国立国語研究所学術情報リポジトリ 近代語文献を電子化するための文字セット

メタデータ	言語: jpn
	出版者:
	公開日: 2020-03-18
	キーワード (Ja):
	キーワード (En):
	作成者:
	メールアドレス:
	所属:
URL	https://doi.org/10.15084/00002764

近代語文献を電子化するための文字セット

高田 智和 (国立国語研究所 理論・構造研究系)1

1. はじめに

本稿では、近代語文献を電子テキスト化する際に準拠する文字セットについて検討する。一般的に、紙媒体の文書を電子テキストに写し取る際には、標準化された符号化文字集合に準拠することが多い。かつて、JIS X 0208 文字セットに拠って、『大正新脩大蔵経テキストデータベース』(http://21dzk.l.u-tokyo.ac.jp/SAT/)の前身、『青空文庫』(http://www.aozora.gr.jp/)、『太陽コーパス』など、学術的価値を有する電子テキスト群が作成されたが、その都度、JIS X 0208 文字セットでは表現できない文字(外字)への対処が問題となっていた(安永(1998)、下田・師(1999)、富田(2000)、池田・白井・高田(2002)、高田(2002)、田中(2005)、當山(2009)、須永・堤・高田(2011)など。2000年に制定された JIS X 0213 は、外字の問題を解消し、「現代日本語を符号化するために十分な文字集合」を目指して開発された規格である。JIS X 0208 が、第 1 水準漢字、第 2 水準漢字、非漢字(X0208 非漢字と呼ぶ)の約 6,800字であるのに対して、JIS X 0213は、JIS X 0208 文字セットを拡張し、第 3 水準漢字、第 4 水準漢字、非漢字(X0213 非漢字)を加え、約 11,000字の文字セットとなっている。

JIS X 0213 文字セットに準拠して作られた『現代日本語書き言葉均衡コーパス』では、のべ 99.96%の文字を符号化できることが確認されている(高田ほか(2009))。『現代日本語書き言葉均衡コーパス』での運用実験により、JIS X 0213 は「現代日本語を符号化するために十分な文字集合」であることが実証されたと言えよう。

しかし、現代から時代を遡って、歴史的な日本語文献に対して、JIS X 0213 文字セットがどの程度の有効性を持ち得るのかは、まだ試みられていない。以下に、『太陽コーパス』の JIS X 0208 外字を対象として、JIS X 0213 による再符号化結果を報告し、近代語文献の電子テキスト化における JIS X 0213 文字セットの有効性について見通しを述べる。

2. 『太陽コーパス』の文字処理

『太陽コーパス』は 2005 年に公表された近代語研究向けのコーパスである。1895(明治 28)年から 1928(昭和 3)年まで発行された総合雑誌『太陽』のうち、1895 年、1901 年、1909 年、1917 年、1925 年の各 12 冊、全部で 60 冊分の記事を電子テキスト化し、引用、文、振り仮名などのタグを付けた電子化コーパスである。

『太陽コーパス』の文字処理では、JIS X 0208 文字セット (第 1 水準漢字・第 2 水準漢

_

¹ ttakada@ninjal.ac.jp

字・非漢字、約6,800字)に準拠して、『太陽』の記事の入力処理が行われている。『太陽コーパス』の開発は、JIS X 0213 が制定される以前から着手されており、開発当時の処理水準を考慮すれば、妥当な選択であったと思われる。

雑誌『太陽』では主に旧字体が使われているが、「羽」のように新字体「羽」との字体差がわずかであるものは、新字体で表現することを許容している。字体粒度の許容範囲は、原則として JIS 規格の包摂規準にしたがっている。また、JIS 規格の包摂規準を拡張させて、独自の包摂も行っている(田中(2005))。語彙研究や文法研究での活用が主に期待されているコーパスであるから、字体の忠実な再現は行っていない。

『太陽コーパス』における JIS X 0208 文字セットの使用状況をまとめたものが表 1 である。表の各セルで、上段がのベ字数、下段が異なり字数である。『現代日本語書き言葉均衡コーパス』では、第 2 水準漢字は異なり 1,311 字の使用にとどまっていることから、近代日本語書き言葉では多種多様の漢字を用いていることが裏付けられる。

水準	1895	1901	1909	1917	1925	計
第1水準漢字	1,395,861	1,285,938	1,107,698	1,012,692	848,430	5,650,619
	[2,687]	[2653]	[2,648]	[2,623]	[2,633]	[2,721]
第2水準漢字	221,940	200,767	172,761	168,536	124,880	888,884
	[2,558]	[2,275]	[2,057]	[1,922]	[1,757]	[2,864]
非漢字	1,688,600	1,651,017	1,556,071	1,454,064	1,447,583	7,797,335
	[291]	[279]	[278]	[268]	[287]	[318]
計	3,306,401	3,137,722	2,836,530	2,635,292	2,420,893	14,336,838
	[5,536]	[5,207]	[4,983]	[4,813]	[4,677]	[5,903]

表 1: 『太陽コーパス』の JIS X 0208 による符号化の内訳

『太陽コーパス』では、JIS X 0208 文字セットで表現できない文字を、外字、踊字、合字、小書の 4 種のタグを用いて表現している。JIS X 0213 は JIS X 0208 を拡張したものであるから、外字、踊字、合字、小書の各タグで表現された文字類を、JIS X 0213 の拡張領域(第 3 水準漢字・第 4 水準漢字・X0213 非漢字、約 4,000 字)と突き合わせることで、『太陽コーパス』の JIS X 0213 文字セットによるカバー率をおおよそ求めることができるであろう。

JIS X 0213 の特長の一つに、JIS X 0208 では符号位置を区別せず包摂されていた字体を、符号位置を区別して表現できる点があげられる。例えば、「徳」の旧字体「徳」や、「鴎」の康熙字典体「鷗」などの異体漢字が該当する。厳密さを追求すれば、JIS X 0208 では包摂するとされていた異体漢字や、『太陽コーパス』で独自に包摂した異体漢字についても、雑誌『太陽』で実際に使われている字体を確認し、JIS X 0213 の拡張領域の符号位置で表現できるかどうかを検討しなくては、JIS X 0213 文字セットによる真のカバー率を求める

ことはできない。しかし、今回の調査目的は、近代語文献の電子テキスト化における JIS X 0213 文字セットの有効性について見通しを得ることであるから、『太陽コーパス』の開発において、JIS X 0208 文字セットで表現できないとされた文字だけを調査対象としても、調査目的が達せられるものと判断した。

さて、外字、踊字、合字、小書、各夕グの使用状況は表 2 のとおりである。次節以降、各夕グについて、JIS X 0213 による再符号化結果を示す。

表 2: 『太陽コーパス』の JIS X 0208 外字の内訳

タグ名	のベタグ数
外字	5,507
踊字	18,019
合字	8,554
小書	187
計	32,267

3. 『太陽コーパス』の JIS X 0213 による再符号化

3.1 外字タグ

『太陽コーパス』の外字タグの運用方法は2種類ある。一つは代用で、もう一つは本当の外字である。代用はJIS 規格の包摂規準では処理できないが、JIS X 0208 に採録された文字の異体漢字である場合に用いる。外字はJIS X 0208 の文字と異体関係が認められず、全くの別字の場合に用いる。それぞれタグの使用例を示す。なお、タグの属性「文字番号」は文字鏡番号である。

(例1)【代用】 潔

<外字 文字番号="001721">熈</外字>々たる明治二十八年の新旭光は 〔t189501〕

(例2)【外字】 礴

内には浩然たる正氣の磅<外字 文字番号="024597"> **=**</外字>するところ禁 ぜんと欲して能はざるあり。[t189501]





外字タグの再処理結果をまとめたものが表 3 である。JIS X 0208 を用いたときに外字タ

グで表現したものの 77.7%が、JIS X 0213 を用いると符号化することができる。

表3:外字タグの再処理結果

	のベ字数	異なり字数
第2水準漢字	2	2
第3水準漢字	2,685	446
第 4 水準漢字	1,371	426
X0213 非漢字	220	38
(小計)	4,278	912
X0213 外字	1,229	579
計	5,507	1,491

第 2 水準漢字で表現できる文字は、次の 2 字 (表 4 参照)である。これらはコーパス開発時の単純なバグであろう。

表4:第2水準漢字で表現できるもの

文字番号	字形	面区点	度数
037880	蹶	1-77-12	1
047275	鷙	1-83-25	1

第 3 水準漢字で表現できる文字は、次の 446 字 (表 5 参照)である。『太陽コーパス』での使用度数順に示す。

表 5:第3水準漢字で表現できるもの

文字番号	字形	面区点	度数
048824	龐	1-94-86	242
016085	欵	1-86-31	239
042124	雞	1-93-66	138
001721	淵	1-14-55	90
051106	开	1-84-17	71
056009	厲	1-14-84	64
035458	詹	1-92-08	54
050005	咜	1-14-88	37
023439	睁	1-88-85	34
032700	虚	1-91-45	29
053267	噶	1-15-20	25

040346	鉸	1-93-13	24
001244	儞	1-14-45	20
019395	燄	1-87-64	20
041315	鬨	1-93-49	19
004276	嘻	1-15-18	17
019890	牖	1-87-69	17
016408	殂	1-86-38	16
024597	礴	1-89-18	16
053201	摹	1-84-88	16
011865	抅	1-84-72	15
012081	挘	1-84-77	15
011617	戢	1-84-66	14

017165	汴	1-86-52	14
027733	縈	1-90-16	14
038785	迤	1-92-52	14
038791	迨	1-92-53	14
000076	丰	1-14-06	13
005113	埈	1-15-47	13
010174	徧	1-84-34	13
015155	楣	1-85-86	13
040272	鈹	1-93-07	13
045469	鬂	1-94-27	13
000774	倘	1-14-30	12
009808	弴	1-84-22	12
018251	潢	1-87-13	12
050268	夑	1-87-67	12
008295	嵓	1-47-85	11
010149	徜	1-84-33	11
015163	楨	1-85-88	11
000597	侔	1-14-22	10
001115	僦	1-14-40	10
009610	弇	1-84-19	10
010094	徉	1-84-32	10
027069	糕	1-89-86	10
028039	纍	1-90-24	10
037868	蹰	1-92-39	10
045430	髹	1-94-26	10
059830	笻	1-89-60	10
000471	你	1-14-13	9
002930	厓	1-14-82	9
014552	枘	1-85-54	9
015843	櫧	1-86-25	9
016752	毗	1-86-44	9
019012	烘	1-87-42	9
024232	确	1-89-06	9
025635	窻	1-89-54	9
028810	翔	1-90-35	9
035370	詎	1-92-04	9
039198	邈	1-92-58	9

003528	咡	1-14-94	8
004394	嗳	1-15-23	8
012181	捥	1-84-80	8
014005	晷	1-85-32	8
017526	涇	1-86-75	8
020916	珉	1-87-89	8
021253	璣	1-88-28	8
022383	瘙	1-88-53	8
027247	紈	1-89-90	8
034513	褰	1-91-84	8
038930	逭	1-92-56	8
003523	咖	1-14-93	7
004407	噲	1-15-25	7
010496	怳	1-84-45	7
010771	惋	1-84-51	7
010815	倘	1-84-54	7
010949	愜	1-84-56	7
011833	扯	1-84-71	7
012105	挹	1-84-78	7
021062	琦	1-88-06	7
021270	璨	1-88-31	7
026734	籙	1-89-79	7
028952	耦	1-90-38	7
031565	蒞	1-91-13	7
032805	虬	1-91-50	7
040223	鈐	1-93-05	7
059130	凞	1-87-58	7
000418	份	1-14-09	6
004631	囉	1-15-31	6
006297	娓	1-15-81	6
007559	尫	1-47-62	6
010679	悞	1-84-50	6
018164	漪	1-87-06	6
018965	炷	1-87-40	6
024580	礟	1-89-16	6
025458	窅	1-89-50	6
026077	独	1-89-66	6

027750 線 1-90-17 043909 線 1-94-07 043965 線 1-94-08 001375 穴 1-14-50 003341 吧 1-14-86 003770 唯 1-15-06 003898 明 1-15-09 007047 孽 1-47-55 008028 响 1-47-74 010803 切 1-84-53 011015 ⑫ 1-84-59 011961 牧 1-84-74 012808 撃 1-84-92 015076 収 1-85-92 017132 汗 1-86-49 017186 穴 1-86-54 017699 森 1-86-54 018010 減 1-87-01 018948 太 1-87-39 019137 炊 1-87-48 021065 环 1-88-07	6 6 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
043965 1-94-08 001375 元	6 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
001375 交	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
1-14-86	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
1-15-06	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
1-15-09	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
1-47-55 008028	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
1-47-74 0108028	5 5 5 5 5 5 5 5 5 5 5 5
010803 1-84-53 011015 世 1-84-59 011961 拖 1-84-74 012808 擎 1-84-92 015076 根 1-85-83 015272 射 1-85-92 017132 汗 1-86-49 017186 沅 1-86-54 017699 森 1-86-86 018010 演 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5 5 5 5 5 5 5 5
011015 1-84-59 011961 1-84-74 012808 1-84-92 015076 根 1-85-83 015272 材 1-85-92 017132 汗 1-86-49 017186 元 1-86-54 017699 森 1-86-86 018010 演 1-87-01 018948 炫 1-87-39 019137 燃 1-87-48	5 5 5 5 5 5 5 5
1-84-74 1-84-74 012808 繁 1-84-92 015076 概 1-85-83 015272 前 1-85-92 017132 汗 1-86-49 017186 沅 1-86-54 017699 森 1-86-86 018010 演 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5 5 5 5 5 5
012808 擎 1-84-92 015076 概 1-85-83 015272 樹 1-85-92 017132 汗 1-86-49 017186 元 1-86-54 017699 森 1-86-86 018010 減 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5 5 5 5 5 5
015076 概 1-85-83 015272 樹 1-85-92 017132 汗 1-86-49 017186 元 1-86-54 017699 森 1-86-86 018010 減 1-87-01 018948 炫 1-87-39 019137 燃 1-87-48	5 5 5 5 5
015272 樹 1-85-92 017132 汗 1-86-49 017186 元 1-86-54 017699 森 1-86-86 018010 減 1-87-01 018948 炫 1-87-39 019137 燃 1-87-48	5 5 5 5
017132 汗 1-86-49 017186 沅 1-86-54 017699 森 1-86-86 018010 滇 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5 5 5
017186 流 1-86-54 017699 森 1-86-86 018010 減 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5
017699 森 1-86-86 018010 減 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	5
018010 減 1-87-01 018948 炫 1-87-39 019137 焮 1-87-48	
018948 炫 1-87-39 019137 城 1-87-48	_
019137	5
70	5
021065 琨 1-88-07	5
- NJ	5
023689 瞪 1-88-91	5
024175	5
025556	5
026844 数 1-89-81	5
029700	5
033077	5
036819	5
039974 1-92-88	5
040280 鉀 1-93-10	5
041383 1-93-52	5
041446	5
042940 鞮 1-93-79	5
042997 鞺 1-93-80	5
043716 額 1-94-06	5
046597	5

056060	煒	1-87-54	5
000572	佾	1-14-20	4
000601	侗	1-14-23	4
004403	噻	1-15-24	4
005179	埭	1-15-50	4
005558	壔	1-15-67	4
006002	奭	1-15-74	4
008106	峴	1-47-77	4
008488	嶠	1-47-89	4
009398	庾	1-84-13	4
012778	猧	1-84-93	4
014224	曛	1-85-42	4
015679	檞	1-86-22	4
017319	泫	1-86-62	4
017609	涿	1-86-80	4
017639	淖	1-86-82	4
018328	澌	1-87-16	4
018833	濼	1-87-35	4
019883	牕	1-87-68	4
021865	畯	1-88-42	4
023307	眴	1-88-80	4
023310	眶	1-88-81	4
024110	砭	1-88-93	4
026313	篙	1-89-70	4
027246	紇	1-89-89	4
027371	絁	1-90-01	4
028727	翟	1-90-32	4
037452	跎	1-92-33	4
039301	邢	1-92-63	4
040365	銈	1-93-14	4
040857	鏽	1-93-39	4
041740	隄	1-93-60	4
047012	鵙	1-94-62	4
056422	鍈	1-93-25	4
056428	鑊	1-93-41	4
057094	氲	1-86-48	4
000420	仿	1-14-10	3

001895	刓	1-14-60	3
001999	剉	1-14-62	3
003422	呍	1-14-87	3
004797	圊	1-15-33	3
004889	圯	1-15-36	3
005167	埤	1-15-49	3
008548	皪	1-47-92	3
008622	巋	1-47-93	3
009844	彀	1-84-25	3
010354	忡	1-84-40	3
011201	憍	1-84-61	3
012624	摭	1-84-91	3
013862	昱	1-85-21	3
013939	晙	1-85-27	3
013991	晳	1-85-31	3
015094	楂	1-85-84	3
017168	汶	1-86-53	3
017323	泮	1-86-63	3
017408	洮	1-86-67	3
017809	湄	1-86-89	3
018026	滎	1-87-02	3
018174	漳	1-87-08	3
018319	澈	1-87-15	3
018814	灞	1-87-33	3
020845	玕	1-87-83	3
021324	瓈	1-88-35	3
021717	筲	1-88-41	3
022979	盌	1-88-72	3
026032	筠	1-89-63	3
026592	簳	1-89-75	3
027087	糕	1-89-87	3
027757	縕	1-90-18	3
028663	翎	1-90-30	3
029614	腊	1-90-47	3
030606	艴	1-90-60	3
034220	袪	1-91-73	3
034285	裊	1-91-74	3

034499	褧	1-91-83	3
037646	踠	1-92-36	3
038412	輭	1-92-46	3
038542	轔	1-92-48	3
038700	辵	1-92-51	3
040475	鋠	1-93-19	3
040640	鍪	1-93-30	3
042575	靚	1-93-75	3
042728	靳	1-93-77	3
043047	騹	1-93-81	3
045985	魦	1-94-34	3
046105	鮞	1-94-40	3
047714	麤	1-94-76	3
048480	鼳	1-94-84	3
048886	龢	1-94-89	3
053113	儆	1-14-42	3
056019	壒	1-15-65	3
000497	佈	1-14-14	2
002497	勻	1-14-75	2
003590	哆	1-15-02	2
003908	喆	1-15-10	2
004175	嘈	1-15-16	2
004299	噉	1-15-19	2
004562	嚬	1-15-29	2
005190	埵	1-15-51	2
005592	壠	1-15-69	2
006206	姝	1-15-80	2
006469	婾	1-15-86	2
006533	媧	1-15-89	2
006941	孖	1-47-54	2
008846	帔	1-84-09	2
008851	帘	1-84-10	2
009058	幩	1-84-11	2
009079	幞	1-84-12	2
010529	框	1-84-47	2
011940	拄	1-84-73	2
012125	捃	1-84-79	2

012494	搢	1-84-87	2
012912	擻	1-85-05	2
013040	攩	1-85-07	2
013796	昉	1-85-13	2
013903	晌	1-85-25	2
013952	晡	1-85-29	2
016275	歧	1-86-36	2
017369	洎	1-86-66	2
017412	洱	1-86-68	2
018416	澶	1-87-21	2
018939	炤	1-87-38	2
021018	琊	1-88-02	2
021067	琪	1-88-08	2
021102	瑇	1-88-16	2
021361	瓚	1-88-37	2
022395	瘞	1-88-54	2
022529	癋	1-88-58	2
022601	癤	1-88-59	2
022630	瘻	1-88-61	2
024586	礱	1-89-17	2
027276	紓	1-89-91	2
027485	綃	1-90-06	2
027856	繇	1-90-20	2
032418	斬	1-91-38	2
033079	蛼	1-91-55	2
033372	螈	1-91-60	2
033578	蟖	1-91-65	2
034523	褲	1-91-85	2
034544	襀	1-91-87	2
035069	觥	1-91-91	2
039658	都	1-92-82	2
040173	釤	1-92-94	2
040446	鋋	1-93-16	2
040808	鏞	1-93-36	2
041052	鑱	1-93-44	2
041451	閿	1-93-55	2
044912	騮	1-94-16	2

045657	竇	1-94-31	2
046171	鯁	1-94-42	2
046803	鴞	1-94-57	2
047034	鷑	1-94-63	2
050918	丰	1-14-05	2
053621	岱	1-92-61	2
056156	茁	1-90-76	2
056180	荇	1-90-82	2
056254	葳	1-91-11	2
056374	藿	1-91-37	2
056390	蘸	1-91-44	2
057018	檉	1-86-19	2
058024	槩	1-86-03	2
000382	仡	1-14-08	1
000499	佉	1-14-15	1
000789	倧	1-14-31	1
001846	刁	1-14-58	1
001899	刕	1-14-61	1
002781	丰	1-14-79	1
003302	吒	1-14-85	1
003446	呫	1-14-89	1
003540	咩	1-15-01	1
003755	唫	1-15-05	1
005419	墉	1-15-60	1
005470	墩	1-15-63	1
005747	變	1-15-72	1
006329	娣	1-15-82	1
007098	宓	1-47-56	1
007257	寘	1-47-57	1
008010	岺	1-47-73	1
008209	崧	1-47-81	1
008502	幡	1-47-91	1
009972	彤	1-84-29	1
009983	죓	1-84-30	1
010498	怵	1-84-46	1
012041	挊	1-84-76	1
012596	摠	1-84-90	1

012787 012900	擄	1-85-01	1
012900		1 00 01	1
	擷	1-85-04	1
013836	昺	1-85-16	1
013852	眗	1-85-19	1
013860	貼	1-85-20	1
014089	暠	1-85-36	1
014451	杈	1-85-50	1
014805	桲	1-85-67	1
015502	樾	1-86-11	1
015556	橛	1-86-15	1
016061	欬	1-86-30	1
016139	歆	1-86-32	1
016750	毖	1-86-43	1
017027	氐	1-86-47	1
017136	汜	1-86-50	1
017289	泔	1-86-60	1
017353	洄	1-86-65	1
017421	洹	1-86-69	1
017578	涬	1-86-79	1
017646	淛	1-86-84	1
017978	溱	1-86-93	1
018139	殭	1-87-04	1
018169	漯	1-87-07	1
018338	澔	1-87-18	1
018811	瀕	1-87-32	1
019174	煇	1-87-51	1
019291	熇	1-87-60	1
019304	熒	1-87-61	1
019519	燾	1-87-65	1
019521	爀	1-87-66	1
020431	狻	1-87-75	1
020643	獐	1-87-80	1
020861	玠	1-87-85	1
020869	玦	1-87-87	1
020874	玫	1-87-88	1
020962	珣	1-87-93	1
021015	琇	1-88-01	1

021071	琮	1-88-11	1
021073	琰	1-88-13	1
021122	瑷	1-88-18	1
021242	璜	1-88-26	1
022152	核	1-88-46	1
022297	瘂	1-88-49	1
022317	瘈	1-88-50	1
022634	癯	1-88-62	1
023167	盼	1-88-77	1
023392	睍	1-88-83	1
023523	睺	1-88-88	1
023541	瞀	1-88-89	1
024147	硃	1-89-01	1
024342	碭	1-89-10	1
024409	磤	1-89-11	1
024545	礜	1-89-15	1
025068	稃	1-89-44	1
025601	窳	1-89-53	1
025998	筎	1-89-62	1
026071	筯	1-89-65	1
026889	粔	1-89-83	1
027101	糝	1-89-88	1
027345	紱	1-89-94	1
027555	綦	1-90-09	1
028454	羗	1-90-28	1
029670	腠	1-90-48	1
029929	膻	1-90-52	1
032512	繁	1-91-43	1
032820	虵	1-91-51	1
033494	螭	1-91-62	1
034292	裎	1-91-75	1
034457	褚	1-91-82	1
034816	覔	1-91-88	1
035968	譔	1-92-18	1
035976	譙	1-92-19	1
036404	豨	1-92-23	1
036878	賸	1-92-27	1

037473	跗	1-92-35	1
037887	蹻	1-92-40	1
039413	郗	1-92-69	1
039476	郯	1-92-72	1
039630	鄧	1-92-80	1
039684	鄴	1-92-83	1
040447	鋌	1-93-17	1
040940	鐳	1-93-40	1
041022	鐮	1-93-42	1
041053	鑲	1-93-45	1
041283	閣	1-93-47	1
041330	閬	1-93-50	1
041367	閶	1-93-51	1
041430	闋	1-93-53	1
041467	闚	1-93-57	1
041650	哑	1-93-59	1
041849	隫	1-93-63	1
042230	雯	1-93-69	1
043191	轀	1-93-83	1
043269	韴	1-93-85	1
043357	頊	1-93-87	1
043463	頫	1-93-89	1
043599	顒	1-93-92	1
043600	顓	1-93-93	1
043614	顗	1-94-01	1
044779	騂	1-94-12	1
044780	騃	1-94-13	1

045597	鬠	1-94-29	1
045969	魞	1-94-32	1
046046	鮄	1-94-37	1
046071	鮏	1-94-39	1
046249	鮭	1-94-44	1
047204	鷀	1-94-66	1
048837	龔	1-94-87	1
050021	濹	1-87-25	1
053024	蒯	1-91-14	1
053212	异	1-84-18	1
056028	廋	1-84-15	1
056061	燁	1-87-62	1
056068	璘	1-88-25	1
056083	磷	1-89-14	1
056138	芾	1-90-69	1
056142	苕	1-90-72	1
056206	莩	1-90-90	1
056209	莿	1-90-91	1
056220	菪	1-91-03	1
056272	蒴	1-91-15	1
056273	蒺	1-91-16	1
056380	蘐	1-91-40	1
056382	蘘	1-91-42	1
057035	概	1-86-21	1
058230	觯	1-94-46	1
065501	籍	1-89-78	1

第4水準漢字で表現できる文字は、次の426字(表6参照)である。『太陽コーパス』での使用度数順に示す。

表 6:第4水準漢字で表現できるもの

文字番号	字形	面区点	度数
018019	滊	2-79-06	252
034969	覰	2-88-42	21
025263	穌	2-83-03	17
056318	葽	2-86-81	17

057871	熳	2-80-01	13
028121	欿	2-84-66	12
056006	儛	2-03-04	12
006573	媳	2-05-70	11
013431	斅	2-13-72	11

000763	倐	2-01-57	10
008848	帕	2-08-83	10
032958	蛃	2-87-41	10
039805	좸	2-90-33	10
025324	穠	2-83-08	9
026062	筩	2-83-48	9
033703	燭	2-87-92	9
039044	遻	2-89-93	9
007253	寖	2-08-07	8
012586	摚	2-13-41	8
012932	攄	2-13-58	8
019727	爹	2-80-13	8
025582	窬	2-83-17	8
041329	閫	2-91-56	8
050013	徤	2-12-24	8
056286	蓰	2-86-65	8
056355	薼	2-87-04	8
005805	夤	2-05-29	7
011826	扭	2-12-93	7
012694	撑	2-13-47	7
013406	斁	2-13-70	7
018777	灔	2-79-53	7
028367	罽	2-84-80	7
038473	轀	2-89-67	7
041601	陁	2-91-67	7
059756	熖	2-79-88	7
001798	凴	2-03-20	6
002689	匾	2-03-48	6
004502	嚕	2-04-45	6
011562	戕	2-12-83	6
016647	殺	2-78-04	6
018092	滹	2-79-10	6
024043	矻	2-82-28	6
028971	耬	2-85-09	6
029758	膄	2-85-45	6
036786	賒	2-89-12	6
037149	趕	2-89-23	6

040542	錕	2-91-07	6
041362	閴	2-91-57	6
056323	蕞	2-86-82	6
000964	傖	2-01-77	5
003874	啡	2-04-08	5
004350	噞	2-04-39	5
008458	嶒	2-08-63	5
008477	嶙	2-08-66	5
011216	憗	2-12-72	5
015945	檷	2-15-85	5
018062	滙	2-79-07	5
019179	煊	2-79-84	5
023532	睽	2-82-11	5
028167	罄	2-84-70	5
032882	蚜	2-87-34	5
033873	蠲	2-88-02	5
033968	䘏	2-88-04	5
037648	踢	2-89-38	5
041425	闈	2-91-59	5
044024	飣	2-92-46	5
044179	餖	2-92-55	5
045359	髠	2-93-19	5
051023	壳	2-05-22	5
002281	劘	2-03-32	4
005746	夔	2-05-28	4
006527	媢	2-05-68	4
007811	屨	2-08-20	4
011960	拕	2-13-03	4
014483	杓	2-14-34	4
015347	槖	2-15-30	4
018239	潚	2-79-21	4
020774	獷	2-80-55	4
021777	畎	2-81-26	4
027368	終	2-84-17	4
027532	網	2-84-33	4
027775	繅	2-84-51	4
032833	虺	2-87-29	4

034257	裀	2-88-18	4
034299	裒	2-88-19	4
035873	謭	2-88-73	4
035934	譃	2-88-74	4
037750	蹋	2-89-44	4
039260	選	2-90-04	4
040001	醨	2-90-40	4
040957	鐻	2-91-44	4
042577	靛	2-91-94	4
042957	鞱	2-92-08	4
043881	颷	2-92-35	4
043921	颶	2-92-38	4
044278	餲	2-92-63	4
044453	饜	2-92-73	4
048477	鼷	2-94-69	4
053339	瞢	2-82-16	4
058529	蟵	2-87-81	4
065717	社	2-88-09	4
065718	袮	2-88-10	4
000680	俏	2-01-52	3
003579	咿	2-03-85	3
003808	啊	2-04-05	3
004974	坨	2-04-68	3
006214	姣	2-05-49	3
006215	姤	2-05-50	3
006720	嬀	2-05-81	3
006776	嬴	2-05-84	3
006932	子	2-05-87	3
007058	宄	2-05-93	3
007804	屧	2-08-19	3
008172	崐	2-08-45	3
009125	幫	2-08-92	3
011158	慼	2-12-68	3
012096	挲	2-13-11	3
014580	枰	2-14-44	3
014759	桅	2-14-64	3
015449	樝	2-15-45	3

015660 括 2-15-62 3 3 016856 找 2-78-12 3 3 020286				
1020286 3元 2-80-60 3 3 3 3 3 3 3 3 3	015660	檑	2-15-62	3
3020895 末 2-80-66 3 3 3 3 3 3 3 3 3	016856	搖	2-78-12	3
78	020286	狁	2-80-30	3
78	020895	玷	2-80-66	3
2-83-16 3 3 029175 1	022453	瘳	2-81-69	3
19 19 19 19 19 19 19 19	025126	稞	2-82-92	3
1968 1978 2-85-75 3 3 3 3 3 3 3 3 3	025561	窣	2-83-16	3
1	029175	聵	2-85-14	3
19 2-87-54 3 3 3 3 3 3 3 2 5 5 3 3 3 3 3 3 3 3	030449	艋	2-85-75	3
1988 1988 2-87-58 3 3 3 3 3 3 3 3 3	032583	蘼	2-87-21	3
1988 2-87-66 3 3 3 3 3 3 3 3 3	033147	蜟	2-87-54	3
035660 記 2-88-65 3 3 3 3 3 3 3 3 3	033202	蜺	2-87-58	3
37865 3	033320	蝲	2-87-66	3
1044316 1	035660	諗	2-88-65	3
045653	037865	蹯	2-89-49	3
045681	044316	餼	2-92-67	3
045861 離 2-93-32 3 3 3 046890 幅 2-94-14 3 3 3 3 3 3 3 3 3	045653	鬫	2-93-28	3
046890	045681	鬳	2-93-29	3
047195 1	045861	魋	2-93-32	3
048261 電 2-94-62 3 3 3 3 3 3 3 3 3	046890	鵂	2-94-14	3
057479 日本日本 2-94-68 3 3 057773 日本日本 2-85-72 3 3 3 3 3 3 3 3 3	047195	鸍	2-94-32	3
057773 腕 2-85-72 3 3 078863	048261	黿	2-94-62	3
078863 流 2-89-31 3 3 3 002190 腕 2-03-30 2 2 003701 唉 2-03-94 2 2 003910 階 2-04-13 2 2 004069 階 2-04-16 2 2 005253 埃 2-05-01 2 2 006568 姓 2-05-69 2 2 006701 焼 2-05-80 2 007051 學 2-05-91 2 2 009127 幅 2-08-93 2 010319 ぱ 2-12-30 2	057479	齬	2-94-68	3
Mathematical Representation Mat	057773	艉	2-85-72	3
1003701 100	078863	疏	2-89-31	3
003910 嗜 2-04-13 2 004069 嗜 2-04-16 2 005253 埃 2-05-01 2 006232 垣 2-05-52 2 006568 基 2-05-69 2 006701 煉 2-05-80 2 007051 掌 2-05-91 2 009127 情 2-08-93 2 010319 忒 2-12-30 2	002190	劂	2-03-30	2
004069 塔 2-04-16 2 005253 埃 2-05-01 2 006232 垣 2-05-52 2 006568 基 2-05-69 2 006701 嫌 2-05-80 2 007051 型 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	003701	唉	2-03-94	2
005253 埃 2-05-01 2 006232 垣 2-05-52 2 006568 基 2-05-69 2 006701 無 2-05-80 2 007051 型 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	003910	喈	2-04-13	2
006232 垣 2-05-52 2 006568 基 2-05-69 2 006701 無 2-05-80 2 007051 章 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	004069	嗒	2-04-16	2
006568 基 2-05-69 2 006701 無 2-05-80 2 007051 掌 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	005253	堠	2-05-01	2
006701 無 2-05-80 2 007051 章 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	006232	姮	2-05-52	2
007051 学 2-05-91 2 009127 情 2-08-93 2 010319 式 2-12-30 2	006568	媱	2-05-69	2
009127 幅 2-08-93 2 010319 式 2-12-30 2	006701	嫵	2-05-80	2
010319 2-12-30 2	007051	孿	2-05-91	$\overline{2}$
	009127	幬	2-08-93	2
011106 協 2-12-67 2	010319		2-12-30	2
	011106	慠	2-12-67	2

014758	喜 肾 身 火 光	2-12-73 2-13-04 2-13-31 2-13-68	2 2 2
012359 013286 月 014758	_	2-13-31	
013286 身 014758 材	_		2
014758	_	2-13-68	
<u> </u>	光	2 10 00	2
015295 精		2-14-63	2
	盍	2-15-27	2
017283	仂	2-78-33	2
017786	Ė	2-78-82	2
018920 オ	ţ	2-79-63	2
019337	逢	2-79-92	2
019935 生	刃	2-80-18	2
020730	僉	2-80-49	2
022631	部	2-81-77	2
023466 불	星	2-82-07	2
024392 1	鬼	2-82-48	2
024495 位	車	2-82-58	2
024713 ក្	麦	2-82-70	2
026299	1	2-83-63	2
026314	Ê	2-83-65	2
026458 🎉	箑	2-83-71	2
026801	差	2-83-79	2
027489	更	2-84-30	2
028006	進	2-84-58	2
028143	并	2-84-68	2
030107 身	具	2-85-57	2
030564	虜	2-85-82	2
034945	見	2-88-41	2
035344		2-88-57	2
036168	黨	2-88-84	2
036207 🗳	含	2-88-88	2
038099	果	2-89-55	2
039794 4	M	2-90-32	2
040296	生	2-90-60	2
040937	蜀	2-91-42	2
041945 身	ち	2-91-79	2
043812	伎	2-92-33	2
044199	肴	2-92-58	2

044480	饟	2-92-74	2
044967	驄	2-93-03	2
045535	鬐	2-93-24	2
045543	鬒	2-93-25	2
047507	鸝	2-94-49	2
048274	鼂	2-94-63	2
048361	鼙	2-94-67	2
053340	懵	2-12-81	2
053365	澘	2-79-24	2
056074	瘼	2-81-68	2
056234	萹	2-86-38	2
056311	蔴	2-86-74	2
056387	蘧	2-87-18	2
056394	虀	2-87-23	2
056433	韡	2-92-10	2
056462	棻	2-14-86	2
057354	箐	2-83-52	2
057486	酹	2-90-37	2
066039	艗	2-93-77	2
000107	\	2-01-08	1
000510	佔	2-01-36	1
000725	俲	2-01-55	1
000843	偎	2-01-65	1
001051	僇	2-01-85	1
001189	儈	2-03-01	1
001797	凳	2-03-19	1
002307	茄	2-03-35	1
002380	勑	2-03-41	1
003254	旦	2-03-67	1
003387	呃	2-03-72	1
003437	呦	2-03-74	1
003535	咦	2-03-80	1
003585	峒	2-03-86	1
003790	啁	2-04-02	1
003804	商	2-04-04	1
003889	啽	2-04-11	1
004138	嗩	2-04-27	1

004194	嘐	2-04-31	1
004625	囅	2-04-51	1
004926	坌	2-04-65	1
004988	坳	2-04-70	1
005248	堞	2-04-94	1
005433	墔	2-05-14	1
005698	夆	2-05-25	1
006461	婺	2-05-63	1
006575	媵	2-05-71	1
006600	媿	2-05-73	1
006655	嫠	2-05-76	1
006734	嬈	2-05-82	1
006912	孌	2-05-86	1
007433	尃	2-08-13	1
008186	崫	2-08-46	1
008267	嵂	2-08-53	1
008532	嶲	2-08-68	1
008549	嶹	2-08-71	1
008702	巤	2-08-77	1
008843	帒	2-08-82	1
009134	幮	2-12-01	1
009311	庥	2-12-03	1
010320	忓	2-12-31	1
010403	忼	2-12-37	1
010661	悕	2-12-46	1
010718	悰	2-12-49	1
010933	偈	2-12-59	1
011021	愺	2-12-63	1
011244	憨	2-12-75	1
012013	拽	2-13-05	1
012050	挍	2-13-07	1
012054	挐	2-13-08	1
012101	挵	2-13-12	1
012127	捄	2-13-16	1
012148	捎	2-13-17	1
012238	掄	2-13-23	1
012265	挣	2-13-24	1

012510	搯	2-13-40	1
012587	摛	2-13-42	1
012678	撇	2-13-46	1
012802	擋	2-13-50	1
012857	擤	2-13-55	1
014013	晻	2-14-09	1
014266	曬	2-14-21	1
014495	柿	2-14-36	1
014936	棖	2-14-94	1
014990	棷	2-15-05	1
015378	槳	2-15-39	1
015851	櫬	2-15-75	1
015868	櫳	2-15-78	1
016517	殛	2-78-01	1
016983	氊	2-78-15	1
017085	Ÿ	2-78-17	1
017157	汭	2-78-24	1
017320	泬	2-78-39	1
017397	酒	2-78-45	1
017398	洧	2-78-46	1
017546	涘	2-78-59	1
017593	涴	2-78-67	1
017827	湏	2-78-88	1
018177	漶	2-79-16	1
018489	濚	2-79-36	1
018709	瀹	2-79-46	1
018881	灾	2-79-59	1
019018	烜	2-79-73	1
019069	煮	2-79-75	1
019210	煞	2-79-86	1
019410	燋	2-80-03	1
019774	牂	2-80-15	1
020430	狺	2-80-36	1
020495	狰	2-80-40	1
020646	獒	2-80-47	1
020754	獯	2-80-53	1
020804	玁	2-80-56	1

020857	玞	2-80-62	1
021060	琤	2-80-78	1
021248	璡	2-81-02	1
022154	痏	2-81-46	1
023224	眗	2-81-91	1
023867	稍	2-82-20	1
024396	磌	2-82-49	1
024399	磎	2-82-50	1
024643	祊	2-82-65	1
024693	祧	2-82-69	1
025407	朰	2-83-11	1
025440	寉	2-83-13	1
026057	筦	2-83-47	1
026136	箑	2-83-53	1
027370	絀	2-84-18	1
027377	絢	2-84-19	1
027438	絪	2-84-25	1
027636	緗	2-84-41	1
027660	緦	2-84-43	1
027854	繅	2-84-55	1
027960	繳	2-84-56	1
028027	類	2-84-60	1
028370	罾	2-84-81	1
028639	翃	2-84-90	1
028860	耊	2-85-03	1
028880	耑	2-85-06	1
029508	脞	2-85-33	1
030206	舄	2-85-62	1
030438	艅	2-85-73	1
030477	艏	2-85-77	1
032852	蚍	2-87-32	1
033137	蜙	2-87-53	1
033208	蜾	2-87-59	1
033213	蝀	2-87-60	1
033268	蝘	2-87-63	1
033570	蟎	2-87-80	1
033682	蠁	2-87-90	1

033745	蠓	2-87-93	1
034106	衩	2-88-11	1
034353	裱	2-88-25	1
035750	諼	2-88-66	1
036120	曹	2-88-81	1
036435	豭	2-89-03	1
037475	跙	2-89-28	1
037617	踔	2-89-35	1
038190	軔	2-89-59	1
038721	辻	2-89-76	1
038988	潰	2-89-92	1
039822	酡	2-90-34	1
039926	醎	2-90-39	1
040031	醮	2-90-41	1
040106	釃	2-90-44	1
040184	釮	2-90-50	1
040205	鈁	2-90-51	1
040216	鈊	2-90-52	1
040291	鉊	2-90-59	1
040351	鉼	2-90-71	1
040506	鋹	2-91-03	1
040556	錝	2-91-08	1
040642	鍫	2-91-18	1
040770	鏁	2-91-32	1
041663	陡	2-91-68	1
041895	隳	2-91-77	1
042864	鞚	2-92-05	1
043917	颻	2-92-37	1
043964	飇	2-92-41	1
043973	飋	2-92-42	1
044022	飡	2-92-45	1
044202	餛	2-92-59	1
044637	駉	2-92-82	1
044665	駔	2-92-83	1
044677	駙	2-92-84	1
044928	騶	2-93-02	1
045124	骯	2-93-07	1
·			

鰶	2-93-73	1
鰷	2-93-74	1
鱦	2-93-94	1
鳲	2-94-02	1
鶏	2-94-23	1
鷃	2-94-34	1
鵩	2-94-44	1
麨	2-94-55	1
鼁	2-94-65	1
鼗	2-94-66	1
齁	2-94-72	1
齅	2-94-73	1
齕	2-94-76	1
昌	2-08-79	1
亹	2-01-20	1
帮	2-08-86	1
倜	2-01-59	1
嫚	2-05-74	1
蔾	2-86-79	1
嬴	2-87-91	1
鄺	2-90-29	1
脯	2-80-17	1
暐	2-14-11	1
暱	2-14-16	1
	· 鰷 軈 乃 鴉 鷃 鹏 数 龍 製 熟 腕 乾 后 月 日 長 月 日 長 日 長 日 長 日 長 日 長 日 長 日 日 日 日	 線 2-93-74 線 2-93-94 鳥 2-94-02 鳴 2-94-23 鳥 2-94-34 変 2-94-55 電 2-94-65 シ 2-94-66 身 2-94-72 泉 2-94-73 む 2-94-76 日 2-08-79 空 2-01-20 井 2-08-86 傷 2-01-59 臭 2-86-79 麻 2-80-17 麻 2-90-29 麻 2-80-17 中 2-14-11

056062	燕	2-80-07	1
056069	瓞	2-81-11	1
056112	芃	2-85-90	1
056116	羊	2-85-91	1
056179	荄	2-86-16	1
056229	萑	2-86-34	1
056281	蓏	2-86-59	1
056326	蕡	2-86-83	1
056375	蘀	2-87-13	1
056458	鱝	2-93-84	1
057196	砉	2-82-32	1
057313	聃	2-85-11	1
057315	聱	2-85-13	1
057368	籯	2-83-81	1
057591	虓	2-87-25	1
057637	髒	2-93-15	1
057641	骷	2-93-09	1
057850	菔	2-86-29	1
065806	欽	2-90-48	1
067184	鰦	2-92-61	1
079011	黟	2-93-71	1
079145	路	2-89-41	1
079566	潔	2-79-39	1

X0213 非漢字で表現できる文字は、次の 38 字(表 7 参照)である。『太陽コーパス』での使用度数順に示す。

表 7: X0213 非漢字で表現できるもの

文字番号	字形	面区点	度数
063004	â	1-09-56	51
063017	é	1-09-63	30
063020	ê	1-09-64	21
063070	ŭ	1-10-68	13
063003	ä	1-09-58	10
063034	î	1-09-68	10
063068	ü	1-09-81	10

062960	Ś	1-10-05	9
063018	è	1-09-62	8
063065	ţ	1-10-55	6
063069	û	1-09-80	5
063051	Ö	1-09-76	4
063060	Ś	1-10-16	4
069682	7	1-03-28	4
063002	à	1-09-54	3

063014	Ç	1-09-61	3
062833	æ	1-09-60	2
062845	œ	1-11-10	2
063048	ñ	1-09-71	2
063052	ô	1-09-74	2
063053	Ŏ	1-08-87	2
063055	Ō	1-09-94	2
063221	Ι	1-13-21	2
062588	S	1-06-57	1
062952	Ô	1-09-43	1
062968	Ü	1-09-50	1
063001	á	1-09-55	1

063005	ă	1-10-41	1
063009	å	1-09-59	1
063023	ē	1-09-93	1
063030	ĥ	1-10-65	1
063031	í	1-09-67	1
063033	ï	1-09-69	1
063049	Ó	1-09-73	1
063066	ú	1-09-79	1
063222	II	1-13-22	1
063223	Ш	1-13-23	1
063224	IV	1-13-24	1

3.2 踊字タグ

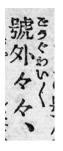
JIS X 0208 には「二の字点」や「くの字点」がないため、『太陽コーパス』では「々」「ゝ」で代用し、踊字タグを使って表現している。次の例は、平仮名「る」の繰り返し符号に使われた「二の字点」を「ゝ」で代用した例である。

(例3) 意を戰况に注がる<踊字 種類="二字点">></踊字>の致す處なりと雖も、[t189501]



また、繰り返し符号が複数連続する場合にも踊字タグが使われている。この場合は、JIS X 0208 で表現できないからタグを用いているわけではなく、繰り返し符号で表現された語が何であるのかを、タグの属性で示すのが目的である。

(例4) 號外<踊字 値="號外"> ママ</踊字> [t189511]



踊字タグの再処理結果をまとめると表 8 のようになる。JIS X 0213 を用いると、『太陽コーパス』に出現するすべての繰り返し符号を表現することができる。なお、X0208 非漢字に分類されたものは、(例 4) に類するものである。

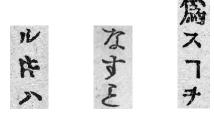
表 8: 踊字タグの再処理結果

	のベ字数	異なり字数	用例
X0208 非漢字	1,444	4	々、ヽ、ゝ、ゞ
X0213 非漢字	16,575	4	二の字点、くの字点(上、上濁点、下)
X0213 外字	0	0	
計	18,019	8	

3.3 合字タグ

『太陽コーパス』では、仮名合字等に対して、合字が表わす語形を入力し、合字 タグを付与している。合字タグの使用例を示す。

- (例 5) 郵便切手ヲ代用スル<合字>トキ</合字>ハ五厘又八壹錢切手ニ限ル [t189501]
- (例6) 其用意をなす<合字>こと</合字>肝要なるべし。[t189501]
- (例 7) 條約締盟國ハ日本國中何レノ地二於テモ通商貿易ヲ爲ス<合字>コト </合字>ヲ得領事裁判權ハ三年以内ニ之ヲ廢止シ〔t189501〕



合字タグの再処理結果をまとめると表 9 のようになる。『太陽コーパス』に出現する合字に対しては、JIS X 0213 はあまり効果がない。

表 9: 合字タグの再処理結果

	のベ字数	異なり字数	用例
X0213 非漢字	23	2	コト、より
X0213 外字	8,531	5	こと、トキ、トモ、かしこ、まゐらせ候
計	8,554	7	

3.4 小書タグ

『太陽コーパス』では、主として小書きの仮名を表わすために、小書タグを設けている。外来語表記のための片仮名が目立つ。

(例8) メルボ<小書>ル</小書>ン〔t189510〕

(例9) 亞弗利加の中心に深進せしリヴ<小書>ヰ</小書>ンストン〔t189510〕



小書タグの再処理結果をまとめると表 10 のようになる。小書タグのうち、JIS X 0213 を用いて符号化できるものは、のベ字数で 3 割程度にとどまる。なお、X0208 非漢字に分類した小書きのワは、コーパス開発時の単純なバグであろう。

 のベ字数
 異なり字数
 用例

 X0208 非漢字
 1
 ワ

 X0213 非漢字
 56
 7
 八、ヒ、フ、ヘ、ホ、ム、ル

 X0213 外字
 130
 8
 キ、コ、ヰ、ヱ、ヲ、八、忘、

 計
 187
 12

表 10:小書タグの再処理結果

なお、片仮名以外の文字について小書タグを使ったものには、次のようなものがある。(例 10)は伏せ字、(例 11)は割注に相当するものと思われる。

(例 10) 亞細亞洲にては 、 に在り、近南洋にては <小書> </小書 > <小書> </小書 大島に在り、[t189510]

(例 11) 大兒子、小兒子、大則以王、小則以覇、大小王<小書>忘</小書>覇< 小書>八</小書>兒子〔t190114〕

4. おわりに

『太陽コーパス』の文字関連タグに対する再処理結果を加え、JIS X 0213 による符号化をまとめると表 11 のようになる。のベ字数では、JIS X 0208 によるカバー率が 99.79% であったのに対して、JIS X 0213 では 99.93% となる。カバー率の増

加はわずかに 0.15 ポイント程度である。しかし、異なり字数では、JIS X 0208 のカバー率 79.16%に対して、JIS X 0213 は 92.06% と、12.90 ポイントも増加している。

『現代日本語書き言葉均衡コーパス』の書籍(生産実態サブコーパスと流通実態サブコーパス)では、JIS X 0208 JIS X 0213 によるカバー率の増加は、のベ字数で平均 0.03 ポイント、異なり字数で平均 7.75 ポイントであるから、現代語を扱うときよりも、JIS X 0213 を用いる効果が大きい。『太陽コーパス』のような近代語文献を電子化する際には、JIS X 0208 文字セットよりも、JIS X 0213 文字セットの方が適していると言えよう。

表 11:『太陽コーパス』の JIS X 0213 による符号化

	のベ字数	異なり字数
第1水準漢字	5,650,619	2,721
第 2 水準漢字	888,886	2,864
X0208 非漢字	7,798,780	318
第 3 水準漢字	2,685	446
第 4 水準漢字	1,371	426
X0213 非漢字	16,874	51
(小計)	14,359,215	6,826
X0213 外字	9,890	592
計	14,369,105	7,418

最後に、JIS X 0213 文字セットで表現できない 592 字を、『太陽コーパス』での使用度数順に表 12 に示す。

表 12: JIS X 0213 外字

文字番号	字形	度数
合字こと		8,487
小書ヰ		105
T003	卷	62
047550	墭	55
019790	牃	32
067692	蔡	27
058594	颺	24
合字まゐら	せ候	22
020571	猺	21
012384	揪	20

18
冉 17
L 14
14
14
11
11
10
10
10
10

020803	羅	9
小書ヲ	PME	9
010846	臾	8
015158	繒	8
020483	裸	8
014585	絽	7
063059	r	7
小書ヱ	. •	7
004658	囒	6
021193	璀	6
065520	維	6
074686	漆	6
T025	闃	6
013004	擅	5
020687	獐	5
059905	景	5
066170	圧	5
074461	皙	5
小書		5
合字トモ		5
001706	凘	4
002225	劊	4
011737	扃	4
025000	秪	4
028693	翖	4
030766	芾	4
034535	褺	4
035801	謏	4
037152	趋	4
037993	躝	4
042302	霉	4
044479	饞	4
048306	鼉	4
048652	齨	4
050664	養	4
058129	耦	4
077953	曹	4

080074	烈	4
000575	侁	3
001171	僿	3
004183	嘍	3
004259	嘵	3
004981	坭	3
005403	塿	3
007191	宼	3
009463	廕	3
012241	掇	3
012712	撙	3
013751	旰	3
015675	檛	3
018819	灡	3
019378	爉	3
020335	狑	3
020672	獘	3
024552	礟	3
030411	舺	3
033980	衅	3
034974	覻	3
039942	醔	3
043609	頼	3
044963	鸄	3
056123	芠	3
059815	磵	3
071309	垔	3
071516	粤	3
071768	禁	3
075506	疄	3
083904	衙	3
T002	燄	3
合字かしこ		3
001134	僰	2
001584	家	2
001953	刱	2
002041	剚	2

003478	咂	2
003673	哷	2
003806	啉	2
004202	嘕	2
004242	嘮	2
004252	囆	2
004455	嚅	2
005133	烀	2
005261	堦	2
005404	墀	2
007285	寠	2
008263	崿	2
008753	耆	2
008760	巹	2
011115	憗	2
011525	戇	2
012124	捂	2
012505	搪	2
012595	摟	2
012728	撦	2
012897	擵	2
012991	攙	2
013668	旐	2
013984	晾	2
014070	睺	2
014163	曀	2
015492	傤	2
017171	汸	2
017492	浲	2
017784	渰	2
019202	煳	2
019430	燖	2
020274	犽	2
020336	狭	2
020377	狥	2
020722	猥	2
020961	珢	2

022784	皤	2
023619	瞕	2
023633	矔	2
024669	袂	2
025778	娛	2
026318	篘	2
026455	簌	2
027483	綁	2
027831	縏	2
028824	翢	2
029426	胭	2
032114	薕	2
032411	蘃	2
033770	盡	2
035233	訏	2
035466	註	2
036843	賬	2
037787	蹗	2
039791	酛	2
040185	釯	2
040350	鉻	2
040514	錁	2
041574	阱	2
044678	駚	2
044710	駧	2
045468	鬁	2
045524	鬎	2
047440	鷽	2
049191	欝	2
049534	脴	2
050020	澚	2
050031	蕡	2
056037	撐	2
056114	市	2
056412	躉	2
057125	塍	2
057128	螣	2
·	_	·

	_	T
059535	崅	2
059665	朥	2
059875	舿	2
062581	ά	2
067039	菈	2
067352	搉	2
077217	遅	2
081515	髩	2
085829	笲	2
T007	磴	2
T016	烱	2
T020	終	2
T105	令	2
T117	鑢	2
T125	魚卵	2
T126	蹼	2
T135	40	2
T136	ch.	2
T137	Q	2
T138	chr-	2
T146	(H)	2
000484	伻	1
001119	僨	1
001154	僱	1
001505	顚	1
001955	刲	1
002094	剮	1
002217	剿	1
002683	鱼	1
002854	四1	1
003235	另	1
003322	咿	1
003368	吵	1
003382	呷	1
003477	咁	1
003628	哎	1
003756	唬	1

003811	昏	1
003940	喔	1
004041	嗃	1
004139	嗶	1
004190	嘏	1
004195	嘑	1
004348	噙	1
004979	坫	1
005080	垧	1
005208	堃	1
005323	塍	1
005326	塏	1
005425	墎	1
005467	墦	1
005560	壖	1
006124	坳	1
006235	姱	1
006361	娬	1
006392	婇	1
006405	婑	1
006579	媸	1
006787	嬝	1
007957	岝	1
007958	岞	1
008175	崒	1
008326	幎	1
008367	嵲	1
008391	嵽	1
008460	臺	1
008556	嶻	1
008579	巑	1
008758	巹	1
009052	幙	1
009081	幠	1
009245	庋	1
009605	弆	1
010394	伢	1

010675	悜	1
010751	惂	1
010766	倭	1
010839	悠	1
010937	愔	1
011167	憀	1
011225	憜	1
011475	幱	1
012131	捆	1
012316	揁	1
012453	搊	1
012491	搠	1
012672	摔	1
012693	撐	1
012951	攉	1
012957	攊	1
013627	匑	1
013664	旃	1
013685	旖	1
013686	腌	1
013691	胺	1
013780	旿	1
013963	唇	1
014045	暌	1
014817	桻	1
015232	樕	1
015356	槤	1
015718	檯	1
015903	欃	1
016017	昳	1
016476	殑	1
017166	泠	1
017169	汥	1
017179	公	1
017510	浼	1
017594	泽	1
017596	涶	1

017685	淰	1
017714	渀	1
017793	渹	1
017853	猖	1
018256	潥	1
018320	澉	1
018388	澦	1
018392	滋	1
019134	焭	1
019379	熸	1
019896	潢	1
019987	牳	1
020054	犄	1
020067	牾	1
020251	犵	1
020391	狪	1
020394	狫	1
020452	犲	1
020456	猅	1
020487	猔	1
020618	猱	1
020680	獜	1
020692	獠	1
020879	玭	1
020988	珵	1
020990	珷	1
021041	琗	1
021208	璊	1
021239	璚	1
021564	甇	1
021749	畁	1
021828	畡	1
022451	瘲	1
022685	皂	1
022787	皡	1
022815	皭	1
023120	盰	1

023178 明	-
023555 版	-
024088 何	-
024178 何 1 024200 徑 1 024249 孫 1 024330 偏 1 024369 頂 1 024568 偏 1 024613 價 1 024734 福 1 024737 福 1 025094 椏 1 025094 椏 1 025368 羅 1 025966 節 1 026109 芹 1 026132 ፫ 1 026309 第 1 026742 澤 1	-
024200 優 1 024249 係 1 024330 偏 1 024369 傾 1 024568 偏 1 024613 償 1 024734 福 1 024737 福 1 024761 煙 1 025094 煙 1 025368 電 1 025438 宅 1 026109 洋 1 026132 茂 1 026309 第 1 026742 準 1	
024249 係	
024330 偏 1 024369 傾 1 024568 偏 1 024613 微 1 024734 超 1 024737 確 1 024761 型 1 025094 粒 1 025368 確 1 025438 宅 1 025966 労 1 026109 洋 1 026132 茂 1 026309 第 1 026742 準 1	-
024369 頂 1 024568 個 1 024613 頂 1 024734 個 1 024737 個 1 024761 煙 1 025094 粒 1 025368 電 1 025438 宅 1 025966 第 1 026109 洋 1 026132 度 1 026309 第 1 026742 準 1	-
024568 日 1 1 1 1 1 1 1 1 1	-
024613 1	-
024734 超 1 024737 福 1 024761 理 1 025094 程 1 025368 程 1 025438 宅 1 025966 第 1 026109 月 1 026132 度 1 026309 第 1 026742 算 1	-
024737 福 024761 煙 025094 粒 025368 福 025438 宅 025966 第 026109 并 026132 定 026309 第 026742 算	-
024761 禪 025094 粒 025368 權 025438 笔 025966 節 026109 并 026132 箎 026309 第 026742 算	-
025094 極 025368 權 025438 笔 025966 第 026109 并 026132 분 026158 基 026309 第 026742 算	-
025368 權 025438 宅 025966 第 026109 詳 026132 別 026158 董 026309 第 026742 津	-
025438 笔 025966 第 026109 拜 026132 E 026158 筆 026309 第 026742 準	-
025966 第 026109 第 026132 第 026158 第 026309 第 026742 第	-
026109 笄 1 026132 箎 1 026158 箠 1 026309 第 1 026742 籜 1	-
026132 焼 026158 筆 026309 菊 026742 準	-
026158 筆 1 026309 第 1 026742 章 1	-
026309 第 026742 算	-
026742	-
1 -1	
026985 矮 1	-
1.~	-
027439 条 1	-
027499	-
027538	-
027701	-
028350 異 1	-
028440	-
028919	-
028978	-
029544 腹 1	-
029582 肾 1	-
029593	
029627 据 1	-
029683 睽 1	

029895	膴	1
030055	臡	1
030202	鳥	1
030279	舏	1
030472	艎	1
030818	芥	1
031589	蔝	1
031616	蒶	1
032354	藯	1
032871	蚖	1
032941	妲	1
033255	蝏	1
033359	婇	1
033367	娘	1
033415	螘	1
033612	嬌	1
033747	蟲	1
034034	衎	1
034339	裩	1
034355	棭	1
034470	裦	1
035019	觗	1
035283	暑	1
035292	訬	1
035384	詖	1
035596	誶	1
035636	諉	1
035656	諕	1
035947	譊	1
035991	譍	1
036050	濤	1
036515	豾	1
036713	貺	1
036747	賌	1
036955	賸	1
036960	贑	1
037004	赩	1

037084 数			
037303 1	037084	趂	1
037508 B	037124	趍	1
037568 一次 1 037568 一次 1 037716 一次 1 1 037806 一次 1 037847 一次 1 037881 一次 1 037886 一次 1 037962 一次 1 039296 小が 1 039851 一計 1 039874 一計 1 039998 1 040104 1 1 040162 1 1 040166 1 1 040226 1 1 040253 1 1 040262 1 040338 1 1 040437 1 0404474 1 0404474 1 0404474 1 040489 1 0404798 1 040798 1 040798 1 040921 1 1 040938 1 041008 1 041008 1 041008 1 041078 1 1 041078 1 1 041078 1 1 041078 1 1 1 1 1 1 1 1 1	037303	趙	1
037568 一次 1 037716 一次 1 1 037806 一次 1 037847 一次 1 037881 一次 1 037886 一次 1 037962 一次 1 039296 小	037508	跍	1
037716 2	037560	脠	1
037806 22	037568	踄	1
037847 職	037716	蹀	1
037881 25	037806	蹚	1
037886	037847	獙	1
037962 職	037881	蹙	1
038414 解	037886	蹺	1
039296 分	037962	躐	1
039851	038414	輅	1
039874 画	039296	邠	1
039998 1 040104 1 1 040162 3 1 1 040166 3 1 1 040226 3 1 1 1 1 1 1 1 1 1	039851	酧	1
040104 瞬 1 040162 到 1 1 040166 到 1 1 040226 到 1 1 040253 計 1 040262 計 1 040338 計 1 040369 計 1 040429 到 1 040437 到 1 040474 到 1 040489 到 1 040798 計 1 040921 到 1 040938 到 1 041008 到 1 041045 到 1 041078 到 1 1 041078 到 1 1 1 1 1 1 1 1 1	039874	酺	1
040162 到	039998	鰺	1
040166 氨 1 040226 丘 1 040253 鉢 1 040262 姉 1 040338 鉮 1 040369 竚 1 040429 迚 1 040437 愛 1 040437 愛 1 040489 鉱 1 040798 鎌 1 040865 釧 1 040921 錢 1 040938 釧 1 041008 貸 1 041078 賃 1 041078 賃 1	040104	釂	1
040226 任 1 040253 休 1 040262 姉 1 040338 鉮 1 040369 슑 1 040429 迚 1 040429 迚 1 040437 貸 1 040474 貸 1 040489 鉱 1 040798 鎌 1 040865 釧 1 040921 鍵 1 040938 釧 1 041008 貸 1 041045 賃 1 041078 賃 1	040162	釟	1
040253 休 1 040262 姉 1 040338 绅 1 040369 銌 1 040429 銼 1 040429 銼 1 040437 鉸 1 040474 鈫 1 040489 鉱 1 040798 鎌 1 040865 鐁 1 040921 錢 1 040938 釧 1 041008 貸 1 041078 貨 1 041078 貨 1	040166	釞	1
040262 姉 1 040338 婶 1 040369 슠 1 040429 銼 1 040429 銼 1 040437 錗 1 040474 殳 1 040489 鉱 1 040798 鎌 1 040865 獅 1 040921 錢 1 040938 鍋 1 041008 貸 1 041045 賃 1 041078 賃 1	040226	鈓	1
040338 轉 1 040369 第 1 040429 建 1 040437 發 1 040474 鏡 1 040489 鉱 1 040798 錐 1 040865 獅 1 040921 錢 1 040938 劉 1 041008 增 1 041078 鎖 1 041078 鎖 1	040253	鉢	1
040369 第 1 040429 建 1 040437 疑 1 040474 疑 1 040489 鉱 1 040798 鎌 1 040865 獅 1 040921 錢 1 040938 劉 1 041008 增 1 041078 鋼 1	040262	鈰	1
040429 姓 1 040437 委 1 040474 女 1 040489 鉱 1 040798 鎌 1 040865 釧 1 040921 錢 1 040938 釧 1 041008 貸 1 041078 貨 1 041078 貨 1	040338	鉮	1
040437 愛 1 040474 愛 1 040489 並 1 040798 嫌 1 040865 嫌 1 040921 遂 1 040938 貸 1 041008 貸 1 041078 貸 1 041078 貸 1	040369	銌	1
040474 競 1 040489 鉱 1 040798 鎌 1 040865 鐁 1 040921 錢 1 040938 鍋 1 041008 貸 1 041078 賃 1 041078 賃 1	040429	銼	1
040489 鉱 1 040798 鎌 1 040865 釧 1 040921 錢 1 040938 釧 1 041008 賃 1 041078 賃 1	040437	鋄	1
040798 鎌 1 040865 釧 1 040921 錢 1 040938 鍋 1 041008 增 1 041045 賃 1 041078 鍋 1	040474	鋟	1
040865 鐁 1 040921 錢 1 040938 鍋 1 041008 貸 1 041045 賃 1 041078 鍋 1	040489	銰	1
040921 鍵 1 040938 鋼 1 041008 1 041045 1 041078 1	040798	鏕	1
040938 調 1 041008 1 041045 1 041078 1	040865	鐁	1
041008 1 041045 1 041078 1	040921	714	1
041045 鏡 1 041078 鋼 1	040938	缥	1
041078	041008	鑙	1
PP PP	041045	鑬	1
041991 博 1	041078	鋼	1
U41441 1	041221	閈	1

041262	閣	1
041289	閡	1
041428	闉	1
041723	煙	1
041810	嘕	1
043414	頒	1
043598	顑	1
043605	龥	1
043813	颴	1
043931	嬺	1
043955	飅	1
044084	飵	1
044126	餁	1
044236	餧	1
044294	餵	1
044362	饈	1
044799	腕	1
044847	騘	1
044921	騳	1
044962	騄	1
045305	髗	1
045841	魈	1
046039	鮀	1
046059	鮇	1
046095	觪	1
046119	魣	1
046407	鰶	1
046856	鴋	1
047937	黈	1
048352	薿	1
049388	砉	1
050022	琑	1
050860	紿	1
050971	嫍	1
053173	摔	1
053192	䓫	1
053438	嫫	1

053607	摒	1
053614	甿	1
054924	翶	1
056017	塟	1
056035	搽	1
056080	矱	1
056118	芑	1
056132	芰	1
056150	苯	1
056248	葠	1
056357	蘅	1
056359	槧	1
056398	蜹	1
056402	蠛	1
056407	襴	1
056410	豶	1
057215	磙	1
057361	簉	1
057461	鞵	1
057530	埕	1
057851	蔻	1
058025	嘅	1
058029	彴	1
058660	蔲	1
059216	喓	1
059236	嗐	1
059247	嗸	1
059271	嚑	1
059618	灎	1
059803	硘	1
059921	螩	1
060001	滛	1
060081	屻	1
061327	糙	1
061593	尞	1
062583	療 カ o	1
062587	Ò	1
·		

062589	ΰ	1
063047	ņ	1
063084	ŷ	1
065057	朅	1
065167	哬	1
065183	耔	1
065417	疝	1
065538	膵	1
065646	蘪	1
065972	嫪	1
066007	鱘	1
066167	蝄	1
066767	鯮	1
067720	啣	1
067843	姬	1
068171	鼓	1
068275	堇	1
070734	区	1
072382	剜	1
072508	廢	1
075331	虔	1
076687	縤	1
077626	秖	1
078439	摩	1
078618	錇	1
079135	燈	1
082432	暨	1
082959	巹	1
085185	癥	1
085262	涿	1
086272	沔	1
086719	吰	1
096371	Ō	1
096381	D	1
096461	₩	1
096551	*	1
T001	嘘	1

T004	瑂	1
T005	曙	1
T006	喜	1
T008	踱	1
T009	電	1
T010	华	1
T011	媬	1
T012	酒產	1
T013	寧	1
T015	怒	1
T017	倏	1
T018	辦	1
T019	液	1
T021	表	1
T022	變	1
T023	遐	1
T024	釰	1
T102	罋	1
T103	圩	1
T104	膠	1
T106	怒	1
T107	屧	1
T108	祭	1
T109	瑭	1
T110		1
T111	*	1
T112	姄	1

T113	囊	1
T114	鋅	1
T115	芸蕉	1
T116	滙	1
T118	稲	1
T119	繙	1
T120	桥	1
T121	*	1
T122	湉	1
T123	们	1
T124	妙	1
T127	蘉	1
T128	噹	1
T129	岖	1
T130	儒	1
T131	腹	1
T132	漉	1
T133	獯	1
T134	\Rightarrow	1
T139	\$	1
T140	グー	1
T145	M	1
小書八		1
小書キ		1
小書コ		1
小書忘		1

文 献

池田証寿・白井純・高田智和(2002)「宋版漢字字体の処理」『京都大学大型計算機センター第 69 回研究セミナー報告 東洋学へのコンピュータ利用』pp.49-62

下田正弘・師茂樹(1999)「大正新脩大蔵経データベース(SAT)における外字問題」 『人文学と情報処理』25, pp.35-43

須永哲矢・堤智明・高田智和(2011)「明治前期雑誌の異体漢字と文字コード」『人文科学とコンピュータシンポジウム論文集 「デジタル・アーカイブ」再考 いま 改めて問う記録・保存・活用の技術 』pp.381-388

- 高田智和(2002)「漢字処理と『大字典』」、『訓点語と訓点資料』109, pp.99-107
- 高田智和・小林正行・間淵洋子・大島一・西部みちる・山口昌也(2009)『JIS X 0213:2004 運用の検証(国立国語研究所内部報告書 LR-CCG-09-01)』国立国語研究所
- 高田智和(2011)「現代日本語コーパスにおける文字処理」『人間文化研究情報資源共有化研究会報告集』2, pp.31-40
- 田中牧郎(2005)「漢字の実態と処理の方法」『雑誌『太陽』による確立期現代語の研究 『太陽コーパス』研究論文集 (国立国語研究所報告 122)』博文館新社,pp.271-292
- 當山日出夫(2009)「『内村鑑三全集』デジタル版の文字処理について」『東洋学へのコンピュータ利用第 20 回セミナー』京都大学人文科学研究所附属漢字情報センター, pp.5-18
- 富田倫生(2000)「青空文庫と外字」、『人文学と情報処理』26, pp.23-30 安永尚志(1998)『国文学研究とコンピュータ』勉誠社