

(報道発表資料)

独立行政法人国立国語研究所

大規模書き言葉コーパスのオンライン試験公開 ～KOTONOHA「現代日本語書き言葉均衡コーパス」～

独立行政法人国立国語研究所（東京都立川市、所長：杉戸清樹）は、現在構築中の『現代日本語書き言葉均衡コーパス』のデータの一部、約 1000 万語分をインターネット上で試験公開します。（<http://www.kotonoha.gr.jp/demo/>）

国立国語研究所は、明治から現代にいたる日本語の電子化資料をコンピュータ上で公開しようとする KOTONOHA 計画を基幹的なプロジェクトとして推進しています。『現代日本語書き言葉均衡コーパス』（コーパスとはコンピュータ上に蓄積された大規模な言語資料のこと）はその一環として昨 2006 年度から構築を開始したものであり、2011 年の完成時には 1 億語を超える量の現代日本語の書き言葉データとして公開する予定です。

代表するに足る傑出した自然の	風景	地を指定する国立公園、国立公園	環境白書_平成6年版(各論)	環境庁企画調整局計画調査室	大蔵省印刷局	1994年6月
道庁県立自然公園 優れた自然の	風景	地であって、都道府県が条例の定	観光白書_昭和62年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1987年5月
。5 通信施設(1) 郵便1	風景	入通信日付印の押印 観光地、史	観光白書_昭和56年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1981年6月
国有林野のうち森林を主体とした	風景	が優れ、かつ、国土保全機能及び	環境白書_昭和60年版	環境庁企画調整局企画調整課	大蔵省印刷局	1985年5月
一般大衆に公開し、魚釣り、散策	風景	観賞等の場として提供するもので	観光白書_昭和63年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1988年5月
緑に恵まれ、各地に多くの優れた	風景	地や野生鳥獣を見ることができる	観光白書_昭和59年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1984年5月
名勝、郷土民俗等をデザインした	風景	入通信日付印を約3,800郵便	観光白書_昭和60年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1985年5月
ある。(i) 国立公園我が国の	風景	を代表するに足る傑出した自然の	観光白書_昭和59年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1984年5月
ける国際協力(実験農場での研修	風景)	農業白書_平成9年度(図説)	農林統計協会	(財)農林統計協会	1998年5月
公園 都道府県の風景を代表する	風景	地を指定する都道府県立自然公園	環境白書_平成6年版(各論)	環境庁企画調整局計画調査室	大蔵省印刷局	1994年6月
(ii) 自然公園は優れた自然の	風景	地であり、国民にとって自然との	観光白書_昭和62年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1987年5月
光・文化の振興に寄与する。iii	風景	入通信日付印の押印 観光地、史	観光白書_平成3年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1991年5月
れるフリーマーケットは日常的な	風景	となりつつあります。最近「フ	循環型社会白書_平成14年版	環境省大臣官房廃棄物・リサイクル対策部循環型社会推進室	ぎょうせい	2002年5月
ドライブ)16.3%、「自然・	風景	鑑賞)11.5%の順となってい	観光白書_平成5年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1993年6月
管理 自然公園は、優れた自然の	風景	地を保護することを目的の一つと	観光白書_昭和62年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1987年5月
名勝、郷土民俗等をデザインした	風景	入通信日付印を郵便局に備え付け	観光白書_平成3年版	総理府内閣総理大臣官房審議室	大蔵省印刷局	1991年5月

試験公開の目的

『現代日本語書き言葉均衡コーパス』をだれもが利用できるコーパスとするためには、全サンプルに著作権処理を実施する必要があります、その総数はおよそ 3 万件に及ぶものと想定されます。しかし昨今では、著作権保護、個人情報保護意識の高まりを反映して、著作権者との連絡にかかるコストが著しく増大する傾向にあり、著作権処理の成否がプロジェクト全体に大きく影響する状況となっています。

今回オンライン試験公開を開始する目的のひとつは、著作権者の方々にこのデモサイトを利用

していただき、ご提供いただくサンプルが実際にどのような形で利用されるかについて理解を深めていただくことにあります。また、不特定多数の方々に対して KOTONOHA 計画で開発中のコーパスに関する情報を提供することも、もうひとつの大切な目的です。

今回公開するデータ

現時点で検索することのできるデータは、各省庁刊行の白書のデータ約 500 万語分と、ヤフー株式会社提供の「Yahoo!知恵袋」のデータ約 500 万語分の合計約 1000 万語です。それぞれについてももう少し詳しく説明します。

白書データの母集団は 2001 年から 2005 年の間に発行された白書と過去 30 年間に継続して発行されつづけた白書の全体です。この母集団から無作為に約 500 万語分を抽出しています。サンプルは白書の対象にしたがって 9 個のカテゴリー（「安全」「科学技術」「外交」「環境」「教育」「経済」「国土交通」「農林水産」「福祉」）に分類されており、カテゴリーを限定した検索も可能です。同様に白書の刊行年を限定した検索も可能です。白書は政府の刊行物ですが、やはり著作権は存在しますので、「著作権フリー」を宣言しているごく一部の白書を除くサンプルについては関係省庁から書面ないし口頭で利用許諾をいただいています。

一方「Yahoo!知恵袋」は、参加者同士で知識を教えあうことを目的とした、Q&A 形式のナレッジコミュニティサービスです。この種のデータについては、従来から言語研究上の重要性が指摘されてきていましたが、今回、ヤフー株式会社のご厚意により、『現代日本語書き言葉均衡コーパス』のデータの一部として利用できるようになりました。現在公開されている「Yahoo!知恵袋」データの総量はそれだけで 1 億語を超える膨大なものですが、今回はそこから 500 万語分の質問と回答を無作為に選択して公開対象としました。

白書のデータと「Yahoo!知恵袋」のデータとはいろいろな面で対照的です。白書が硬い書き言葉のひとつの典型であるのに対して、「Yahoo!知恵袋」はかなりくだけた話し言葉的な特徴を示します。表記や文法についても、白書の文書がよく統制されているのに対して知恵袋は自由奔放といえそうです。

今後の展開

本試験公開サイトには、今後とも著作権処理が完了したデータを追加してゆく予定です。今後追加を予定しているデータには以下のものがあります。

- 国会会議録 (500 万語程度)
- 新聞記事 (100 万語程度)
- 文芸作品 (500 万語程度)

これらが追加されることによって、本試験公開サイトは次第に均衡コーパス（以下の「背景説明」参照）として整備されてゆくこととなります。

背景説明：『現代日本語書き言葉均衡コーパス』

現在、世界的潮流として、コーパスを利用した言語研究が活発におこなわれています。このような研究には、対象とする言語の様々なジャンルのデータを偏りなく収集した「均衡コーパス」

を整備する必要がありますが日本語に関するコーパスの整備は、英語圏はもとより中国語や韓国語に比べても著しく遅れているのが現状です。『現代日本語書き言葉均衡コーパス』はこのような問題を解消して我が国の言語研究のインフラを整備するために企画されたものです。またそれと同時に、常用漢字の見直しの問題、公用文の表記法の問題、教育用基本語彙の選択の問題など、日本語の言語政策にかかわる諸々の問題を検討するためにも活用される予定です。

出版(生産実態)サブコーパス 2001～2005年に出版された書籍、雑誌、新聞 3500万語	図書館(流通実態)サブコーパス 東京都の13自治体以上の図書館に収蔵されている書籍 対象期間:1986-2005 3000万語
特定目的(非母集団)サブコーパス ウェブ上の文書、白書、教科書、国会会議録、ベストセラー等 対象期間はさまざま、最長30年。3500万語	

図に示したように、『現代日本語書き言葉均衡コーパス』は三つのサブコーパスから構成されています。左上の出版サブコーパスは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、そこから合計で3500万語分のテキストを無作為に抽出するものです。

このサブコーパスは「どれだけ書き言葉が生産されたか」という観点から構築したコーパスです。書籍の販売部数には大きな差が生じますが、このサブコーパスでは

100万部を超えるベストセラーも数百部しか売れなかった本も同じ重みで扱われている点に特徴があります。

これに対して、実際にある程度広い範囲に出回ったことが判明している書籍だけを対象としているのが、図右上の図書館サブコーパスです。このサブコーパスでは、東京都下の公立図書館に所蔵されている蔵書のうち、少なくとも13館以上に収蔵されている図書(34万冊ほど)を母集団として、そこから無作為にサンプルを選んでいきます。図書館サブコーパスのもうひとつの特徴は、対象期間が20年間におよんでいることです。

最後に、図下部の特定目的サブコーパスは、国立国語研究所が種々の研究活動で必要とするデータのうち、上述のふたつのサブコーパスでは十分にデータが集まらないものを格納するものです。例えば公共性の高い文章における外来語や難解語の言い換え提案を行うための資料として白書の本文をデータとして格納しています。また、現代の日本語を考察するためには欠かすことのできないインターネット上の文書のデータもこのサブコーパスに格納しています。

謝 辞

出版サブコーパスおよび図書館サブコーパスのサンプリングには国立国会図書館、日本図書館協会、立川市立中央図書館、都立図書館、一橋大学図書館にご協力をいただいています。知恵袋データはヤフー株式会社からご提供いただきました。文芸作品の著作権処理には、日本文藝家協会、日本推理作家協会、日本児童文学者協会、日本児童文芸家協会、日本ペンクラブにご協力いただいています。新聞データについては、読売新聞、毎日新聞、産経新聞、共同通信社にご協力いただいています。

■本件に関するお問い合わせ先

独立行政法人国立国語研究所 研究開発部門 言語資源グループ
グループ長 前川喜久雄
電話：042-540-4515
メール：kikuo@kokken.go.jp