

朝日新聞における辞書非掲載漢字の出現状況

横山 詔一* 笹原 宏之*
YOKOYAMA Shoichi SASAHARA Hiroyuki

野崎 浩成** 米田 純子*
NOZAKI Hironari YONEDA Junko

国立国語研究所*
The National Language Research Institute
115 東京都北区西が丘3丁目9番14号
3-9-14, Nishigaoka, Kita-ku, Tokyo, 115

名古屋市立大学**
Nagoya City University
467 名古屋市瑞穂区瑞穂町山の畑1
Yamanohata 1, Mizuho-cho, Mizuho-ku, Nagoya, 467

Email : [yokoyama, sasa, jyoneda]@kokken.go.jp, hnozaki@ccews1.cc.nagoya-cu.ac.jp

【1】はじめに

JIS漢字には、読みや意味や出典がはっきりしない音義未詳字、通称「幽霊文字」が含まれているとの指摘が漢字・辞書研究者らの間でなされてきた。本研究は、1993年版朝日新聞記事全文データベース(以下、CD-HIASK'93(朝日新聞社・紀伊國屋書店・日外アソシエーツ、1994))を活用して、約11万件の記事のなかに「幽霊文字」や1978年以前の漢字辞書には掲載されていない「辞書非掲載漢字」がどの程度の頻度で出現するかを調査した。

今回、調査対象に新聞メディアを選択した理由は以下の3点による。

1. 現代日本において流通している文字の実態を包括的・網羅的に精査しようとする場合、大量データを含んだ電子化されたコーパスが必要となる。朝日・毎日・日経といった新聞社が発行している記事全文データベースは、比較的安価で市販されており、しかもCD-ROM化されている。
2. 字数が大量であり、さまざまな分野の用字を含むうえに、多くの人々の目に触れ、かつ影響力の大きいものは、新聞であろう。その発行部数は、『読売新聞』のように1000万を超える

ものもある。また、日本国内に限らず、海外にも我が国の新聞記事はインターネット経由で流れしており、日本語を解する世界中の人々に読まれているようである。

3. 新聞記事全文データベースに格納されているテキストデータを分析するには、著作権者である新聞社の使用許諾が必要である。今回、我々は朝日新聞社からCD-HIASK'93のテキストデータを光磁気ディスクに複写・蓄積する許諾を得たので、朝日新聞の記事を分析することにした。

文字調査にあたっては、上で述べたように、紙媒体上の記事だけではなく、今や新聞社各社のWebのホームページから発信されている記事も射程におさめたうえで議論すべき時代が到来したといつても過言ではない。その意味で、電子メディアとしての新聞記事CD-ROMを研究対象とすることは、世界中に流通している日本語の文字の実態を探る糸口としてたいへん重要であろう。

さて、『朝日新聞』は国内第2位の発行部数を誇っており、その紙面に、いわゆる「拡張新字体」(常用漢字表外字に表内字の新字体を準用した字体。当用漢字以前から存在した略体も含む)を大幅

に採用しているほか、固有名詞での異体字の使用を制限していたなど、他紙とは一線を画した漢字使用を行っている点で注目される（笹原・横山・米田・野崎, 1997）。

CD-HIASK'93 に入っている新聞記事データは、通信社からの記事のほか、テレビ欄や天気予報欄、広告欄などを含んでおらず、紙面のすべてを包含するものではない。しかし、収録字数は、1 年間 分で約 5500 万字と他の文字資料に比べて桁違いに多い。また、異なり字種も 4400 を超えており、3000 代に落ちつくことの多いサンプリングによる漢字頻度調査と比較して、異例の数を示している（笹原ほか, 1997）。

本研究の目的は、JIS 漢字（「JIS X 0208」(1978 初版 1997 改正)に収める漢字）における「幽霊文字」を初めとする従前の漢和辞典に掲載されていなかった「辞書非掲載漢字」に焦点を絞って CD-HIASK'93 と紙面（縮刷版、朝日新聞社：1993 年 1 月～12 月）との比較照合を行い、「幽霊文字」や「辞書非掲載漢字」の新聞紙面上での挙動を追跡するシステムを開発することである。

【2】調査対象

CD-HIASK'93 における 5462 万余字を対象とし、そのうち検索する漢字は次のような基準により選択した。

1. JIS 漢字における幽霊文字

「幽霊文字」とは、JIS 漢字（1978）選定の際に、原典から転記する際の誤写により生じたとみられる漢字字体のこと（笹原, 1997）。なお、次の説明は、矢印記号を挟んで、左から順に「JIS 漢字

の原典」「JIS 漢字表」「他の資料」を示すものである（笹原ほか, 1997）。

- **広義**：字体が過去に存在しなかったもの
？ → 「蜀」
- **狭義**：字体が過去に存在したもの
同字の場合…暗合
閨 → 「閨」 ← 閨
別字の場合…衝突
□(木+品) → 「楓」 ← 楓

2. 幽霊文字と辞書非掲載漢字

- 本研究では、幽霊文字を「JIS 漢字選定時に、原典から転記する際に誤写されて変化・派生した、とみられる字体」と定義する。ただし、その字種が漢和辞典にある字体と暗合・衝突した字体については調査を継続中である（笹原ほか, 1997）。
- このような事情に配慮して、やや広めに『新字源』や『大漢和辞典』にない字体をも視野に入れて、129 字を選び、それを調査対象とした。その結果、純然たる広義・狭義の幽霊文字のほかに、漢和辞典に載せられることのほとんどなかった「辞書非掲載漢字」（国字・異体字を含む）や、頻度の低い漢字をも対象に含めた（笹原ほか, 1997）。

【3】方 法

縮刷版で辞書非掲載漢字を検索するには、当該の漢字を含む K W I C とその出典情報が必要である。K W I C プログラムの開発は、awk、sed、Perl などのスクリプト言語を用いて行い、データ処理はワークステーションを活用した（野崎・笹原・米田・横山, 1997）。

本プログラムは、調査対象となる文字

(以下、「ターゲット文字」と呼ぶ)を読み込み、入力されるテキストファイル中に出現するターゲット文字について、KWICを作成する。KWICには、記事番号、見出し、行番号などの出典情報が付加されており、新聞紙面上の記事との照会が可能である。入力データとなるテキストファイルはテキスト形式でありさえすれば、任意の文書の取り扱いが可能である。ターゲット文字には、先に述べた幽霊文字と辞書非掲載漢字の候補を併せて129字用いた。図1にKWIC作成プログラムの処理手順を示す(野崎ほか, 1997)。

【4】結果と考察

1. 非出現字

上記の129字について、CD-HIASK'93を検索したところ、出現しなかった字が95種ある。これらは、いずれも国立国語研究所が1966年に実施した新聞に対する漢字調査(1976)でも出現していない。むろん、他の資料には用例があるものがほとんどである(笹原ほか, 1997)。

寫 世 吁 啼 嘸 囉 噴 壇 埋 地
塙 塚 塗 堤 壇 墻 壤 壤 僕 崔
峩 峰 嵐 嵐 嵐 嵐 嵐 嵐 嵐 嵐
坤 捺 捺 捺 捺 捺 捺 捺 捺 捺
榾 榾 榾 榾 榾 榾 榾 榾 榾 榾
榶 榶 榶 榶 榶 榶 榶 榶 榶 榶
榷 汗 汗 汗 汗 汗 汗 汗 汗 汗
榶 笛 笛 笛 笛 笛 笛 笛 笛 笛
榷 舶 舶 舶 舶 舶 舶 舶 舶 舶
榷 酒 酒 酒 酒 酒 酒 酒 酒 酒
榷 駄 駄 駄 駄 駄 駄 駄 駄 駄

2. 出現字

一方、CD-HIASK'93に出現した字は34種(353回)ある。国立国語研究所の

1966年調査でも出現していた字は「築」「礦」(石+廣)「蒐」「峓」「馬」の5種(笹原ほか, 1997)。

事 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇
𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇 𠂇
簀 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵
簀 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵 𣎵

ちなみに、第1水準の2560「礦」と、第2水準の6672「礦」とは、1978JISでは正字体と拡張新字体の関係であったが、1983年のJIS漢字改正において、入れ替えがあった字体である。そのため、「礦」(石+廣)が調査の対象であったが、「礦」(石+廣)も確認のために合わせて調査した。

3. ゲタ文字

本研究ではJIS外字の文字・記号を表すのに用いられることが多い「ゲタ文字：=」も検索してみた。ゲタ文字の総度数は876あり、紙面と対照したところ以下のようになった。

漢字：延べ856、異なり43

記号：延べ20、異なり6

ゲタ文字に置換された漢字の一覧表から、CD-HIASK'93には不正ゲタ文字とでも言うべきものが浮き彫りになった。

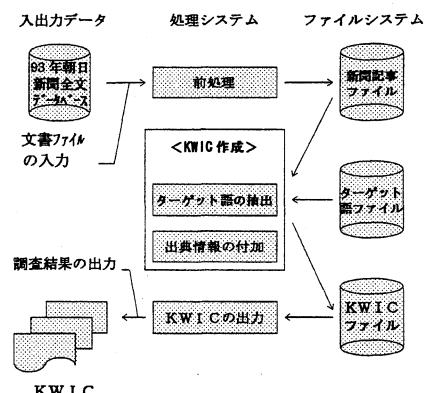


図1. KWIC作成プログラムの処理手順

その典型例を示す（横山・笹原・米田，1997）。

● 不正ゲタ文字の例

83JISに存在するのに、なぜかゲタ文字になっているケース。

槙原投手 → =原投手

遙かな → =かな

「槙」は度数230で相当目立つ。ゲタ文字に置換された文字のなかで出現順位が第2位である。この事実はCD-HIASK'93を扱う際に注意すべき側面を如実に示したものである（横山ほか，1997）。

【5】まとめ

本研究は、辞書非掲載漢字（幽霊文字を含む）が新聞記事のなかでどのような振る舞いをしているのかを悉皆的に精査した。

その結果、国立国語研究所の1966年調査で出現していた5種の漢字（「築」「礦」（石+廣）「蒐」「峠」「馬」）以外に、29種の漢字が実際の新聞紙面上に登場したことが明らかになった。これら併せて34種の漢字の出現頻度は350を越えており、これまでの文字調査では得られなかった新たな資料を手にできた。

またゲタ文字の調査により、電子メディアを資料として扱う際に注意を払うべきことがらの一端が明らかになった。

本研究の知見は、今後の漢字符号化研究の基礎資料として役立つことが期待される。

引用文献

- 朝日新聞社・紀伊國屋書店・日外アソシエーツ（1994）『CD-HIASK'93』
- 国立国語研究所（1976）『現代新聞の漢字』、集英出版
- 笹原宏之（印刷中）「字体に生じる偶然の

一致」『日本語科学』創刊号、国立国語研究所

- 笹原宏之・横山詔一・米田純子・野崎浩成（1997）「文字資料としての『朝日新聞』紙面とCD-ROM—「JIS X0208」における辞書非掲載漢字を中心にして」、シンポジウム人文科学における数量的分析(2)、文部省統計数理研究所
- 野崎浩成・笹原宏之・米田純子・横山詔一（1997）「幽霊文字調査のためのKWIC作成プログラムの開発」、シンポジウム人文科学における数量的分析(2)、文部省統計数理研究所
- 横山詔一・笹原宏之・米田純子（1997）「朝日新聞CD-ROMに出現するゲタ文字の分析」、シンポジウム人文科学における数量的分析(2)、文部省統計数理研究所

参考文献

- 横山詔一・野崎浩成（1996年3月）「朝日新聞CD-ROMによる漢字頻度基準表の作成と数量分析」、人文科学における数量的分析シンポジウム、文部省統計数理研究所、pp.11-14
- 横山詔一・野崎浩成（1996年9月）「コーパスを利用した日本語環境の分析」、日本行動計量学会第24回大会特別セッション、pp.138-139
- 野崎浩成・横山詔一（1996年9月）「新聞と雑誌における漢字使用頻度の分析－心理学での材料統制の観点から－」、日本行動計量学会第24回大会、pp.266-267
- 横山詔一・野崎浩成・米田純子（1996年9月）「新聞の漢字使用順位に影響する要因の分析」、計量国語学会第40回大会、p.7
- 野崎浩成・横山詔一ほか（1996年9月）「漢字使用頻度の時代的变化に関する考察」、計量国語学会第40回大会、p.8

＜謝辞＞

本研究は、横山・野崎・米田に対する以下の文部省科研費の援助を受けた。①「国際社会における日本語についての総合的研究」（創成的基礎研究、代表者：水谷修）、②「インターネットにおける学術漢字の符号化に関する基礎的研究」（重点領域研究、代表者：斎藤秀紀）、③「海外日本語学習リソース提供システムの実験研究」（国際学術研究、代表者：柳澤好昭）。

また、研究の推進に当たり、朝日新聞社電子電波メディア局著作権チームマネジャーの杉野信雄氏には多大なるご理解と貴重なご助言をいただいた。

紙面とCD-ROMの綿密な照合作業は佐藤朋子氏と藤川美穂氏に、83JIS漢字符号については太田幸代氏にお手伝いいただいた。ここに記して厚く感謝申し上げる。