

『日本語歴史コーパス』における原文KWIC表示機能の実装

著者	小木曾 智信, 岡 照晃, 中村 壮範, 八木 豊
雑誌名	言語資源活用ワークショップ発表論文集
巻	2
ページ	252-257
発行年	2017
URL	http://doi.org/10.15084/00001526

『日本語歴史コーパス』における原文 KWIC 表示機能の実装

小木曾 智信 (国立国語研究所言語変化研究領域)・岡照晃 (国立国語研究所コーパス開発センター)・中村壮範 (マンパワーグループ株式会社)・八木豊 (株式会社ピコラボ)

Implementation of “Original Text KWIC” Display Function in the Corpus of Historical Japanese

Toshinobu Ogiso (NINJAL), Teruaki Oka (NINJAL), Takenori Nakamura (Manpower Japan Co., Ltd.), Yutaka Yagi (Picolab Co., Ltd.)

要旨 日本語史研究の基礎資料は、残された文献に見られる用例であるが、その原文は今日一般に用いられる表記とは大幅に異なる形である場合が少なくない。例えば、『万葉集』は万葉仮名で、キリシタン資料は当時のポルトガル語式のローマ字で表記されている。こうした資料をコーパスとして形態論情報を付与し、現代人に読みやすいものとするためには、原文を校訂して漢字ひらがな交じりにした読み下し本文を用意する必要がある。一方で、読み下し本文では失われてしまう情報も少なくないため、用例には原文を併せて表示することが求められる。『日本語歴史コーパス』では従来、原文情報を保持しつつ必要な修正を行った上で形態論情報を付与して公開してきたが、原文情報の提供方法は限定的だった。今回新たに、コーパス検索アプリケーション「中納言」上で、原文の前後文脈付きで検索結果を表示できる機能を実装した。本発表ではこの原文 KWIC 表示機能について述べる。

1. 『日本語歴史コーパス』における「原文」

過去の時代の日本語を研究するにあたっては、当時使われた用例がほとんど唯一の手がかりであり、それがどのように書かれているのかは日本語研究者にとってきわめて重要な情報である。用例の確認は、根本的には一次資料である原本そのものやその写真・画像にあたることができれば良いが、その一方で、原本のままでは現代人には読みづらく検索ができないため、現行の活字にそのまま直した翻字本文や、表記を読みやすく改め誤りを正した校訂本文が必要とされる。校訂済みの本文は、現代人にとって読みやすいだけでなく、『日本語歴史コーパス』のように形態素解析を施して単語の情報等を付与する際にも適している。このように、一つの日本語史研究資料であっても、原文画像、翻字本文、校訂本文とさまざまな段階があり、それぞれが研究上で必要とされる価値を持っている。

『日本語歴史コーパス』(小木曾 2016) の「平安時代編」では、小学館『新編日本古典文学全集』(新編全集) の校訂本文を底本としている。ここでは、コーパスにとっての「原文」は校訂本文が唯一のものである。ところが、「鎌倉時代編 I 説話・随筆」の『今昔物語集』では、本文が漢字カタカナ交じりであるだけでなく、部分的に漢文の語順で書かれ、返り点が表示されている。そのためこうした作品をコーパス化するにあたっては、底本である新編全集の本文をさらに改変し、形態素解析が可能な通常の語順の漢字ひらがな交じり文に直す必要があった(富士池・田中 2012, 富士池ほか 2013)。また、「明治・大正編 I 雑誌」では、当時出版された雑誌そのものを底本としたため、自ら本文校訂を行う必要があり、そのために本文を修正したほか、漢字カタカナ交じり文の記事はひらがなに直した上でコーパス化を行っている(近藤 2016)。

『日本語歴史コーパス』のテキストは XML で構造化・タグ付けされており、以上のよう

な本文修正については元の様態を保存し再現できるようにタグ付けがなされている。このタグによって再現される元のテキストを「原文文字列」と読んでいる。

2. これまでの「中納言」と原文表示

『日本語歴史コーパス』は Web 上のコーパス検索アプリケーション「中納言」を通じて提供されているが、ここでの検索結果として表示される本文 (KWIC の前後文脈) は、形態素解析の対象となった校訂済みの本文である。したがって、修正前の本文については文脈からは確認できない。しかし、日本語研究のための資料として、調査対象の用例についてはできるかぎり原態を示したいため、検索結果の表に「原文文字列」という列を設け、ここでキーとなった語の「原文」の様態を確認できるようになっている (中納言バージョン 2.2.2.2, 図 1)。

サンプル ID	開始位置	連番	コア	前文脈	キー	後文脈	語彙素読み	語彙素	品詞	原文文字列	振り仮名
60M明六 1875_39001	10800	7490	1	法律の 條款 其 幾 千萬 なるを 知ら ず#其 東洲 に 在 て 唐律 を 以 て	粗	其 備 は れる 者 と し 我 が 國 往 昔 の 律 法 亦 之 に 本 づく#而 て 明律 あり#清律 あり	ホボ	略	副詞	粗	
60M明六 1874_17001	5500	3750	1	其 詳 を 論 ぜ ん に 民選 議院 の 事 は 世 上 に 公論 あり て 其 制 も 亦	畧	分 明 な れ が 今 又 贅 論 せ ず#憲 に は 只 其 財 務 に 預 かる 要件 のみ を 掲 ぐ	ホボ	略	副詞	畧	
60M明六 1874_18006	4010	2730	1	従 て 諸 を 作 り 諸 を 案 じ て 調 を 爲 す の 法 な り#此 法 支 那 に は	畧	に れ 有 り#歐 米 諸 國 に は 殆 ど 精 妙 を 極 む#只 我 邦 に 未 だ 開 闢 す#今 之 を	ホボ	略	副詞	畧	
60M明六 1875_33004	8960	5790	1	に 來 入 し 其 勢 を 變 ぜ ず し て 若 干 歳 月 を 經 び 終 に は 金 紙 の 間	畧	平 均 を 成 す に 至 る べ し#一 た び 平 均 を 成 す に 至 ら ば 是 は 紙 幣 を 燒 却 する	ホボ	略	副詞	畧	
60M明六 1875_41003	5250	3510	1	ん や#勇 を 養 ふ は 武 に あり#我 邦 風 習 太 古 よ り 武 を 重 ん じ 名 を 惜 む#	畧	日 耳 曼 に 似 た り#王 室 中 古 の 文 弱 は 上 たる 者 武 を 鄙 し も こ よ り 而 して 武 士 の	ホボ	略	副詞	畧	
60M明六 1874_05004	14210	9460	1	を 壓 服 する 極 めて 難 く 人 民 常 に 自 由 の 氣 質 を 保 持 せ り#亞 非 利 加 も 亦 其 氣 候	畧	亞 細 亞 の 南 方 と 似 た る を 以 て 其 人 民 の 隷 従 する 概 ね 亦 相 同 じ く 又 亞 米 利 加	ホボ	略	副詞	畧	/ \
60M明六 1874_04001	10480	7220	1	歐 羅 巴 に 於 て は 氣 候 既 に 漸 を 逐 て 強 と 強 と 相 隣 し 其 力	畧	相 比 敵 する を 以 て 甲 國 よ り 乙 國 を 一 舉 し て 壓 服 する に 極 めて 難 し#是 れ 即 ち	ホボ	略	副詞	畧	

図 1 これまでの「中納言」の「原文文字列」表示

前後文脈まで含めた原文は提供されておらず、機能は限定的である。それでも、これまでに公開してきた資料については機能的に十分であったと言える。

3. 「万葉集」「キリシタン資料」と原文

今回、新たに『日本語歴史コーパス』に『万葉集』とキリシタン資料『天草版平家物語』『エソポのハブラス (伊曾保物語)』を追加することとなった。これらの資料は、原文が漢字仮名交じり文ではないため、原文と書き下した校訂本文との差が甚だしく、従来の枠組みでは扱いきれない。

「奈良時代編 I 万葉集」として収録される『万葉集』の原文は周知の通り漢字だけの万葉仮名で書かれており、通常はこれを漢字仮名交じり文に書き下したものを読んでいる。次に例を示す。

金野乃 美草苺茸 屋杼礼里之 兔道乃宮子能 借五百礮所念
秋の野の み草刈り茸き 宿れりし 宇治のみやこの 仮廬し思ほゆ

(7 番歌)

熟田津尔 船乗世武登 月待者 潮毛可奈比沼 今者許芸乞菜
熟田津に 船乗りせむと 月待てば 潮もかなひぬ 今は漕ぎ出でな

(8 番歌)

許等ゝ波奴 樹尔波安里等母 宇流波之吉 伎美我手奈礼能 許等尔之安流倍志
言とはぬ 木にはありとも 愛しき 君が手馴れの 琴にしあるべし

(811 番歌)

許等騰波奴 紀尔茂安理等毛 和何世古我 多那礼乃美巨騰 都地尔意加米移母
言とはぬ 木にもありとも 我が背子が 手馴れの御琴 地に置かめやも

(812 番歌)

また、「室町時代編IIキリシタン資料」として収録される予定の『天草版平家物語』『エソポのハブラス』の原文は当時のポルトガル語ローマ字で書かれており、これも漢字仮名交じり文に書き下したものとともに利用されている。次に例を示す。

VManojô. Qêgueônobô, Feiqe no yurai ga qiqitai fodoni, ara ara riacu xite vo catari are.

QIICHI. Yafui coto de gozaru : vôtata catari maraxôzu.

右馬の允. 検校の坊, 平家の由来が聞きたいほどに, あらあら略してお語りあれ.

喜一. やすいことでござる: おほかた語りませうず.

(平家物語 巻第一)

EVROPA no vchi Phrigiatoyû cunino Troia toyû jôrino qinpeni Amoniato yû fatoga vogiaru. Sono fatoni nauoba Efopoto yûte, yguiô fuxiguina jintaiga vogiattaga, fono jidai Europano tencani cono fitoni mafatte minicui monomo vorinacattato qicoyeta.

エウロパの中ヒリジヤといふ国のトロヤといふ城裡の近辺にアモニヤといふ里がおぢやる。その里に名をばエソポというて、異形不思議な仁体がおぢやったが、その時代エウロパの天下に、この人にまさって醜い者もおりなかったと聞えた。

(エソポのハブラス エソポが生涯の物語略)

これらの資料において、上段に示した「原文」と下段に示した形態素解析対象となる読み下し本文は、既存のサブコーパスのように形態素解析等のために本文を校訂してカタカナをひらがなに直したといったレベルではなく、全く異なる文字種によるテキストとなっている。個々における「原文」は、研究上の利用価値が高く、漢字仮名交じり文では落ちてしまう貴重な情報を含んでいる。

たとえば、『万葉集』の例で言えば、当該例が一字一音の仮名で書かれているのか、漢字を訓読した例なのか、あるいは助詞等を補読したものなのか、という違いは、用例の価値を大きく左右するものである。また、原文が音仮名で書かれていれば上代歴史仮名遣いを確認することも可能であるが、こういった情報は漢字仮名交じりの本文では落ちてしまっている。また、キリシタン資料の原文では、ローマ字によって当時の音形が確認でき、特にオ列長音の開合の別が「ô」「ò」で示されていたりするが、漢字仮名交じりの本文ではこうした情報も確認できない。

このようなことから、『万葉集』やキリシタン資料にとっては前後文脈まで含めて原文が参照できることが望まれる。とくに、漢字仮名交じり文のテキストと対照する形で参照できることが望ましい。

4. 原文 KWIC 表示機能

『万葉集』とキリシタン資料のコーパスの構築にあたっては、当初より原文と形態素解析対象の本文とを別に用意し、それぞれを関連づけるアライメントを行ってきた（山田ほか2015、Oka and Kono 2016）。原文と、形態素解析の対象となった漢字仮名交じりの本文、さらに形態素解析結果である短単位情報は、コーパスを格納した「形態論情報データベース」（小木曾・中村 2014）上でファイル頭からのオフセット値によって相互に関連付けられている。これにより、個々の単語について、前後文脈の原文テキストを出力することが可能になっている。

図2は、新しく公開予定の「中納言」上で前後文脈の原文を表示した例である。従来の漢字平仮名交じりの本文の KWIC（前文脈・キー・後文脈）の下段に、原文の KWIC（原文前文脈・キーの原文文字列・原文後文脈）を表示し、原文を形態素解析対象となった漢字仮名交じりテキストと対照しながら閲覧することが可能になった。検索結果のダウンロード時には、それぞれを別の列としたタブ区切りテキスト形式のデータとしてダウンロードされる。



図2 公開予定の「中納言」の原文 KWIC 表示機能（開発中の画面）

この機能の提供により、新しく公開される『万葉集』とキリシタン資料のデータの利用の幅が大きく広がることになるはずである。

5. 「原文」をめぐる注意点

このようにして提供される「原文」情報について、いくつか利用にあたって注意を要する点がある。

一つ目は、作品・サブコーパスごとに「原文」とされているものの実態が大きく異なることである。それぞれの中身を整理したものを表 I に示す。もともとの資料の性質が大きく異なるうえ、サブコーパスによって底本も違うためやむを得ないことであるが、利用に際しては注意が必要である。たとえば、「平安時代編」に含まれる作品の原文は、原点にまで戻れば、大部分が仮名からなる崩し字で書かれたテキストが原文であるが、底本を新編全集とする『日本語歴史コーパス』ではそこまで遡ることはできない。

二つ目は、漢字平仮名交じりの本文と原文とが、一対一に対応するとは限らないということである。たとえば、『今昔物語集』においては「未」が「未だ〜ず」と読まれるような“再読文字”がある。この場合、原文の一文字が、本文中の離れた2箇所に対応することとなる。返り点が入るような“返読”の箇所でも、対応が2箇所に分かれる場合がある。また、同じよ

表 1 『日本語歴史コーパス』の「原文」

サブコーパス・資料	原文	本文 (漢字ひらがな交じり)
奈良時代編 I 万葉集	新編全集 (原文・万葉仮名)	新編全集 (読み下し文)
平安時代編	新編全集の校訂済み本文 (共通)	
鎌倉時代編 I	今昔物語集	新編全集の本文 (漢字カタカナ交じり)
	その他	新編全集の校訂済み本文 (共通)
室町時代編 I 狂言	『大蔵虎明能狂言集翻刻註解』のテキスト	濁点付与など一部のみ校訂したもの
室町時代編 II キリシタン資料	原典のローマ字	ローマ字から生成した漢字仮名交じり文
明治・大正編 I 雑誌	原典をテキスト化したもの (一部漢字カタカナ交じり、踊り字あり)	原テキストを校訂した漢字仮名交じり文

うな漢文的表記で、読み下した場合に対応する読みがない“置字”がみられることがあり、この場合には原文の文字に対応する本文がないことになる。この逆の場合として、『万葉集』などで多く見られる“補読”がある。たとえば原文「金野乃」を本文「秋の野の」と読むとき、一つ目の「の」は原文に対応する文字がない。

以上のような一対一対応しないものについてもコーパスのデータベース上では問題なく格納されているが、「中納言」上での実際の利用にあたっては対応部分が見当たらなかつたり複数あつたりするために注意を要する。図3は原文 KWIC 部分を拡大したもののだが、枠

10-萬葉 0806_00006	30710	この道にてしおしるや難波の海と名付けけらしも#土(や)も 望しく(ある)	べき	古代に語り継ぐべき名は立てずして#我が背子が着る衣薄し佐保風
10-萬葉 0806_00019	11020	直超乃此徑尔呂師押照哉難波乃海跡名附家良思慕士也母空 母の命凡るかに心に尽して思ふらむその子なれやもますら や望しく(ある)	べき	有万代尔語徳可名者不立之而吾背子我著衣薄佐保風者疾莫吹及家左右
10-萬葉 0806_00020	62850	知智乃夷乃父能美許等波播麻葉乃母能美己等於保呂可尔情尽而念 良牟其子奈礼夜母大夫夜无奈之久 めづらしもかくしにぞ現し明らめ秋立つにに#天地を離らす 日月の極みなく(ある)	可	在尔須重布理於許之投失毛知千尋射和多之劍刀許尔等理波彼安之比奇能八葦布美越左之麻久流情不 降後代乃可多利都具倍久名乎多都倍志母
30-今昔 1100_12014	13950	静めて念ひを専らにして(私)を説き奉らば、必ず其の(利益)は (有)	べき	也と(有)語り伝へたるにや。
30-今昔 1100_12024	17340	然レバ、人若シ急難ニ値ハム時ハ、心ヲ静メテ念ヒテ等ニシテ仏ヲ念 ジ奉ラバ、必ず其ノ利益ハ有	べき	也トナム語り伝ヘタルトヤ。
30-今昔 1100_12034	18510	差せり。#夏の事なれば、土葬(也)と云へども、少(も)香(は)(有) 其ノ上ニニ都婆ヲ起テ、釘括ヲ差セリ。夏ノ事ナレバ、土葬也ト云ヘ モ、少(モ)香(ハ)	べき 可	に、(露)其の臭き香無し。#其の後、七日毎に仏経ヲ供養 有(キ)ニ、露其ノ臭キ香無し。其ノ後、七日毎ニ仏経ヲ供養ス。
60-明六 1874_01002	9180	に(抱)て(乗)せ(枕)事(に)そ(何)なる(べき)事(に)か(有)ら(ぬ)。(有) 若シ、辭シテ參ザラムヤ、強ニ馬ニ抱テ乗セム事コノ何レベキ事カ 有ラム。極テ恐シ	べき 可	事(何)なり(と)思(ひ)臥(伏)たる(に)、上長押ヨリ鼠(ノ)走(渡)る(に)、枕上(に)物 有(キ)事(力)ナト思(ひ)臥(伏)タルニ、上長押ヨリ鼠(ノ)走(渡)る(に)、枕上(に)物 ヲ握(リ)テ見(レ)バ、紙(ノ)破(也)。
60-明六 1874_05005	17010	其内(に)は學者輩出して(洋字)を以(て)和漢(ノ)史(傳)等(を)記(す)る(者) 者(有)る	べき	なれども要する(に)二重(の)傍(た)る(に)を(免)か(ず)是(其(不)利(の)三(なり)此
60-明六 1874_20001	10480	て(全く)同(位)同(等)の(者)となす(が)故(に)又(他)教(他)派(を)容(忍)する(制 度)の(有)る	べき	ナトモ要スルニ二重ノ勞タルヲ免カレズ是其不利ノ三ナリ此三不利ヲ冒シテ從來未ダアラザル處ノ奇法ヲ行 ハント欲スルハ三尺ノ童子ト雖トモ其至難ナルヲ知ルベシ
60-明六 1875_31001	19280	總テ諸教諸ニ共ニ憲法上ニ於テ全く同位同等ノ者トナスカ故ニ又他教 他ニ容(忍)スル(制度)ナル	べき	理(絶)へ(て)は(れ)ら(ず)臣民(ノ)論(述)の(自由)禮(拜)の(自由)其(他)總(て)神(道)ニ(關)係(せ 理(絶)へ(て)之(レ)アラス臣民(ノ)論(述)ノ自由(禮)拜(ノ)自由(其)他(總)テ神(道)ニ(關)係(セ)ル事(ノ)自由(ハ)決(シ)テ政府(ヨリ)授(與)セラ ル(特)權(ニ)アラス
		吾人(な)べて(之)を(信)ふ(に)信(ず)る(は)人情(ノ)自然(に)して(人)道(論) に(然)る(有)る	べき	に(は)り(夫)ノ罪(惡)を(罰)する(は)政府(ノ)職(掌)なり(然)れ(ども)律(法)隨(善)惡(滅)免(ノ)例(あり
		夫(レ)善(惡)ハ吾(人)ナベテ之(ヲ)言(フ)ヲ(憚)ル(ハ)人情(ノ)自然(ニ)シテ人(道)論(ニ) 然(ル)アル	べき	に(は)り(夫)ノ罪(惡)ヲ罰(ス)ルハ政府(ノ)職(掌)ナリ然(レ)トモ律(法)隨(善)惡(免)ノ例(アリ
		する(者)な(れ)ば(本)來(の)理(に)於(て)は(人民)上(に)あり(て)政府(下)に (有)る	べき	が(ゆ)に(な)れ(ども)政府(ハ)人民(を)保(護)する(の)大(權)を(掌)握(せ)ざる(可)ら(ざる)は(以
		(人民)主(ニ)シテ政府(ハ)人民(ノ)爲(メ)ニ存(在)スル(者)ナレハ本(來)ノ理(ニ)於 テハ人民(上)ニアリテ政府(下)ニアル	べき	加(ゆ)クナレトモ政府(ハ)人民(を)保(護)スルノ大(權)を(掌)握(せ)ザル可(ラ)サルヲ以(テ)上(位)ヲ占(ム)ルヲ要(ス)ルナリ

図 3 原文 KWIC が本文と一対多で対応する例

で囲んだ中の「べき」は、原文「可有キ」を本文で「有るべき」と読んでいるため、原文の2箇所に対応する。このような場合には最も前方の対応箇所を「キー」としてとり、後方の対応箇所は原文後文脈中の色つきの括弧で囲んで示す仕様となっている。

5. おわりに

以上のように、原文と読み下し本文との乖離が甚だしい日本語史資料のコーパス公開にあたって、コーパス検索アプリケーション「中納言」に原文 KWIC 表示機能を実装する。これにより、『日本語歴史コーパス』の利用の幅が一段広がることとなった。今後、こうした機能の活用により、コーパスなしでは困難であった新しい研究が実現することを期待したい。

謝 辞

本研究は国立国語研究所の共同研究プロジェクト「通時コーパスの構築と日本語史研究の進展開」、および科研費基盤(A)「日本語歴史コーパスの多層的拡張による精密化とその活用」による成果の一部である。

文 献

- 富士池優美・田中牧郎(2012).「今昔物語集の返読文字について—形態素解析の前処理を通して—」、日本語学会 2012 年度春季大会予稿集、pp.223-228
- 富士池優美・河瀬彰宏・野田高広・岩崎瑠莉恵(2013).「『今昔物語集』のテキスト整形」『第4回コーパス日本語学ワークショップ予稿集』, pp.125-134.
- 小木曾智信・中村壮範(2014).『現代日本語書き言葉均衡コーパス』形態論情報アノテーション支援システムの設計・実装・運用, 自然言語処理, 21(2), pp.301-332.
- 山田祐実・大村舞・鴻野知暁・Kevin Duh・小木曾智信・松本裕治(2015).「万葉集を対象とした原文と読み下し文のアライメント」『第8回コーパス日本語学ワークショップ予稿集』, pp.243-252
- 小木曾智信 (2016).『日本語歴史コーパス』の現状と展望, 國語と國文學, 93(5), pp.72-85.
- Teruaki OKA, Tomoaki Kono (2016). Original-Transcribed Text Alignment for Manyosyu Written by Old Japanese Language, *Language Technology Resources and Tools for Digital Humanities (LT4DH)*, (http://researchmap.jp/mukbfhtwa-2098193/#_2098193 より閲覧可能)
- 近藤明日子(2016).「『明六雑誌コーパス』『国民之友コーパス』の構築—形態論情報を付与した近代雑誌コーパスの設計—」日本語の研究, 12(4), pp.167-174

関連 URL

『日本語歴史コーパス』 http://pj.ninjal.ac.jp/corpus_center/chj/
 コーパス検索アプリケーション「中納言」 <https://chunagon.ninjal.ac.jp/>