

情報処理の用字と用語

渡辺 定久
大岸 洋

(電子技術総合研究所)
(松下電器産業株式会社)

1 概要

情報処理 (Vol.17, No.1~12) の巻頭言, 論文, 資料, 解説など合計 139 件の文書 (以下「情報処理」) を対象に, 漢字と用語の使用状況を調査した。「情報処理」の漢字は, のべ字数が約 33 万字に対して約 1800 字のことでなり構成されており, 雑誌や新聞のことでなりが 3000 字を越えているのにくらべていちじるしくすくな。表外漢字の使用を制限するなどの注意を払えば 1000 字程度の字種で表記上の不都合はほとんどないと考えらる。

用語については, 漢字 2 文字以上の漢字表記語であるをけなして計算機による自動抽出を行なうと共に, 3 字長以上の用語は 2 字長用語による複合語であると考へて, 2 字長用語による分解と統合を行なった。このようにして得られた 2 字長用語は, ことでなり約 5000 語, のべ約 12 万語であり, 漢字の約 80% は 2 字長用語の表記に使用されている。3 字長以上の用語の表記に必要な漢字ののべ字数は, 全体の 2% 以下である。4 字長以上の用語の半数以上が固有名詞である。

2 漢字の数と種類

「情報処理」に出現した文字の数は, 記号, 約物を別にし, 表-1 の通りであり, 漢字についてはのべ 326,992 字が, 1,761 字のことでなりによつてまかなわれている。

この結果を, 国立国語研究所が雑誌と新聞について行なった調査と比較するに表-2 のようになつて, 「情報処理」におけることでなり字種は雑誌や新聞にくらべて, いちじるしくすくなりこがわかる。

この原因が, 「情報処理」が情報処理という限られた分野の論文などによつて占められていることにあることはいうまでもないが, これをさらに用字法の面から見ると, 新聞と雑誌では, 1,300 ~ 1,400 字もの表外漢字が使用されているのに対して, 「情報処理」における表外漢字は 200 字にすぎない。表外漢字がこのようにすくなくすんでいることが, 「情報処理」のことでなり字種をすくなくおさえている原因のひとつであるといえよう。

この表外漢字が「情報処理」の表記に不可欠なものがあるかどうかを見るため, 使用度数の高い表外漢字につい

種類	のべ字数	ことでなり字数
漢字	326,992	1,761
ひらがな	419,919	78
カタカナ	129,306	83
英数字	88,176	72
合計	964,393	1,994

表-1 「情報処理」の文字

調査対象	のべ字数	ことでなり字数		
		当用漢字	表外漢字	合計
情報処理	326,992	1,561	200	1,761
雑誌	280,094	1,835	1,493	3,328
新聞	991,375	1,844	1,369	3,213

表-2 調査対象別漢字の数

て、使用頻度とその順位および主な用例を調べた。結果は表. 3の通りである。この表より明らかのように、この表の外漢字には固有名詞の表記に不可欠な漢字や、情報処理の分野ではごく普通に使用される漢字もふくまれているが、「頃」や「殆」のように、現代風の用字法ではかながまが普通の漢字もすくなくはない。このような漢字の使用を制限したところでは、技術論文の表記に不都合があるとは思えないから、「情報処理」にとって必要な漢字は、表. 2 に示したのよりさらにすくなくしてよいように考えらる。

次に「情報処理」に出現した漢字を累使用度数の各段階ごとに合計および累計すると、表. 4が得られる。使用度数は4000から1の間への範囲にわたってこの表は漢字の使用度数調査に見られる一般的傾向である。この表によると、使用度数の99.2%上位1003字によって全体の99.2%がまかなわれていることが示されているが、これを雑誌と新聞の場合と比較すると表. 5のようになり、「情報処理」における漢字の使用度数分布の上位集中の傾向が、一般の文書の場合より強いことが一層はつきりする。

順位	漢字	度数	主な用例
308	頁	264	次頁参照
439	汎	136	汎用
593	桁	67	有効桁数
601	鞍	62	鞍部点
614	繹	58	痕繹
621	頃	57	日頃
644	閾	49	閾値
759	罫	28	罫(ベキ)
779	殆	25	殆んど
803	巾	23	10ル巾
813	狙	22	狙って
822	尤	21	尤度
848	叉	19	交叉点
"	勾	"	勾配
861	窠	18	語窠
874	箇	17	(固有名詞)
"	播	"	伝播
902	鍵	15	主鍵検索打鍵
915	阪	14	(固有名詞)
954	尖	12	尖点
976	勿	11	勿論
"	捺	"	押捺指紋
"	矩	"	矩形
"	込	"	込る
1003	且	10	且つ

順位	漢字	度数	主な用例
1003	捉	10	捕捉
1029	或	9	或る時
"	綴	"	綴(ゴリ)
"	歪	"	歪(ヒズ)み
"	楕	"	楕円
1057	廻	8	上廻り
"	誰	"	誰ども
"	燈	"	点燈
"	悉	"	悉皆走査
1090	菱	7	(固有名詞)
"	棚	"	棚積
"	坦	"	平坦
1128	料	6	(固有名詞)
"	筭	"	筭ごあり
"	栗	"	(固有名詞)
"	撻	"	撻音
"	稜	"	稜線
1177	鏡	5	(固有名詞)
"	摺	"	摺むこと
"	緻	"	精緻, 緻密
"	蝕	"	写真蝕刻
"	跡	"	跡線
1218	僅	4	僅かに
"	前	"	-前ソ
"	載	"	載く

順位	漢字	度数	主な用例
1218	辻	4	(固有名詞)
"	抹	"	抹消
"	斯	"	斯界
"	曖	"	曖昧
"	瞭	"	明瞭
"	滌	"	活滌
1282	脇	3	脇役線
"	芯	"	芯線
"	凌	"	凌駕
"	剪	"	剪断力
"	吃	"	吃水
"	宛	"	宛の
"	嶺	"	(固有名詞)
"	洩	"	洩透
"	滲	"	滲透
"	莫	"	莫大貌
"	貌	"	全貌
"	迂	"	迂回
"	亘	"	亘る
"	蹊	"	(固有名詞)
"	迥	"	迥る
"	昧	"	曖昧
"	跡	"	跡の目
"	擱	"	擱筆
"	藪	"	藪

表. 3 使用度数の99.2%表外漢字 (○:オ2水準漢字)

「情報処理」で多く使用される漢字の種類が一般文書の場合とどのように異なるかを調べる。図-1、図-2のようになる。図-1は「情報処理」に出現する漢字を雑誌の場合と比較したものであつて、この両者が共通して使用される漢字が1468字、「情報処理」のみに出現する漢字が293字であること、この293字のうち、度数順位が0-499

の範囲にあるものが4字、500-599の範囲にあるものが53字、であることなどを示している。ただし、この集計では、ホ2水準漢字は集計の対象外としたため、雑誌と新聞については固有名詞の表記に用いた場合を除く使用度数が9回(雑誌)、10回(新聞)以上の漢字のみを集計の対象とした。図-3には雑誌と新聞について上と同じ

度数区間	合計		累計		表外漢字				
	異なり	延べ	異なり(%)	延べ(%)	ホ1水準	ホ2水準			
~4001	1	4,570	1	0.1	4,570	1.4			
4000~3001	8	26,747	9	0.5	31,317	9.6			
3000~2001	20	44,296	29	1.6	75,613	23.1			
2000~1001	55	75,627	84	4.8	151,240	46.3			
1000~501	114	80,090	198	11.2	231,330	70.7			
500~401	40	18,034	238	13.5	249,364	76.3			
400~301	41	14,225	279	15.8	263,589	80.6			
300~201	86	21,299	365	20.7	284,883	87.1	1		
200~101	138	20,346	503	28.6	305,234	93.3	1		
100~91	26	2,465	529	30.0	307,699	94.1			
90~81	24	2,043	553	31.4	309,742	94.7			
80~71	33	2,506	586	33.3	312,248	95.5	1		
70~61	20	1,314	606	34.4	313,562	95.9	2		
60~51	37	2,066	643	36.5	315,628	96.5	1	1	
50~41	49	2,226	692	39.3	317,854	97.2		1	
40~31	53	1,820	745	42.3	319,674	97.8			
30~21	86	2,160	831	47.2	321,834	98.4	4	1	
20~10	172	2,563	1,003	57.0	324,397	99.2	12	1	
	26	260	1,029	58.4	324,657	99.3	2		
	9	252	1,057	60.0	324,909	99.4	5		
	8	264	1,090	61.9	325,173	99.4	4		
	7	266	1,128	64.1	325,439	99.5	3		
	6	294	1,177	66.8	325,733	99.6	4	1	
	5	41	205	1,218	69.2	325,938	99.7	4	1
	4	64	256	1,282	72.8	326,194	99.8	8	1
	3	86	258	1,368	77.7	326,452	99.8	13	6
	2	147	294	1,515	86.0	326,746	99.9	22	8
	1	246	246	1,761	100.0	326,992	100.0	59	33

表-4 「情報処理」の漢字の度数分布

計 146 54
合計 200

こを行な。た結果を、比較のため示しておく。

図.1, 図.2を図.3とくらべると「情報処理」における使用字種と雑誌や新聞の使用字種との間には若干の相異があることがわかるものの、「情報処理」に出現する漢字の86.87%は雑誌や新聞にも出現する漢字であり、「情報処理」あるいは雑誌や新聞のみに出現する漢字のそのほかの資料における使用度数順位は低い場合が多い、という面もある。このようなことから考えて、情報処理分野の技術論文を記述するには必要な漢字は、他の分野でもよく使用される漢字に情報処理に特有な漢字を追加したものであるということができる。

	情報処理	雑誌	新聞
上位の10字	10.4%	8.6%	10.6%
200	70.7	52.0	56.1
500	93.3	75.4	79.4
1000	99.2	90.0	93.9
1500	99.9	96.0	98.4

表.5 使用度数分布の比較

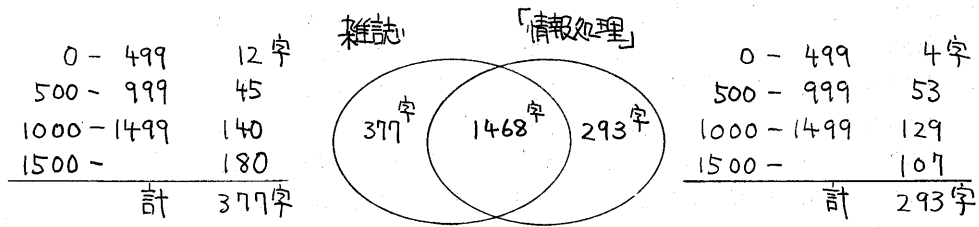


図.1 「情報処理」と雑誌

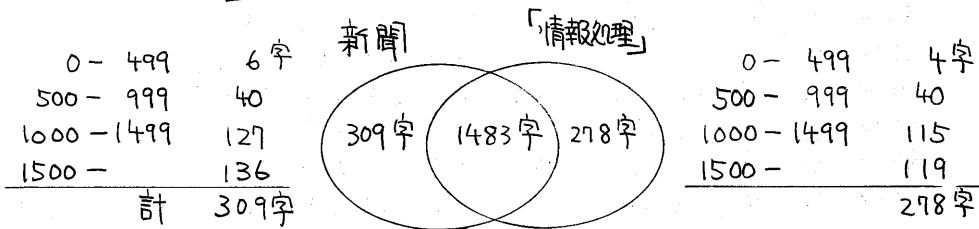


図.2 「情報処理」と新聞

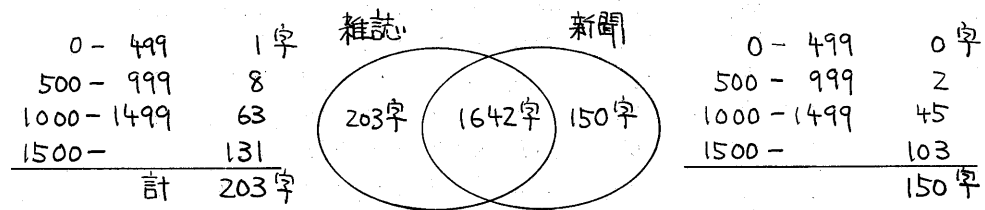


図.3 雑誌と新聞

そのような漢字が具体的にどのようなものであるかを、表.6によつて示した。この表は漢字を「情報処理」における出現順にならべ、各漢字に新聞における使用度数順位を併記したものである。

3 用語について

二ニという用語とは、「情報処理」の文章のなかで、漢字が2文字以上連続してなっている漢字列を指す。二のような漢字列のなかには、「存外難い」、「光陰矢の如し」、「全然違う」、「結構旨い」、「毎回解く」のように、用語ではないものもふくまれる一方、行(行う)、列(列)のような1字漢語や、「読み出し」、「書き込み」のように仮名と漢字の組みあわせが用語になっている例もあるが、このような例はすくなくはない。特に、漢字1字だけの用語は、「行なう」、「示す」、「用いる」、「これに付して」、「大きい」、「少ない」、「及び」など動詞、形容詞、副詞、接続詞が大部分であり、情報処理用語としてより、一般文書のなかでも多用される一般用語でもあるため、情報処理のための用語を考えるには取り上げない方が良くも考えられる。

上で述べたように用語とは凡字の漢字が並んでいる漢字表記語であるとして、「情報処理」に出現する用語の文字数別の語数は表. 7 のようになる。用語のなかには7文字以上20文字の漢字列からなるものもあるが、そのような用語はこの表から除いた。

長さ	ことば語数	のべ語数	のべ字数
1字	1,003	47,483	47,483
2	4,385	71,011	142,022
3	4,483	15,602	46,806
4	5,759	13,001	52,004
5	2,195	3,710	18,550
6	1,243	1,951	11,706
合計	19,068	152,758	318,571

表. 7 用語の長さ別語数

この表によると、ことば語数では

4字長用語が最も多くて3字長用語が二小にすぎ、この二つでことば語数全体の54%に達しているが、3字長以上の用語のなかには、以下の例のように2字長用語による複合語である例が多くふくまれていると考えられる。

- ① 3字長用語 --- 計算④機
- ② 4字長用語 --- 情報④処理
- ③ 5字長用語 --- 大型④計算④機
- ④ 6字長用語 --- 構造④記述④言語

従って、このような分割と統合を行なえば、コンパクトな(専門用語)辞書ファイルが作成できる可能性が大きい。

もっとも、「情報処理」における3字長以上のことば語は13,680語に上るような操作を人手によって行なうのは不可能であったので、われわれは分割の単位を「情報処理」に最初から存在する2字長用語とする計算機処理を試みた。この場合、3字長以上の用語が2通りの方で分割できる場合には、出現度数の多い方の2字長用語によって分割する。たとえば、「自動的」という用語は、「自動」と「動的」として2通りの2字長用語によって2通りの方法で分割可能であるが、実際には使用度数の多い「動的」を分割の単位として使用する。分割の結果、新しく生じた2字長用語があれば、4小も分割の単位として使用する。結果は表. 8 に示す通りである。

長さ	ことば語数	のべ語数	のべ字数
1字	1,086	67,258	67,258
2	5,262	123,152	246,304
3	503	1,183	3,549
4	239	289	1,156
5	47	56	280
6	4	4	24
合計	7,141	191,942	318,571

表. 8 分割・統合後の用語

また、用語を表記するのに必要な漢字の数を、2字長用語への分割と統合を行なう前後について比較すると表.9 のようになり、用語を表記するのに必要な漢字は、分割と統合の結果、ほぼ1/5に減少していることがわかる。

次に、このようにして得られた2字長用語の使用度数分布を示すと、表.10 のようになる。使用度数が1000を越える用語がある一方では、「情報処理」全体を通じて1回しか出現しなかった用語が2000近くあるなど、広い範囲でバラツキしているものの、上位の300語

長さ	分割前	分割後
1字	1,003字	1,087字
2	8,770	10,524
3	13,449	1,509
4	23,036	952
5	10,975	235
6	7,458	24
合計	64,691字	14,331字

表.9 用語の表記に必要な文字数

によるカバー率が60%を越えていることは、一般文書における漢字のことでな

度数区間	合計		累計		割合	
	こたまり	のべ	こたまり	%	のべ	%
~1001	6	7,625	6	0.1	7,625	6.2
~501	27	18,614	33	0.6	26,239	21.3
~401	17	7,493	50	1.0	33,732	27.4
~301	21	7,291	71	1.4	41,023	33.3
~201	60	14,396	131	2.5	55,419	45.0
~101	147	20,526	278	5.3	75,945	61.7
~91	18	1,737	296	5.6	77,682	63.1
~81	44	3,738	340	6.5	81,420	66.1
~71	43	3,218	383	7.3	84,638	68.7
~61	57	3,746	440	8.4	88,384	71.8
~51	75	4,137	515	9.8	92,521	75.1
~41	104	4,722	619	11.8	97,243	79.0
~31	118	4,187	737	14.0	101,430	82.4
~21	195	4,895	932	17.7	106,325	86.3
~11	443	6,568	1,375	26.1	112,893	91.7
10	78	780	1,453	27.6	113,673	92.3
9	75	675	1,528	29.0	114,348	92.9
8	94	752	1,622	30.8	115,100	93.5
7	111	777	1,733	32.9	115,877	94.1
6	158	948	1,891	35.9	116,825	94.9
5	179	895	2,070	39.3	117,720	95.6
4	267	1,068	2,337	44.4	118,788	96.5
3	377	1,131	2,714	51.6	119,919	97.4
2	685	1,370	3,399	64.6	121,289	98.5
1	1,863	1,863	5,262	100.0	123,152	100.0

表.10 2字長用語の使用度数分布

り字種とカバー率の関係に匹敵する。表.10 は使用度数の多い2字長用語の例がある。

順位	用語	順位	用語	順位	用語	順位	用語	順位	用語
0	処理	37	領域	74	分野	111	最大	148	通称
1	現場	38	条件	75	目的	112	必要	149	秘更
2	計画	39	操作	76	的索	113	有関	150	決立
3	情報	40	管理	77	検動	114	設連	151	合号
4	構造	41	研究	78	表示	115	定値	152	討入
5	必要	42	内容	79	形体	116	準心	153	後時
6	機能	43	要素	80	全体	117	中心	154	順差
7	問題	44	態上	81	指指	118	分類	155	近階
8	方法	45	基本	82	評重	119	容分	156	用查
9	文制	46	方記	83	重単	120	内線	157	率観
10	構成	47	方記	84	近手	121	順論	158	様似
11	関係	48	変作	85	手名	122	効複	159	元均
12	言利	49	作技	86	信手	123	定種	160	記種
13	利用	50	技術	87	程概	124	良化	161	上力
14	関可	51	論理	88	念線	125	識成	162	程過
15	定義	52	理操	89	号法	126	書書	163	具初
16	時部	53	在在	90	程概	127	障障	164	複規
17	方入	54	用下	91	記在	128	知生	165	連外
18	便特	55	回路	92	在容	129	成書	166	計境
19	装置	56	味定	93	容実	130	書書	167	
20	集合	57	析置	94	実実	131	障障	168	
21	果現	58	識現	95	単単	132	成書	169	
22	発力	59	現在	96	論論	133	書書	170	
23	照憶	60	現在	97	統分	134	障障	171	
24	般行	61	現在	98	分出	135	成書	172	
25	列計	62	未令	99	率貫	136	障障	173	
26	算	63	割式	100	速来	137	成書	174	
27		64	線較	101	信方	138	成書	175	
28		65	別象	102		139	成書	176	
29		66		103		140	成書	177	
30		67		104		141	成書	178	
31		68		105		142	成書	179	
32		69		106		143	成書	180	
33		70		107		144	成書	181	
34		71		108		145	成書	182	
35		72		109		146	成書	183	
36		73		110		147	成書	184	

表.11 「情報処理」で使用度数の多い2字長用語

