

『現代日本語書き言葉均衡コーパス』の構築

1. 『現代日本語書き言葉均衡コーパス』とは

『現代日本語書き言葉均衡コーパス』は、国立国語研究所が構築を進めている現代日本語の大規模なデータベースです。最終的な規模は1億語以上、本年度から2010年度までの五か年で構築を終え、2011年の春には一般公開する予定です。構築作業の一部は文部科学省科学研究費補助金特定領域研究「日本語コーパス」の補助によって実施しています。

2. KOTONOHA計画

『現代日本語書き言葉均衡コーパス』は、近現代の日本語全体をデータベース化しようとするKOTONOHA計画の一部として位置づけられています。図1はそのKOTONOHA計画の全体像を示しています。図の中央には時間軸が走っており、明治から現代までの時間をあらわしています。また時間軸の上部は書き言葉、下部は話し言葉に該当します。書き言葉を代表するジャンルには「書籍」「新聞」「雑誌」「ウェブ」を、また話し言葉のジャンルとして「モノログ」「対話」「雑談」を認めています。

KOTONOHA計画は、今後数十年の間に一連のコーパスを構築することによって、近現代の日本語の全体像を可能な限り広く、また歪みのない形で記録に残すことをめざしており、これまでに『太陽コーパス』、『日本語話し言葉コーパス』のふたつを公開してきました。図1左上の「太陽」と記された楕円と右下の「CSJ」(Corpus of Spontaneous Japanese)と記された楕円がこれに該当します。また図の下部には、近現代の日本語史における重要な出来事が記入されています。

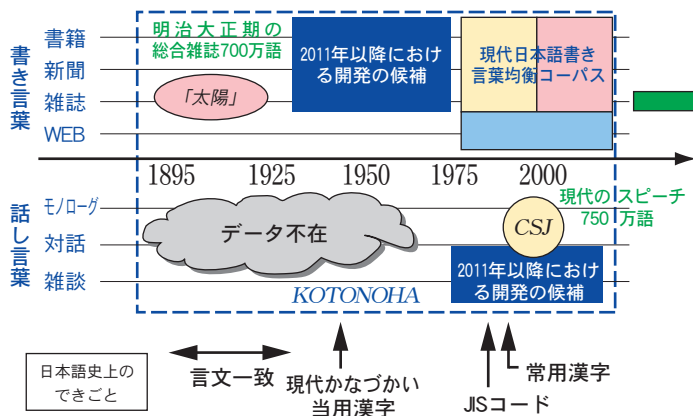


図1 KOTONOHA計画

明治・大正期の「言文一致」運動、終戦直後の「現代かなづかい」と「当用漢字」の告示(1946)、そして1980年前後におけるJIS漢字コードの制定(1978, 1983)と常用漢字の告示(1981)です。『現代日本語書き言葉均衡コーパス』はJISコード制定後の情報化社会における書き言葉のコーパスです。

3. 『現代日本語書き言葉均衡コーパス』の構造

『現代日本語書き言葉均衡コーパス』は三種類のサブコーパスから構成されています。図2左上の「生産実態サブコーパス」は、2001年から2005年の間に出版されたすべての書籍、雑誌、新聞を母集団として、そこから約3500万語を無作為に抽出したコーパスです。テキストの内容による選別は行っていません。図2右上の「流通実態サブコーパス」は東京都下の公共図書館に収蔵されている書籍を母集団として約3000万語をやはり無作為に抽出したコーパスです。最長で過去30年間分の書籍が対象となること、一定数(例えば10館以上)の公共図書館に所蔵されている図書のみが対象となる点で、生産実態サブコーパスの書籍部分とは異なっています。最後に「非母集団サブコーパス」は、特定の母集団を設定することなく、国立国語研究所が実施する研究に必要な書き言葉データを格納しています。政府が刊行する白書、法律、国会の会議録、検定教科書、ベストセラーなどを予定しています。また現代の書き言葉の著しい特徴であるインターネット(WWW)上の書き言葉のデータも重要な対象です。非母集団サブコーパスの規模は約3500万語を予定しています。

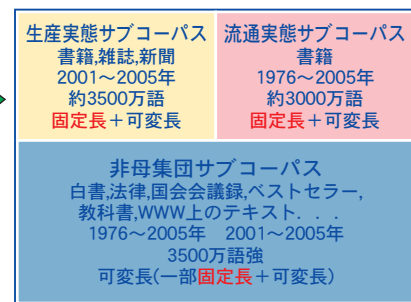


図2 『現代日本語書き言葉均衡コーパス』

4. コーパスの構築方法と進捗状況

以下では『現代日本語書き言葉均衡コーパス』を構築するプロセスを書籍を例にとって説明します。このような手順に沿って、2007年3月時点で、約1500万語相当のサンプルの著作権処理が終了し、そのうち500万語分の電子化が終了しています。

5. 関連URL

KOTONOHA計画

<http://www.kokken.go.jp/kotonoha/>

科学研究費特定領域研究「日本語コーパス」

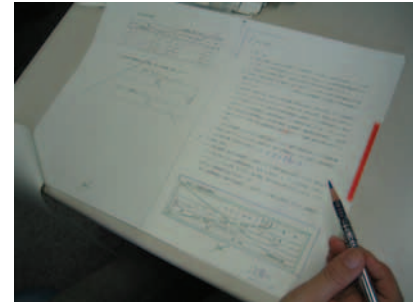
<http://www.tokuteicorpus.jp/>

(前川 喜久雄・山崎 誠)

2001~2005年に出版された書籍のリスト(母集団)

「〇〇全集」1巻	205頁
「△△殺人事件」	100頁
「板橋区史」	600頁
「もうかる投資」	80頁
「ファール昆虫記」	200頁
「××自叙伝」	300頁
……	
等々	数十万冊

無作為に数万ページを抽出



サンプルとなるテキストの決定

```
<sample id="OW1X_00000" version="0.1">
  <article>
    <hierarchy>
      <title>
        第1部 内外均衡に向かった昭和53年度経済<br/>
      </title>
      <titleBlock>
        <title>
          第1章 昭和53年度の日本経済<br/>
        </title>
        <titleBlock>
          <title>
            一その推移と特徴一<br/>
          </titleBlock>
          <body>
            <title>
              第2節 内外均衡の背景<br/>
            </title>
            <paragraph>
              <sentence> 53年度中にみられた内外均衡回復に向けての動きは、
              <sentence> 以下では、それらの動きの重要な背景として、<circle des
              をとりあげてみよう。</sentence>
            </paragraph>
          </body>
        </titleBlock>
      </hierarchy>
    </article>
  </sample>
```

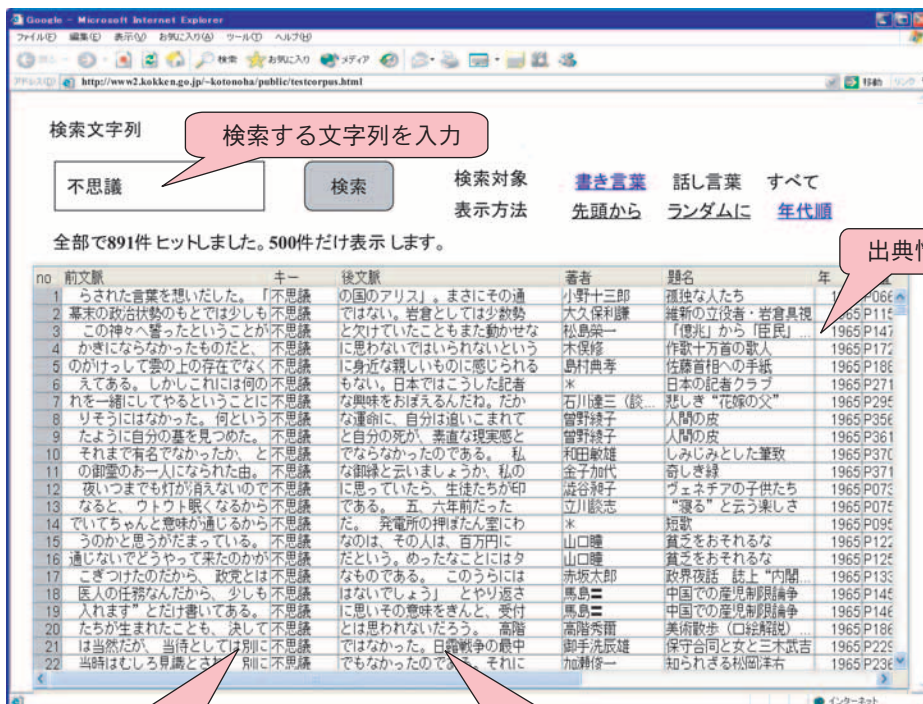
一律1000字の固定長サンプルと、節や章に対応する可変長サンプル(平均4000字程度)の二種類を作成

著作権処理

著作権者から利用許諾を取得

サンプルをコンピュータに入力。文字や文書構造に関する情報をXML(データ記述言語)で表現。さらに品詞情報や文節間の係り受け関係についての情報を付与。

サーバーに搭載してインターネット上で公開



先行する文脈(20字程度)

後続する文脈(20字程度)

新「ことば」シリーズ20『文字と社会』

国立国語研究所編、ぎょうせい刊

■新「ことば」シリーズとは

日本語に関する興味・関心を一般の方々に持っていただく目的で、国立国語研究所は『新「ことば」シリーズ』という普及書を毎年約6万部発行し、ほぼ全国の学校や公民館に無償配布しています。2007年3月30日発行の第20号は「文字と社会」を取り上げます。



本書は、一般書店で購入できます。
新「ことば」シリーズ20『文字と社会』
ぎょうせい、A5判112ページ 定価500円(税込)

■この号のテーマ

私たちの暮らしや社会は「文字」に支えられています。日本全国で、様々な文字が用いられています。文字は人間の「脳」や「心」にも影響を及ぼしているようです。ゲーム機に搭載された漢字ドリルが高齢者の皆さんに歓迎されていることから明らかなように、文字に対する世の中の関心は高いものがあります。

この号は「文字と社会」のほかに「人間」とのかかわりをめぐる話題も提供し、心豊かな言語生活を楽しむために作られました。

文化や学術研究の話題に加えて、全国の行政窓口などで行われている実務（戸籍業務など）にも直結した解説を準備しました。

■内容と執筆陣は

1. 巻頭エッセイ「私の文字生活」
阿刀田高氏（作家）によるエッセイ。
2. 座談会「放送現場の文字、心に伝わる文字、強い言葉」
文字は生きて変化しつつある「ことばの小宇宙」です。その不思議な魅力についての鼎談^{てい}を掲載しました。
・『パカの壁』の著者である養老孟司氏（解剖学者）
・桜井洋子氏（NHKアナウンサー）
・杉戸清樹（司会、国立国語研究所長）
座談会で出た話題をいくつか紹介しましょう。

【問】NHKアナウンサーが読んでいるニュース原稿は、「縦書き」でしょうか、それとも「横書き」でしょうか？

【問】手書きの効用は？

【問】「クオリア」とは何でしょうか？

（答えは座談会の記事のなかにあります。）

3. 問答形式の解説「ことば事情」

「放送と漢字」（柴田実氏：NHK放送文化研究所）、「教科書の文字」（小椋秀樹：国立国語研究所）、「公共サービスの文字」（高田智和：国立国語研究所）の3本立て。

「放送と漢字」は、次のような疑問に答えます。

【問】最近のテレビ放送は、話している内容が文字で書かれたり、説明が文字になっていたりして文字が使われることが多くなったような気がしますが、時代とともに変わっているのでしょうか？

「教科書の文字」は以下のような話題を取り上げます。学校教育の現場で、漢字学習の指導などに役に立つものと思います。

【問】小学校で、どういう順序で漢字を習うかについて決まりはあるのでしょうか？

【問】印刷された文字と手書きの文字とでは、字形に違いがありますが、印刷された文字と同じように書かないといけないのでしょうか？

「公共サービスの文字」は住民基本台帳や戸籍の文字について国立国語研究所が行った最新の研究成果を紹介します。

【問】「戸籍の文字を「正字」に改めます」という通知をもらったことがあります。これは何のことでしょうか？

【問】「戸籍統一文字」や「住基統一文字」とは何でしょうか？

4. 問答集「ことばの質問BOX」

次のような15の質問に新野直哉（国立国語研究所）など数名の専門家がお答えします。

【問】小学生の娘が作文で「きょうわたのしかったです。」と書いたところ、「わ」に×をつけられてしまいました。なぜ「わ」と書いてはいけないのでしょうか。

【問】「ヤマダハルオ」のローマ字表記は、「Haruo Yamada」がよいのでしょうか。「Yamada Haruo」がよいのでしょうか。

5. コラム

「文字と失語症」（横山詔一：国立国語研究所）と「オーノ氏の懊悩に満ちた選択」（宇佐美洋：国立国語研究所）の2本。

6. 外国人から見た日本語「ことばと社会」

ドイツー日本研究所の研究者であるP. バックハウズ氏が、東京の山手線28駅で調査したデータを「公共文字と日本の多言語化－東京の言語景観を事例に」に掲載しました。

あわせて、ドイツー日本研究所でのインタビュー記事も御覧ください。

（横山 詔一）

「ことば」フォーラム報告

第30回「日本語の中の外来語と外国語 ～ 新聞、雑誌、テレビ」

■ 140名が参加

第30回「ことば」フォーラムが、2007年2月24日（土）の午後、国立国語研究所2階講堂で開催されました。参加者は140名でした。

■ マスメディアで使用される外来語と外国語の現在と将来を考える

前半は次の3件の講演がありました。

①朝日新聞社校閲センターの福田亮氏は「新聞記事の外来語」について、新聞で使用された外来語の現状を中心に述べられました。なかでも、「新語・新用法の定着するまで」、「現れては消えていく外来語」での、見落とされがちな外来語の生態という視点は新鮮でした。また戦前から現在までの「朝日新聞の外来語使用に対する態度」では新聞社の表記上の配慮に関する報告がありました。

②伊藤雅光は「雑誌に見られる外来語と外国語」について、研究所の語彙調査データに基づいて、外来語と外国語の過去と現在の比較データを紹介し、さらに将来の状況について考えました。この40年間における雑誌の外来語と外国語の増加率はめざましく、将来この傾向が助長されるであろうことを予想し、外来語だけではなく、外国語に対しても今のうちに目配りをしておく必要性について述べました。

③NHK放送文化研究所の塩田雄大氏は「放送における外来語の過去と現在」について、放送における外来語使用の規範と実態の変遷について説明されました。規範では昔は個別に指定していたものから、理念の提示、そして言い換え案の提示という、柔軟な方式に変化してきていることが指摘されました。また、実態では外来語表記の変化とアルファベット語の増加の問題について取り上げられました。

フォーラムの後半に行われた参加者との質疑応答では、「外来語表記のコレと変化」や「タテ書き・ヨコ書きとアルファベット語増加の関係」などの質問が多く寄せられ、外来語と外国語に対する関心の高さがうかがわれました。講演者が質問への回答と解説を行い、「外来語表記の規範意識の変化」などをテーマにディスカッションをし、議論を深めることができました。



（伊藤 雅光）

新刊

1. 新「ことば」シリーズ20『文字と社会』

2007年3月／ぎょうせい／A5判横組み112ページ／税込500円

2. 『日本語科学』21

2007年4月／国書刊行会／B5判横組み160ページ／税込3,150円

連携大学院プログラムの発展を目指して～海外日本語教育セミナーの開催～

国立国語研究所では、政策研究大学院大学、国際交流基金日本語国際センターとの連携大学院プログラム「日本語教育指導者養成プログラム（修士課程）」及び「日本語文化研究プログラム（博士課程）」を平成13年10月より実施しています。海外で活躍する日本語教育指導者を対象とする、このプログラムは、これまでに16か国43名の修了生を輩出し、各地の日本語教育において指導者的な立場となる人材を育成してきました。

プログラムの一層の発展を目指し、昨年11月にウズベキスタン共和国において「日本語教育セミナー」を連携3機関共同で開催し、広報活動と修了生の追跡調査を実施しました。

■NIS諸国における日本語教育

NIS諸国では、カザフスタン共和国、キルギス共和国、ウズベキスタン共和国、ウクライナに7名の修了生が活躍しています。また、これらの国からの応募者も増加傾向にあります。

NIS諸国における日本語教育は、ソ連崩壊後の教育の自由化に伴い、ほとんどの機関が1991年以降に開始しています。カザフスタン共和国など、中央アジアの国々では、多くの国立大学に日本語専攻コースがあり、初等中等教育機関にも日本語教育の裾野が広がっています。

ウクライナにおける日本語教育は、19世紀にまで遡るようですが、本格的な日本語教育は、最高学府であるキエフ国立大学東洋学部の日本語学科開設

(1990年)に始まるといわれています。

首都キエフだけではなく、地方都市でもそれぞれ独自に日本語教育が行われています。



ウズベキスタンの日本語教室の掲示物

これらの国々では、中央アジア日本語弁論大会や全CIS (Commonwealth of Independent Statesの略、ロシア連邦及びNIS諸国)日本語弁論大会を共同で開催しており、国を越えた日本語教育のネットワークを通じて情報や意見の交換が行われています。

■日本語教育セミナーの開催

今回の日本語教育セミナーは、在ウズベキスタン日本国大使館、ウズベキスタン日本語教師会、ウズベキスタン日本人材開発センターの協力を得て、11月3日と4日の2日間にわたり、首都タシケントで開催されました。セミナー初日は連携大学院プログラム及び連携機関の紹介、2日目は連携大学院プログラムの教官2名による講演及び修了生2名の研究発表を行いました。

セミナーには、ウズベキスタン共和国及び周辺国の日本語教育関係者、連携大学院プログラムへの参加を検討している方々及び修了生の合計約50名の参加者がありました。連携機関の活動についての紹介、修了生の現状報告・研究発表など、多角的な情報を提供することによって、文章では伝えられない、プログラムの特色を理解していただくことができました。質疑では、現地の日本人教師からもプログラムへの期待の声がかげられました。



日本語教育セミナーの講演に耳を傾ける参加者

セミナーに前後して、周辺国における修了生の追跡調査を行いました。機関により事情は様々ですが、専門性の高い教師の不足、外国語教育学分野の研究の遅れ、等が共通の問題として指摘されました。連携大学院プログラムにとっても、これらの問題への対応が今後の課題のひとつとなるでしょう。

なお、本プログラムに関する詳しい情報は、国立国語研究所のホームページに掲載されています。

(福永 由佳)

さくら

春になるとサクラが咲き、日本各地のサクラの名所が賑わいます。サクラは、万葉の昔から詩歌に詠まれ、愛でる花として親しまれてきました。最近では女の子の名前で、さくらちゃんが好まれています。

「桜（櫻）」の字を漢和辞典で引いてみると、一番初めに「ゆすらうめ」（桜桃）という意味が出てきます。「さくら」は、漢字が持つ本来の意味（古典中国語での意味）とは違って日本だけで通用する意味（国訓）として、「ゆすらうめ」の次に出てきます。つまり、古い時代の中国語では、「桜（櫻）」の字はユスラウメを表す文字で、わたしたちが親しんでいるあのサクラを表す文字ではないのです。

もともとユスラウメを表す文字を、サクラの文字として使っていることに違和感を覚え、我慢できなくなった人もいます。江戸時代の花陰散人という人は、自分でサクラを表す文字を作っていました。「木」と「色」と「香」を組み合わせた文字です。

木色香

花陰散人は『櫛字説』という本を書いて、自作のサクラの文字を解説しています。内容を簡単に紹介します。

日本を代表する花であるサクラに専用の文字がないのはよろしくない。サクラは色も香りもとても良い花であるから、「木」と「色」と「香」を組み合わせ、サクラの文字にしよう。そもそも「櫛」の字は、私の夢に菅原道真公が現れて、サクラを表わす専用の文字がないことを悲しみ、サクラに「櫛」の字を使うようお告げを下されたのだ。

花陰散人は『櫛字説』を天満宮に奉納し、この文字が受け入れられることを祈りました。しかし、この文字は全く普及しませんでした。サクラは今でも「桜（櫻）」のままです。

江戸時代には、サクラは日本固有の花だと信じられていました。サクラは固有の花だから、それを表すための固有の文字が必要だと、花陰散人は考えたのでしょうか。でも、サクラは日本固有の花ではありませんでした。そして、花陰散人（江戸時代）や菅原道真（平安時代）が親しんだサクラ—ヤマザクラは、わたしたちが思い浮かべるあのサクラ—ソメイヨシノとは品種が違います。

「桜（櫻）」の字は、古い時代の中国語ではユスラウメを表し、昔の日本ならばヤマザクラ、現代ではソメイヨシノに代表されるサクラを表す文字です。植物の漢字は、時代と地域によって表すものが変わり、複雑です。
(高田 智和)

「ことばビデオ」シリーズ: 紹介ビデオが国研ホームページに 〈近日公開〉

国立国語研究所では、「ことばビデオ」シリーズ〈豊かな言語生活をめざして〉を平成13年度から17年度まで、年に1巻ずつ制作してきました。このビデオシリーズは、私たちが日ごろ何げなく使っている日本語の様々な側面に光を当て、映像と音声を通じて、言葉について考える機会を提供することを目的としています。

今までに作成したビデオは、全国の視聴覚ライブラリー等で視聴することができますし、製作会社から購入することも可能です。しかし、これまでは、実際にビデオを手に取り、再生をするまで、その内容、特に登場人物による生き生きとした言葉のやりとりを実感することができない、という課題がありました。その結果、「ことばビデオ」について、名前は知っていても、実際に見たことはないという方々が少なくなかったのです。

このたび、「ことばビデオ」の各巻について、その



一部を短く編集し、当研究所のホームページに掲載することになりました。クリック一つで、映像や音声の一部を見ることができます。これにより、「ことばビデオ」の中身に触れ、そのねらいを直接感じ取り、活用を広げていただけたらと願っています。

「ことば」フォーラム

「ことばと映像」をテーマに2回の「ことば」フォーラムを開催します。

国立国語研究所では、研究活動で得られた成果を学校教育・日本語教育・生涯学習などで広く活用していただくために、「ことばビデオ」という映像作品を制作してきました。その作品をめぐる、話しことばと映像の関係について考えます。

第32回「ことば」フォーラム 「映像作品から話しことばを考える」

「ことばビデオ」のもととなった調査データを紹介し、教育現場での活用の事例と可能性について考えるとともに、広く話しことばと映像の関係を考えます。

日 時：2007年6月30日（土）
午後1時30分（1時開場）～4時30分

場 所：国立国語研究所・講堂（定員180名）

講 師：ことばと映像

品田 雄吉（多摩美術大学名誉教授・映画評論家）

「ことばビデオ」の情報源

尾崎 喜光（国立国語研究所）

日本語教育で映像を使うと

小河原 義朗（北海道大学留学生センター）



- 多摩モノレール「高松駅」下車徒歩5分
- 立川バス 立川駅北口2番のりばより、立川バスで「自治大学校・国立国語研究所」下車徒歩2分

- ・入場無料・事前申し込み制。定員になり次第、締め切ります。
- ・手話通訳を御希望の方は、開催日の1週間前までに御連絡ください。

【申し込み方法】氏名・連絡先を下記まで御連絡ください。

国立国語研究所「ことば」フォーラム係 TEL：042-540-4300(代) FAX：042-540-4456

ホームページ (<http://www.kokken.go.jp/>) から申し込みができます。

予告

第33回「ことば」フォーラム 「映像作品から話しことばを考える—国語教育の現場で—」

日 時：2007年11月2日（金）午後

場 所：アクロス福岡 大会議室（定員200名）福岡市中央区天神1-1-1

講師予定：中神 智文（福岡県立朝倉高等学校）、杉戸 清樹（国立国語研究所）ほか

*詳細は、次号でお知らせします。

表紙のことば

表紙の写真は、沖電気社製の「漢字テレタイプライター」のキーボードです。全体では右の写真のような形をしています。

これは大型電子計算機に漢字情報を入力するための機械で、キーボードとペダルとの組み合わせで、一つのキーから4種類の文字が入力できました。入力データは紙テープにパンチングされて記録されます。

電子計算機が導入されたのは昭和41年で、まだ理系の研究室でも計算機を持つことが珍しかった時代ですが、国語研究所ではこうした機器を活用し、どのような語彙が多く使われているかを分析するなど、計量言語学と呼ばれる分野の確立に貢献しました。

