	1.05	2,368	0.33
	605	2.36	9,500	1.32
	235	0.91	3,943	0.55
タ行計	923	3.60	10,732	1.49
ナニノネ	875	3.41	33,964	4.61
	435	1.70	11,254	1.57
	33	0.14	574	0.12
	78	0.30	1,273	0.17
ナ行計	317	1.23	7,012	0.97
ハヒフヘホ	934	3.64	25,695	3.58
	692	2.70	13,437	1.87
	771	3.01	22,083	3.08
	52	0.20	322	0.04
ハ行計	439	1.71	10,139	1.41
マミムメモ	752	2.93	26,423	3.69
	914	3.57	21,232	2.95
	291	1.13	8,594	1.20
	43	0.16	258	0.03
マ行計	438	1.71	13,261	1.85
ヤユヨ	672	2.62	33,931	4.74
	207	0.80	1,685	0.23
	445	1.74	10,366	1.43
	1,325	5.17	52,102	7.27
ラリルレロ	17	0.06	65	0.00
	73	0.30	500	0.06
	2	0.00	6	0.00
	10	0.03	15	0.00
ラ行計	26	0.10	80	0.01
ワ	291	1.13	13,405	1.87
	291	1.13	13,405	1.87
調査より除外したデータ			21,216	2.95
合計	25,605	100.	715,851	100.

表2 姓の頻度分布

順位	累計パーセント
~ 50	27.81
~ 100	37.21
~ 200	48.40
~ 300	55.35
~ 500	63.80
~ 1000	74.59
~ 2000	83.61
~ 3000	87.89
~ 5000	92.30
~ 10000	96.66
~ 15000	98.38
~ 20000	99.21
~ 25000	99.91

図1 姓の頻度分布 (第百生命)

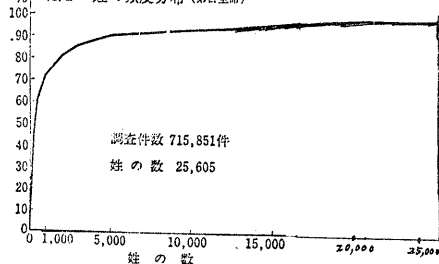


表 3 姓の頻度分布 (東邦生命)

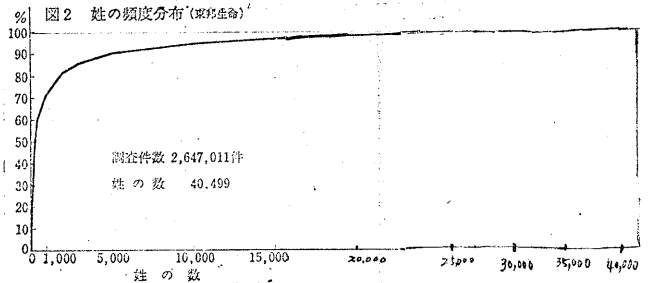
	姓の種類		件数	
	数	比率	数	比率
アイ	1,613	3.98	125,379	4.73
イウ	2,427	5.93	291,335	7.60
エオ	1,483	3.66	56,163	2.12
カ	458	1.13	25,381	0.95
ク	2,138	5.27	191,556	7.23
カ行計	8,119	19.98	589,814	22.66
カ	2,850	7.03	157,192	5.93
キ	1,107	2.73	65,642	2.48
ク	1,207	2.95	68,296	2.58
ケ	172	0.43	1,943	0.05
コ	1,657	4.09	114,420	4.32
カ行計	6,993	17.26	406,893	15.57
サ	1,257	3.10	173,697	6.56
シ	1,778	4.39	90,829	3.43
ス	709	1.75	88,033	3.32
セ	473	1.17	21,787	0.82
ソ	347	0.85	9,296	0.35
サ行計	4,566	11.27	383,635	14.49
タ	2,007	4.95	191,411	7.23
チ	430	1.06	12,095	0.45
ツ	321	0.79	36,432	1.37
テ	373	0.92	13,861	0.52
ト	1,381	3.41	43,049	1.62
タ行計	5,112	12.62	296,848	11.21
ナ	1,311	3.23	123,452	4.66
ニ	683	1.70	43,225	1.63
ノ	160	0.39	3,954	0.15
ネ	129	0.31	4,772	0.18
ヌ	500	1.23	23,459	0.88
ナ行計	2,789	6.88	198,942	7.51
ハ	1,392	3.43	84,808	3.58
ヒ	1,052	2.59	49,091	1.85
フ	1,135	2.80	83,353	3.14
ヘ	112	0.27	1,512	0.05
ホ	685	1.69	40,030	1.51
ハ行計	4,376	10.80	268,776	10.15
マ	1,105	2.72	94,750	3.58
ミ	3,384	8.38	79,657	3.00
ム	905	2.24	30,638	1.15
メ	105	0.25	1,121	0.04
モ	622	1.48	46,418	1.75
マ行計	5,711	14.10	252,584	9.54
ヤ	1,052	2.59	121,772	4.60
ユ	315	0.77	6,839	0.25
ヨ	635	1.56	61,407	2.32
ヤ行計	2,002	4.94	190,018	7.17
リ	54	0.13	228	0.03
ロ	165	0.40	1,791	0.03
レ	37	0.09	143	0.00
ル	21	0.05	38	0.00
ル	73	0.18	390	0.01
ラ行計	352	0.86	2,590	0.09
ワ	479	1.18	46,911	1.77
合計	40,499	100.	2,647,011	100.

表 4 頻度上位60までの姓 (東邦生命)

順位	苗	字	順位	苗	字
1	佐	藤	31	太	田(大田)
2	鈴	木	32	三	浦
3	高	橋	33	岡	田
4	伊	藤	34	村	上
5	渡	辺	35	藤	田
6	斎	藤	36	長	谷川
7	田	中	37	中	島
8	小	林	38	坂	井(荒井)
9	佐	々	39	工	藤
10	山	本	40	坂	本
11	中	村	41	前	田
12	加	藤	42	小	野
13	吉	田	43	青	木
14	阿	部(安部)	44	山	下
15	山	田	45	金	子
16	木	村	46	近	藤
17	松	本	47	竹	田(武田)
18	井	上	48	新	井(荒井)
19	山	口	49	中	野
20	林		50	松	田
21	菊	池(菊池)	51	中	川
22	橋	本	52	竹	内(武内)
23	森		53	田	村
24	清	水	54	柴	田(芝田)
25	遠	藤	55	千	葉
26	池	田	56	上	田(植田)
27	山	崎	57	藤	井
28	後	藤	58	西	村
29	石	川	59	福	田
30	小	川	60	岡	本

表 5 頻度上位60までの姓 (東邦生命)

順位	姓	件数	順位	姓	件数
1	サトウ	54,689	31	オノタ	7,957
2	スズキ	42,163	32	ミウラ	7,815
3	タカハシ	36,045	33	オカダ	7,738
4	イトウ	28,470	34	ムラカミ	7,731
5	ワタナベ	23,396	35	フジタ	7,598
6	サイトウ	27,603	36	ハセガワ	7,554
7	タナカ	25,825	37	ナカジマ	7,352
8	コバヤシ	22,393	38	サカイ	7,231
9	ササキ	21,197	39	タドウ	7,146
10	ヤマモト	19,709	40	サカモト	7,056
11	ナカムラ	18,873	41	マエダ	7,044
12	ホトケ	18,376	42	ノノ	6,859
13	ホンダ	18,345	43	アベ	6,840
14	アベ	16,330	44	ヤマシタ	6,794
15	ヤマダ	15,176	45	カネコ	6,744
16	カムラ	11,195	46	コドウ	6,709
17	マツモト	11,991	47	タケダ	6,661
18	イノウエ	11,762	48	アライ	6,447
19	ヤマダチ	11,585	49	サカノ	6,239
20	ハヤシ	10,911	50	マツダ	6,257
21	カケチ	10,429	51	ナカガワ	6,149
22	ハシモト	10,232	52	タケウチ	6,097
23	モリ	10,036	53	タムラ	6,006
24	シメズ	9,900	54	シバタ	5,968
25	ニンドウ	9,519	55	チバ	5,943
26	イケダ	9,484	56	ウエダ	5,911
27	ヤマザキ	8,005	57	フジイ	5,838
28	ゴトウ	8,535	58	ニムラ	5,701
29	イシカワ	8,315	59	フタダ	5,657
30	オガワ	8,159	60	オカモト	5,480



70万から264万件にデータを増加させても増加した割合で姓の数は増えず、累積曲線はほぼ一致していた。このことからカナの姓の数は約4万位又はこの増加に若干の増減はあっても大きな影響は与えないと思われる。

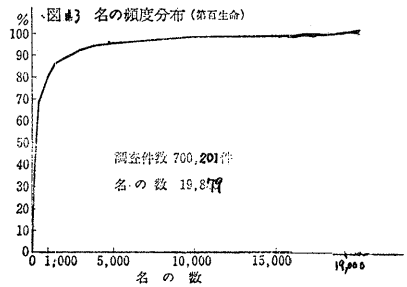
名についての統計資料

表06名の分布 (第百生命)

	男		女		男		女	
	名前の数	%	名前の数	%	件数	%	件数	%
アイ	408	2.05	125	0.62	15,515	2.27	7,205	1.02
エ	554	2.78	151	0.75	13,724	1.95	2,643	0.37
イ	194	0.97	52	0.26	590	0.84	812	0.11
ウ	233	1.27	88	0.44	6,693	0.95	5,243	0.71
エ	184	0.92	101	0.50	2,782	0.39	191	0.02
ア行計	1,593	8.01	517	2.60	40,048	5.71	16,105	2.30
カ	854	4.29	234	1.17	34,473	4.96	9,461	1.35
キ	607	3.03	236	1.18	19,166	2.73	12,824	1.83
ク	313	1.57	71	0.35	6,229	0.88	2,113	0.30
ケ	354	1.83	60	0.30	16,871	2.40	3,325	0.54
コ	471	2.36	146	0.73	14,605	2.08	976	0.13
カ行計	2,809	14.13	747	3.75	91,349	13.04	29,199	4.17
サ	524	2.66	160	0.80	12,791	1.82	7,324	1.04
シ	1,228	6.32	330	1.62	45,324	6.61	8,032	1.14
ス	307	1.54	129	0.64	5,519	0.74	4,035	0.58
セ	393	1.97	54	0.32	10,319	1.47	3,061	0.43
ソ	187	0.94	49	0.24	1,759	0.25	358	0.05
サ行計	2,719	13.67	722	3.63	77,142	11.01	22,860	3.26
タ	902	4.59	224	1.12	45,308	6.47	6,413	0.91
チ	411	2.05	185	0.93	4,036	0.57	6,539	0.94
ツ	357	1.79	132	0.66	8,586	1.22	2,470	0.35
テ	300	1.50	67	0.33	9,713	1.38	2,728	0.38
ト	853	4.29	210	1.05	32,681	4.65	11,417	1.63
タ行計	2,913	14.65	819	4.11	100,352	14.33	29,667	4.23
ナ	362	1.82	111	0.55	3,872	0.55	2,379	0.33
ニ	97	0.48	22	0.11	304	0.04	24	0.01
ノ	12	0.06	5	0.02	18	0.01	56	0.01
ネ	15	0.07	1	0.01	48	0.01	1	0.00
ナ行計	228	1.44	64	0.32	14,638	2.00	4,083	0.58
ハ	714	3.59	203	1.02	13,231	2.61	6,548	0.93
ヒ	330	1.66	122	0.61	8,295	1.17	5,639	0.80
フ	581	2.92	153	0.76	42,340	6.04	9,033	1.29
ブ	425	2.14	142	0.71	7,599	1.08	6,939	0.99
パ	83	0.42	6	0.03	356	0.05	11	0.01
ホ	106	0.53	33	0.16	235	0.03	46	0.01
ハ行計	1,528	7.68	456	2.29	59,237	8.45	21,742	3.10
マ	550	2.76	163	0.81	41,325	5.90	9,050	1.29
ミ	510	2.71	245	1.23	21,243	3.03	17,531	2.50
ム	151	0.75	33	0.16	1,481	0.21	514	0.07
メ	37	0.18	22	0.11	287	0.04	235	0.03
モ	254	1.32	53	0.26	3,639	0.51	751	0.10
マ行計	1,547	7.78	516	2.59	67,975	9.70	23,031	4.01
ヤ	289	1.45	93	0.46	11,531	1.64	4,047	0.57
ユ	256	1.28	84	0.42	14,132	2.01	6,084	0.85
ヨ	503	2.53	134	0.67	33,372	4.76	9,915	1.42
ヤ行計	1,048	5.27	311	1.56	59,035	8.43	20,031	2.86
ラ	20	0.10	12	0.06	36	0.01	59	0.01
リ	345	1.73	83	0.41	6,448	0.92	1,895	0.27
ル	12	0.06	13	0.06	37	0.01	337	0.04
レ	56	0.28	21	0.10	525	0.07	1,695	0.24
ロ	38	0.19	4	0.02	368	0.05	30	0.01
ラ行計	471	2.63	133	0.66	7,435	1.06	4,006	0.57
ワ	79	0.39	34	0.17	779	0.11	279	0.03
ワ行計	79	0.39	34	0.17	779	0.11	279	0.03
男女別計	15,421	77.53	4,456	22.38	521,633	74.45	178,568	25.46
合計	名前の数 19,879		100.0		件数 700,201件		100.0	

表07名の頻度分布

順位	累計(%)
~ 50	27.82
~ 100	40.56
~ 200	54.37
~ 300	61.90
~ 500	70.78
~ 1,000	81.34
~ 1,500	86.43
~ 2,000	89.43
~ 3,000	92.98
~ 4,000	94.89
~ 5,000	96.06
~ 10,000	98.47
~ 15,000	99.32
~ 19,879	100.00



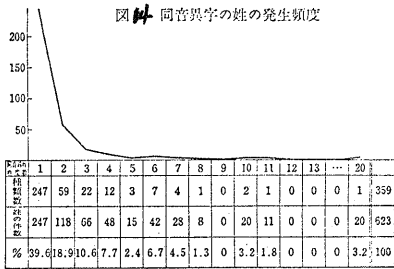
名前について姓の統計と同様なことを行ってみた。名前は世襲のそのでなく個人にとって一代かぎりのものであるから、膨大な種類があるだろうと想像していた。しかし統計を取って調べてみるとこの予想はみごとにはずれてしまった。

名前の種類は姓の種類よりむしろ少ないのである。頻度順に名前を並べ名前の累計%を作成すると表7のようになる。約3,000の名前で92%、1万で98%になる。

親達、子供に立派になつてほしいと希望をこぼるにそかわらぬ、名にも姓以上の先取りがあり、代表的名前ばかりの部分を占めていることがわかる。これに因るを因るのと比べてみると名前のグラフの元が急であることがわかる。第百生命のファイルについて同様の分析を行ってみたがほぼ同じ結果を得た。

3.3 同一発音の姓名の発生頻度

カナ文字ファイルだけの分析では十分であるため一つの発音の姓名についてそれがどの程度の種類を持っているか分析を行った。
 6万3千件の姓を調査し同音異字の発生頻度を表にすると次のようになる。



これによると一つの呼び方について平均約1.7個の異った姓があることがわかる。
 東邦生命の呼び方の件数(姓)4万件と掛け合わせる約6万5千件ということになる。

これから約全国には6万以上の姓があり又これくらいのデータは集めなければならぬ。
 名についても同様の調査を行うと平均約3.8個の異った姓があることがわかった。

東邦生命の呼び方の件数(名)3万5千件と掛け合わせる

と約13万3千ということになる。これから全国には約13万以上の名があり又これくらいのデータが必要であることがわかった。

3.4 姓名マスターデータの収集

このシステムを作成するにあたって一番困難であったのはどのようにしてデータを集めるかであった。これはコンピューターの分析力ではできない仕事でこつこつとした調査と多くの協力者によってデータを集めることができた。この方法について具体的方法、内容は省略する。

	オ一次マスター作成 総件数 異った発音		S.42.6.	オ二次マスター作成 総件数 異った発音		S.42.7.
姓	62,468件	37,212件		姓	71,048件	約4万5千件 * (※調査の名)
名	38,726件	10,655件		名	93,099件	約2万5千件 * (※調査の名)
	オ三次マスター作成 総件数 異った発音		S.42.12.予定			
姓	約12万件	約6万件				
名	約13万件	約3~4万件				

姓名マスターはオ三次マスターまで作成する予定でありこの時点で姓約12万名約13万件になる予定である。

3.5 姓名マスターの変換率の推定

姓名のカナ漢字変換用マスターを上記予定で作成しているがこれがどの程度の変換率か、どのような問題点があるか分析してみる必要がある。そこで約54万件のカナファイルと約228万件のカナファイルについてカナ文字段階での変換率の測定を行った。

姓名のマスターはオ一次マスター使用。

54万件のファイルで実験

1. 姓 98%の変換率

インポートデータ	マフチデータ	アンマフチデータ
547,210件	535,166件	12,044件

姓の種類 21,040件	マ→チ17の種類 16,237件	アンマ→チの種類 4803件
2. 名 96%の変換率		
インポイントデータ 547,093件	マ→チデータ 525,752件	アンマ→チデータ 21,341件
名の種類 13,240件	マ→チ17の種類 6781件	アンマ→チの種類 6,459件

228万件のファイルで実験

1. 姓 96.5%の変換率

インポイントデータ 2746,663件	マ→チデータ 2,166,954件	アンマ→チデータ 79,709件
姓の種類 53,390件	マ→チ17の種類 28,503件	アンマ→チの種類 24,787件

2. 名 95%の変換率

インポイントデータ 2,281,277件	マ→チデータ 2,165,988件	アンマ→チデータ 115,289件
名の種類 36,096件	マ→チ17の種類 9,184件	アンマ→チの種類 26,912件

第一次データによるテストであるが姓名とも95%以上の%を示している。これにより姓名マスターを作成する作業にはあまり先寄りがなく集められていることがわかる。

第一次データは姓6万2千件名3万8千件であるため少し不十分であることも同時にわかる。このテストはカナ文字でのテストで実際変換率はこれより2~3%下るはずである。

1か1オニ次、オニ次データではこの変換率も向上し最終的には実際変換率で次の目標になる予定である。

対象顧客数 500万件 実際変換率 姓 99%
名 98%

3.6 データファイルの問題点

アンマ→チデータを分析してみるとジ、チ、ズ、ツの使われ方がまちまちであるためアンマ→チになった例が多くみつけられた。例えば「カズオ」と「カツオ」があらわれていた。そのため「カズオ」は処理されても「カツオ」はアンマ→チになる等の例がみられた。このようなものの例としては次のようなものがある。

1. ジとチ、ズとツ
2. 清音と濁音の相異 サワ, サワ等
3. 長音, 促音などでの表し方の異い。

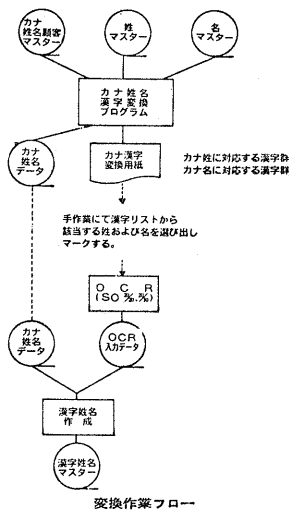
これらのデータのためには カズオはカズオと等しいという処理をほどこすことにより変換率を向上させることにした。

4. 姓名マスターを使用した変換プロセス

カナの姓名から漢字の姓名に変換する際、人間の介入を最小とし誰れでも行える方法としてOCR伝字を使用することを考えた。OCR伝字以外のキオとしては漢字テキストが考えられるがオペレータ一人に対する投資額が増えるのでこれについては

別の面から考えると、ここでは研究の対象からはずした。

4.1. 変換プロセス.



変換作業は三々に分割される。

- (1) 顧客のカナ姓名ファイルと漢字姓および名マスターから「カナ姓名漢字変換プログラム」を使用し、カナ姓に対応する漢字群およびカナ名に対応する漢字群を漢字プリンタを使用して印字する。
- (2) リストされた姓および名の漢字群から手作業にて該当する姓および名を選び出し用紙上にマークをつける。
- (3) 用紙上のマークをOCRで読み取ることによりカナ姓名に対応する漢字を決定する。

なおこのシステムで変換不可能な特殊な姓名に関しては漢字入力装置を使用して別途入力する。

4.2. 変換に使用するOCR用紙.

この用紙の記入方法はゴトウに対応するものに、ヒデオに対応するものにマーク、又は手書きの数字を記入することによって入力する。印刷されたもの以外の姓名については用紙の姓名それぞれの欄に手書きの欄がありそこに記入し漢字入力装置で入力する。

5. おわりに

カナ文字の姓名を漢字の姓名に早く、正確に変換できることがわかった。今後はこの分野でカナ漢字変換の分野で、もっと別の分野のカナ漢字変換を研究してゆきたい。最後にこのシステムを作成するにあたって多大な協力をしてくださった東洋生命、第一生命の方々に深く感謝します。