

企業名のカナ漢字変換システム

田 中 康 仁 (日本ユニバック)
大久保 静 子 (日本ユニバック)

はじめに

漢字システムに於いて重要な問題点の1つは入力作業である。入力作業はカナ文字ファイルから漢字ファイルへの移行作業と毎日毎日発生するトランザクション・データの入力作業である。これらの作業をさらに分析すると漢字の入力作業は次の3つの条件を満たしていなければならない。

この条件は

1. 安 く
 2. 早 く
 3. 正 確 に
- である。

この三つの条件はどんなシステムにも必須であるが、これ以外に次の条件を特につけ加えておきたい。

4. 取 扱 い 易 い

これは1, 2, 3という条件を満たした入力システムでも前近代的な労働条件で人を働かせていたのでは決して良いものではない。又教育とか練習が特に必要なシステムを考えるとこのために要する費用も膨大なものになる。教育とか練習が必要なシステムでは人間の採用に限られるとか特別な手当を支払わなければならない等の問題も発生する。そこでここでは誰れでも取り扱えるシステムについて考えてみることにする。特に企業名についてカナ文字での入力を考えてみる。

この研究は3つの大きな部分から構成されている。

- 第1. 企業名のカナ漢字変換の可能性について
企業名单語の収集
- 第2. 企業名(カナ)の分かち書きについて
企業名单語が十分であるか
- 第3. 実用化への試行
評価(測定)改良

この研究は第3の評価の一部分を残すだけになっている。

I 企業名のカナ漢字変換の可能性について

I.1. 企業名について

全国には一体どのくらいの企業があるのだろうか又それはどのような資料によって判るのだろうかと思ひ調査を行った。総理府統計局では5年毎に事業所統計を行いその内容を業種別、企業名(アイウエオ順)に並べ会社企業名鑑として出版している。最近の統計は昭和49年に実施しその内容が会社企業名鑑として出版されている。昭和50年3月に発行されたこの資料によると企業数は270,697件事業所数は323,637件である。これら企業も常にその姿を変えている。新しく起る企業、企業と企業の合併による名称の変更、倒産する企業、……時代の流れに従って企業名も少しずつ変化している。

このため個々の企業毎にカナ文字名と企業名の漢字を対応させることは変動す

る企業名を取り扱うには不便な事が判る。さらにカナ文字の表記方法には曖昧さがある。例えば「工業」は「コウギョウ」とカナ付けをするが「コウギョオ」「コオギョウ」……等と誤っている場合もある。又企業名は姓、名等に比べ桁数が長いという特徴がある。桁数が長くなればなるほどカナ文字の表現種類が増加する。このように考えると企業名のカナ文字に漢字の企業名を1対1に対応させることは不可能に近くなる。

1.2 企業名を単語に分解する

企業名を調べると殆どどの企業が「株式会社」という単語を持っている。さらに幾つかの企業名を頭に思い浮かべると面白いことがわかる。

例えば ○○銀行、田中商店、山田工業大阪支店を単語に分解すると

○○, 銀行, 田中, 商店, 山田, 工業, 大阪, 支店になる。

これら単語は他の企業名にもしばしば使われている。さらにこれらの単語は次のようにも分類される。

1. 人名, 地名
2. 名称 (東洋, 日本)
3. 業種名 (銀行, 保険 ……)
4. 形態名 (支店, 工業, 株式会社 ……)
5. その他の単語

このように分類された単語は次のように合成することも出来る。

$$\text{企業名} = \left\{ \begin{array}{l} \text{人名, 名称} \\ \text{地名} \end{array} \right\} \oplus \{ \text{業種名} \} \oplus \{ \text{形態名} \}$$

更に一般的に書けば次のようになる。

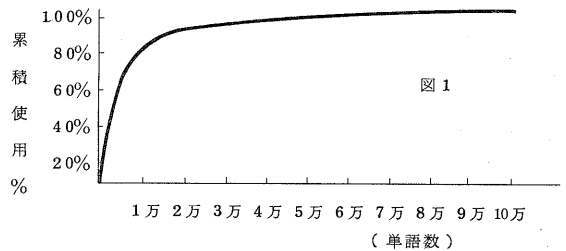
$$\text{企業名} = \{ \text{単語} \} \oplus \{ \text{単語} \} \oplus \{ \text{単語} \} \oplus \dots \oplus \{ \text{単語} \}$$

但し、この単語は企業名に使われているものである。

企業名を単語に分割し個々の単語の使用頻度を調べ使用頻度の高い順に単語を並べその累積使用率を%で表わしてみるとそのグラフは次のようになる。

1.3 企業名单語ファイルの作成

総理府統計局で出版している企業名鑑をデータ収集の基礎資料にした。この資料は全国の全業種の企業名と事務所が記載しており又資本金も一定レベル(300万円以上)に制限しているためデータ収集の対象資料として最適である。



第1回企業名单語収集状況

- | | |
|-------------------------------|---------------------|
| 1. データ収集期間 : S 49年12月~S 50年8月 | 5. 正しいデータ : 98,321件 |
| 2. 総データ件数 : 205,677件 | 6. 異った発音 : 62,115件 |
| 3. 重複データ件数 : 98,761件 | 7. 調査資料 : 会社企業名鑑 |
| 4. 桁数オーバ : 8,685件(漢字5桁以上) | 会社四季報 |

第1回の企業名单語収集で約10万件弱の単語を集めた。但し第1回の企業名单語収集作業では事業所名を省略したので第2回目の作業でこれらの入力をした。

第2回企業名单語収集作業状況をまとめてみると次のようになる。

第2回企業名单語収集状況

- 1. データ収集期間 : S50年9月~12月
 - 2. 入力件数 : 21,306件
 - 3. 第1回データ : 98,231件
 - 4. 重複データ : 16,098件
 - 5. 正しいデータ : 103,439件
 - 6. 異った発音 : 66,341件
 - 7. 調査資料 : 会社企業名鑑 (事業所名)
- 桁数オーバ (カナ10桁以上, 漢字6桁以上削除)

1.5 企業名单語ファイルの構造

収集した企業名单語ファイルの性質を調べるために次の統計資料を作成した。

- 1. アイウエオ順単語数調査 (表1)
- 2. 漢字の桁数分布 (表2)
- 3. カナ文字桁数分布 (表3)
- 4. 同音異字単語の状況 (表4)

これら統計資料によりシステム・デザインに必要な情報が得られる。特に表4は同音の漢字の書き方が異っているものはどの程度あるか調べたものである。頻度6程度で98.13%である。頻度7以上のものは使用頻度の高いものを選べばほぼ満足ゆくシステムができる。これまでの段階で研究の第1段階が終了した。

アイウエオ順単語数統計 (表1)

カナ	件数	%	カナ	件数	%	カナ	件数	%	カナ	件数	%	カナ	件数	%	カナ	件数	%
ア	3,870		サ	3,726		ナ	2,346		マ	3,567		ヤ	1,835		ル	190	
イ	4,140		シ	7,381		ニ	1,915		ミ	2,869		ユ	1,128		レ	483	
ウ	1,909		ス	3,250		ヌ	141		ム	531		ヨ	1,649		ロ	590	
エ	2,260		セ	1,980		ネ	255		メ	710		小計	4,612	4.45	小計	3,124	3.02
オ	3,631		ソ	1,142		ノ	604		モ	1,139		ラ	526		ワ	914	
小計	15,810	15.28	小計	17,479	13.73	小計	5,261	5.08	小計	8,816	8.52	リ	1,335		合計	103,439	100.0
カ	6,366		タ	5,137		ハ	4,128		同音異字単語の状況 (表4)								
キ	4,715		チ	1,717		ヒ	2,872										
ク	1,913		ツ	1,345		フ	3,409										
ケ	1,569		テ	1,919		ヘ	653										
コ	5,548		ト	4,093		ホ	2,004										
小計	20,111	19.44	小計	14,211	5.08	小計	13,066	12.63									

同音異字単語の状況 (表4)

頻度	件数	%	頻度	件数	%
1	51,456	77.56	17	22	0.03
2	7,920	11.94	18	22	0.03
3	2,922	4.40	19	9	0.01
4	1,482	2.23	20	19	0.03
5	819	1.23	21	9	0.01
6	516	0.77	22	7	0.01
7	327	0.49	23	10	0.01
8	221	0.33	24	6	0.01
9	153	0.23	25	9	0.01
10	104	0.15	26	2	0.00
11	69	0.10	27	4	0.00
12	42	0.06	28	1	0.00
13	50	0.07	29	1	0.00
14	38	0.05	30	3	0.00
15	28	0.04	30以上	38	0.05
16	32	0.05	合計	66,341	100.0

漢字の桁数と件数 (表2)

桁数	件数	%
1	1,498	1.448
2	53,234	51.464
3	27,458	26.545
4	15,935	15.405
5	5,314	5.137
合計	103,439	100.0

カナ文字桁数と件数 (表3)

桁数	件数	%
1	105	0.101
2	4,330	4.186
3	25,131	24.295
4	44,778	43.289
5	15,599	15.080
6	6,237	6.029
7	3,582	3.463
8	2,291	2.214
9	1,386	1.340
合計	103,439	100.0

企業名の桁数統計によると平均9桁で最長桁数は43桁である。3桁～29桁のデータでほぼ99.84%を占めている。桁数の長いものはヒラガナ表示のもの英字をヒラガナで表現したものが多く、このデータはファイルのデータの長さを決めるために使われる。(但しこのデータ表の中には株式会社、有限会社等の文字は含まない)

企業名单語の桁数統計は分かち書きした個々の単語の桁数である。最多桁数は4桁である。8桁以内の桁数で95.85%を占めている。9桁以上のデータは分かち書きを正しく行わなかったものがほとんどである。

企業名单語に分かち書きされる項目件数は最多項目数は2である。分かち書きは多くの場合4項目以内でほとんどが占められ、まれには8項目にまで分割されることがある。

同一企業名の発生頻度統計は同じ発音の企業名が全国ではどの程度ありその重複度を調べたものである。19万8千件の企業名に重複がないことがわかる。

企業名カナ名に対応して企業名漢字を対応させることは大変労力のかかることがわかる。(表6参照)

この中で発生頻度の最も多いものはマツモトショウテンであった。

同一企業名の発生頻度(表6)

重複度	カナ見出(件数)	累積%
781～100	13	1.0
～33	129	3.0
～22	336	5.0
～10	1276	10.0
～6	3064	15.0
～4	5955	20.0
～2	1500	25.0
～1	198824	100.0
合計	272667	-

II.3 企業名单語10万語で十分であるか?

分かち書きされたカナ文字企業名を項目ごとに分解しアイウエオ順に分類し企業名单語マスターと照合した。その結果次のことが判った。

1. 調査日時 : S51年9月
2. 分かち書き : 558,978件
項目総数
3. 企業名单語マスターとの照合結果

	件数	種類
マッチ :	477,165件	25,076件
アンマッチ :	81,813件*	41,296件*
合計	558,978件*	66,372件*

*の中には不注意による分かち書きのなされていないデータが23,000件程度入っている。これは種類、件数とも同数入っている。それ故補正して考える必要がある。

企業名のカナ文字単語は約4万件程度で十分であり頻度順分析結果は表7のとおりである。

10万語の収集した企業名单語によるマッチング率は90%(補正後)であった。アンマッチ・データのうち頻度の高いものの入力を行えば98%程度のマッチ率に容易に向上することも明らかになった。

企業名单語分析表(表7)

種類件数	単語累積%	種類件数	単語累積%
2	5%	4,688	80%
4	10%	16,932	90%
20	20%	38,424	95%
64	30%	60,783	99%
167	40%	63,577	99.5%
379	50%	65,815	99.9%
814	60%	66,372	100.0%
1,818	70%		

(補正前)

企業名单語の頻度の高いもの(表8)

1	コウギョウ	7	デンキ
2	ショウテン	8	セイサクショ
3	ケンセツ	9	ジドウシヤ
4	ショウジ	10	ニホン
5	サンギョウ	11	ハンバイ
6	ショウカイ	12	セイサクジョ

長い単語を採用すべきであろう。(有意語の選定)

さらに最長一致法に“一部後退機能”をもたせる。この機能は最長一致法が正しく行われなかった場合に長い単語を採用するものである。

分かち書きについてもほぼ正確に行えることが確認できた。ここまでの段階で第2段階の研究が終了した。

Ⅲ. このシステムの実用化

Ⅲ.1 企業名カナ漢字変換の流れとファイル構成

カナ文字の企業名より漢字の企業名への変換は既に述べた方法により次の3段階をへて行く。

1. オリジナル・データの整備(カナ文字のエラー・チェック)

タナカショウテン , ヤマダコウギョウオオサカシテン

2. これらのカナ文字を企業名カナ単語ファイルによって次のように分かち書きに変える。

(最長一致法とパターン・マッチ法を採用)

タナカ△ショウテン , ヤマダ△コウギョ△ン△オオサカ△シテン

3. 分かち書きに漢字をあてはめ漢字企業名を合成する。

田中商店

山田工業大阪支店

同音異字語の選択はOCR用紙に漢字プリンターでプリント・アウトする方法と漢字ディスプレイを使う2種類がある。OCR用紙による変換方法は次の用紙にプリント・アウトし人間の視覚によって同音異字語の選択を行いマークを付け処理する。

OCR用紙による変換方式 図3

漢字ディスプレイによる変換方式 図4

企業コード 012345678	カナ漢字変換用紙(№3 企業名) No. 23456 ◎ 太枠内は汚さないで下さい。														
企業名 ヤマダコウギョウオオサカシテン															
ヤマダ コウギョウ オオサカ シテン															
山田	興業	大阪	支店												
工業															
鉱業															
該当する漢字の企業名がない場合は下欄に鉛筆で記入して下さい。															
企業名を漢字で書いて下さい。															
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
日本ユニパック 日本ユニパック・ソフト															

漢字ディスプレイ画面

→ 企業コード : 012345678
→ カナ企業名 : ヤマダコウギョウオオサカシテン
→ 分かち書き : ヤマダ△コウギョ△ン△オオサカ△シテン
→ 分かち書きは正しく行われているか? Y, N
← Y
→ 漢字単語の選択をして下さい。
1. 山田 2. 興業 3. 工業 4. 鉱業 5. 大阪 6. 支店
← 1, 3, 5, 6
→ 漢字企業名 : 山田工業大阪支店
→ 変換は正しく行われていますか? Y, N
← Y

漢字ディスプレイ方式は日本ユニパック総合研究所で開発中である。

OCR又はシートによる変換も日本ユニパック総合研究所で開発し、ほぼ終了した。

漢字ディスプレイ方式のプロセスとファイル構成は次のようになる。

企業名カナ・インデックス(4万件), 企業名単語ファイル(10万件)をディスク・ファイルにロードし, カナ・データを読み, 分かち書きを行い, 漢字単語を引き, 漢字ディスプレイで同音語の処理をし, 漢字に変換を行う。

企業名カナ・インデックスと企業名単語は同一化すべきであるが今回の開発では別ファイルとした。

IV. このシステムの評価

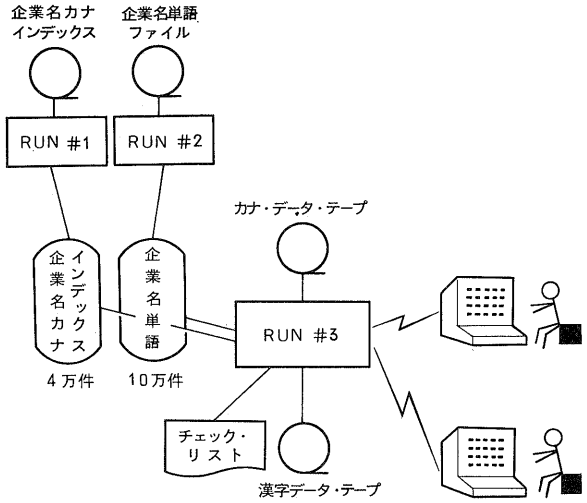
この入力方式はComputer による機械処理を多くし期間と労力を減らす方法である。今後我々の社会では単純労働の安い賃金は得られにくくなっていることを考えるとこの方法は大変良いシステムともいえる。

このような変換方式で中規模の銀行を考え直接入力方式と変換方式の比較を机上で行ってみたいことにした。

前提条件

1. 業 種 : 銀行
2. データ量 : 200 万件 (全データ)
3. 支店数 : 100 支店
4. データの性質 : 姓名データ 8 割 **
企業名データ 2 割

プロセスとファイルの構成



5. 直接入力方式には地方支店への撮影機材運搬費用, 出張旅費は含まず。
6. 変換システムでは支店でのマーク付作業は含まず。
7. 未変換率 2 割と想定している。
(実際には 1 割以下の予定)

直接入力と変換システムの比較

方式	作業工程	作業期間	コスト	データの正確性	チェックシステム	その他
直接入力方式	原票撮影	データ件数	原票撮影 円	入力ミス	原票との照合 エラーリストの訂正 支店への問合せ アンマッチ・リストの訂正	地方支店の場合 原票撮影の機材 運搬出張旅費が さらに必要。
	原票複写	200,000 件	原票複写 円	校正ミス		
	入力	100 店舗	入力 円	不明原票の扱い 等の点で精度 には限度が有 る。		
	校正	17 ヶ月	アンマッチ チェック 円	カナ文字との関 連がとりにく い。		
	訂正入力		ブルーフリ スト作成 円			
	支店への問合せ		TOTAL 円			
	アンマッチチェック		直接入力 1 件費用 を a とする。			
	訂正確認					
変換システム	支店説明会	データ件数	コンピュータ処理 円	基本的にユーザ のファイルに 忠実である。	支店チェックに よるデータは 信頼性が高い 支店別ブルー リストにより データの抜け がチェックで きる。	地方支店のマス ター漢字化も容 容易に行える。
	コンピュータ処理	2,000,000 件	漢字プリント 円	一連作業の中 に支店チェッ クが組み込ま れている。		
	漢字プリント	100 店舗	キーパンチ 円	精度はよい カナ文字との 関連がとりや すい。		
	支店チェック	6~7 ヶ月	補正入力 円			
	キーパンチ		シート 円			
	補正入力		TOTAL 0.3 a 円			
	ブルーリスト検収		変換方式コスト = 0.3 a 円			

前に述べた比較表からも判るように多量のデータについては変換方式が有利であることがわかる。特に次のような点を特徴としてあげておきたい。

1. 漢字ファイル作成期間の短縮
2. 漢字ファイル作成コストの低減
3. 漢字ファイル作成の正確性
4. 現有ファイルの有効利用（カナ文字との連動）

このシステムを作成して判ったことであるが直接入力方式ではカナ文字との連動が全くつかないことがある。

例えば次の入力を考える。

コンピュータのカナ文字： オカヤマアサヒコウトウガッコウ 又は アサヒコウトウガッコウ

原 票： 岡山県立岡山朝日高等学校

直接入力によると： 岡山県立岡山朝日高等学校

変換方式によると： 岡山朝日高等学校又は朝日高等学校

変換方式の方がコンピュータのカナ文字と連動が行われることがわかるであろう。

V. 今後の研究

今後の研究テーマとしては次のようなことを考えている。

1. 企業名单語の分析

企業名单語は姓，名，地名，名称，企業名独特のもの，その他から構成されているのでそれらの重複度合いの調査。

2. 単語の接続

単語には前に接続を許すもの，許さないもの等がある。これにより変換のあいまいさを減らす。

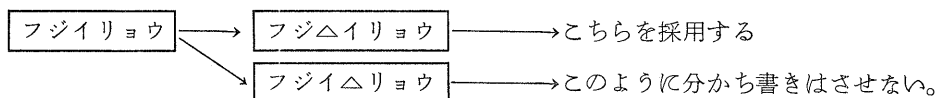
3. あいまいさの除去

ある場合には単語の分割をさけ長い単語を採用する。

ex 接続させる 化学工業

4. 分かち書きのあいまいさの除去

カナ見出しの中に特別のカナ文字列は強制的に変換するように固定する。



5. 表記方法のゆれ

表記方法の種類が多数あるもの

アルミニウム エンジニアリング

アルミニューム エンジニヤリング

6. 企業名を複合語として研究を行う。

これらについて今後調査，研究を続けてゆきたい。

VI. このシステムの研究と開発過程

この研究は2年半の期間をかけて行われた。このように長期にわたった理由はこの作業にたずさわる専門の人がアサインされず筆者の一部分の時間で行われたためである。又このような経験が他に無かったためでもある。

1. '74年 8月 ~ 11月 企業名に使われる単語の分析
2. '74年12月 ~ '75年4月 単語の整理（一次）
3. '75年 5月 ~ 7月 単語辞書の磁気テープ化（一次）
4. '75年 8月 ~ 10月 企業名のカナ漢字変換システムのドキュメント化
5. '75年 8月 ~ 10月 単語の整理（二次）
6. '75年11月 ~ '76年4月 カナ文字統計の準備
7. '76年 1月 SHARE-11 に発表（日本ユニバック総研発行）
8. '76年 5月 ~ '76年9月 カナ文字統計資料作成，シミュレーション
9. '76年10月 ~ '76年12月 ディスプレイを用いたシステム，OCR方式のシステム開発
10. '77年 1月 ~ 今後の研究

VII. お わ り に

企業名のカナ漢字変換システムについての変換率の測定，操作性，運用等については今後の研究にまかせたい。なおこのシステム作成にあたって助言して下さった社会保険庁の上田陸奥夫氏，又このシステム作成のために協力して下さった人々に深く感謝致します。

参考文献

1. 文書情報処理に関する研究48-S004（情報処理開発センター）
2. Computer Report 1974/12 P35~P45 “JAPATICにおける漢字検索システム”
3. 昭和50年度第16回大会講演論文集（情報処理学会）
昭和51年度第17回大会講演論文集（情報処理学会）

資 料

1. 会社企業名鑑（S46年度版，S50年度版）（総理府統計局編集）
2. 会社四季報（東洋経済社）
3. 電話帳 職業別