

Guide to using the KOTONOHA 少納言 (Syoonagon) corpus

- KOTONOHA (少納言/Syoonagon) corpus — demo version of the Balanced Corpus of Contemporary Written Japanese (BCCWJ) [<https://shonagon.ninjal.ac.jp/>]
- English-language overview [<https://clrd.ninjal.ac.jp/en/tool.html>]
- Description of corpus project [<https://link.springer.com/article/10.1007/s10579-013-9261-0>]
Maekawa, Kikuo et al. (2014). Balanced corpus of contemporary written Japanese. *Language Resources and Evaluation* 48(2): 345-371. (Open access journal article.)

The KOTONOHA corpus (Balanced Corpus of Contemporary Written Japanese; BCCWJ) was created by The National Institute for Japanese Language and Linguistics (NINJAL; 国立国語研究所). The Syoonagon interface is publicly available and allows full-text string-match searches. (Access to the morphologically parsed version of the corpus requires a license for the Tyuunagon interface.)

1. Accessing the corpus; Terms of Use

- (1) **Access the Syoonagon corpus** here: [<http://www.kotonoha.gr.jp/shonagon/>]
 - The main page has introductory and background information about the corpus; for an English-language introduction and overview, see [Maekawa et al. \(2014\)](#)
 - Scroll down to the red button at the bottom of the main page and click through to the Terms of Use:



利用条件を読んで少納言を使う

Read the Terms of Use and begin using Syoonagon

- (2) The **Terms of Use** are relatively straightforward; I'll give the key points here, but anyone who wants a fuller translation is welcome to contact me for more information.
 1. (著作権の帰属) — Copyright: NINJAL holds copyright to the corpus, and the sources of the materials sampled in the corpus hold copyright to those materials.
 2. (許諾の範囲等) — Limitations: The use of Syoonagon is limited to research and educational uses; duplication is prohibited; infringement of the rights of the copyright holders is prohibited.
 3. (研究成果の公表) — Publication: Research results and knowledge obtained by using the corpus in accordance with point (2) may be published; please acknowledge use of the Balanced Corpus of Contemporary Written Japanese (BCCWJ).
 4. (免責) — Disclaimers: NINJAL is not responsible for damages resulting from the use of the corpus. NINJAL may change access to Syoonagon at any time without prior notice.
 5. (利用条件の更新について) — These Terms of Use are subject to change.
 - If you agree with the Terms of Use, click OK to enter the corpus.

2. Searching the corpus

- (3) Overview — please read (4)–(6) for details.
- (a) Type or paste your search term in the text box. This must be an exact-string match (no options or alternatives are available — do a separate search for each distinct string you want to find).
 - (b) You can optionally specify the preceding or following context for your search term. This function allows you to specify alternatives at the single-character level! Use the format `^[ABC]`, with or without spaces, to mean “A or B or C”. (This may be most useful for finding conjugated forms of verbs.)
- (4) At the top of the search page is a text box. Type or paste your search term here. Note that the text you type will be matched exactly: 子供 <kodomo> ‘child’ and こども <kodomo> ‘child’ are two different searches.

<p>検索条件</p> <p>検索文字列: <input type="text"/></p> <p>こちらをクリックすると正規表現を使用して前後の文脈を指定できます。</p> <p><input type="button" value="検索"/> <input type="button" value="クリア"/></p>	<p>Search criteria</p> <p>Search string: _____</p> <p>(see below)</p> <p>Search Clear</p>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------

- (5) The third line in the view above says:
- [こちら](#)をクリックすると正規表現を使用して前後の文脈を指定できます。
- Click [here](#) to specify preceding and/or following contexts, using regular expressions.
- If you click, you see the following version of the search interface:

検索条件

検索文字列:

前文脈: [前後文脈の指定について](#)

後文脈:

- Here, lines 3 and 4 say:
 - 前文脈: Preceding context
 - 後文脈: Following context

(6) The link to the right side in the above view says:

前後文脈の指定について About specifying the preceding and following context

- If you click, you go to a screen with the following information. (Translations paraphrased; please see the actual web site for the search examples referred to in the chart below.)

前後文脈の指定について	Specifying preceding/following context
<p>前文脈および後文脈では、コーパス内において検索文字列の前後に現れる文字列を検索条件として指定できます。</p> <p>例えば、以下に示すように、検索文字列に「銀行」、前文脈に「日本」を指定して検索すると、「銀行」を含む文のうち、前に「日本」が現われている文のみを検索結果として表示します。</p>	<p>You can specify text strings that appear before and after the search string as search criteria.</p> <p>For example, as shown below*, you can specify the preceding context 「日本」 <nihon> 'Japan' for the search string 「銀行」 <ginkoo> 'bank'. This will return only those sentences containing 「銀行」 <ginkoo> preceded by 「日本」 <nihon>.</p> <p>(*see actual web site for search examples)</p>
<p>また、前文脈および後文脈は、正規表現を使用して指定できます。以下に示す検索条件を指定して検索すると、「取」を含む文のうち、直後に「ら」から「っ」のいずれかが現われている文のみを検索結果として表示します。</p>	<p>Preceding and following contexts can be specified using regular expressions. If you search using the search criteria shown below*, this will return sentences containing 「取」 <to[r]> 'take' in which the following character is one of [らりるれろっ] <ra ri ru re ro <i>small-tu</i>>.</p>

3. Restricting your search by genre or time period

(7) The first set of black bars below the search box allows you to restrict your search to particular genres of text. Select/deselect the relevant options. To include all subcategories of a particular genre, check the box in the left corner of the corresponding black bar.

メディア/ジャンル Media/genre

- Buttons say 全てのチェックを外す | 全てにチェックを入れる
Deselect all | Select all
- Click on the plus (+) symbol to expand the subcategories (as shown below)
- Screenshots of each category, with translations provided except where noted:

<input checked="" type="checkbox"/> <input type="checkbox"/> 書籍 (1971~2005)	
<input type="checkbox"/> 日本十進分類法 (NDC):	<input type="checkbox"/> 0 総記 <input type="checkbox"/> 1 哲学 <input type="checkbox"/> 2 歴史 <input type="checkbox"/> 3 社会科学 <input type="checkbox"/> 4 自然科学 <input type="checkbox"/> 5 技術・工学 <input type="checkbox"/> 6 産業 <input type="checkbox"/> 7 芸術・美術 <input type="checkbox"/> 8 言語 <input type="checkbox"/> 9 文学 <input type="checkbox"/> 分類なし
Books (1971-2005)	
Nippon Decimal System (NDC):	0 General 1 Philosophy 2 History 3 Sociology 4 Natural sciences 5 Tech/engineering 6 Industry 7 Arts 8 Language 9 Literature Unclassified

<input checked="" type="checkbox"/> <input type="checkbox"/> 雑誌 (2001~2005)	
<input type="checkbox"/> 総合:	<input type="checkbox"/> 総記/マスコミ <input type="checkbox"/> 一般 <input type="checkbox"/> 家庭/生活 <input type="checkbox"/> 児童 <input type="checkbox"/> 娯楽/芸能 <input type="checkbox"/> レジャー/趣味 <input type="checkbox"/> スポーツ
<input type="checkbox"/> 教育・学芸:	<input type="checkbox"/> 教育 <input type="checkbox"/> 学習/語学 <input type="checkbox"/> 文学/芸術 <input type="checkbox"/> 社会科学 <input type="checkbox"/> 自然科学 <input type="checkbox"/> 人文科学
<input type="checkbox"/> 政治・経済・商業:	<input type="checkbox"/> 政治/外交 <input type="checkbox"/> 経済/経営 <input type="checkbox"/> 金融/財政 <input type="checkbox"/> 商業/消費者 <input type="checkbox"/> 国勢/民力
<input type="checkbox"/> 産業:	<input type="checkbox"/> 農林水産 <input type="checkbox"/> 運輸/通信
<input type="checkbox"/> 工業:	<input type="checkbox"/> 機械 <input type="checkbox"/> 電気機/電子
<input type="checkbox"/> 厚生・医療:	<input type="checkbox"/> 厚生 <input type="checkbox"/> 医学

Magazines/journals (2001-2005)	
Comprehensive:	General/mass comm. General Home/lifestyle Juvenile Entertainment/arts Leisure/hobby Sports
Education/Liberal arts:	Education Language Literature/arts Social sciences Natural sciences Human sciences
Politics/Economics/Business:	Politics/diplomacy Economics/management Finance Business/consumers National resources
Industry:	Agriculture/forestry/fisheries Transportation/communication
Manufacturing:	Machines Electronics
Public welfare/Medicine:	Public welfare Medicine

<input checked="" type="checkbox"/> <input type="checkbox"/> 新聞 (2001~2005)	
<input type="checkbox"/> 全国紙:	<input type="checkbox"/> 朝日新聞 <input type="checkbox"/> 読売新聞 <input type="checkbox"/> 毎日新聞 <input type="checkbox"/> 産経新聞
<input type="checkbox"/> ブロック紙:	<input type="checkbox"/> 北海道新聞 <input type="checkbox"/> 中日新聞 <input type="checkbox"/> 西日本新聞
<input type="checkbox"/> 地方紙:	<input type="checkbox"/> 河北新報 <input type="checkbox"/> 新潟日報 <input type="checkbox"/> 京都新聞 <input type="checkbox"/> 神戸新聞 <input type="checkbox"/> 中国新聞 <input type="checkbox"/> 高知新聞 <input type="checkbox"/> 琉球新報

Newspapers (2001-2005)	
National:	Asahi Yomiuri Mainiti Nikkei
Blogs (regional online news):	Hokkaidoo Chuuniti Nisi-Nihon
Regional:	Kahoku Niigata Kyooto Koobe Tyuugoku Kooti Ryuukyuu

白書 (1976~2005)

White papers (1976-2005)

'White papers' are official (typically government) research and policy reports. There are many subcategories listed here; please let me know if you need the subcategories translated.

教科書 (2005~2007)

<input type="checkbox"/> 国語:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 数学:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 理科:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 社会:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 外国語:	<input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 技術家庭:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 芸術:	<input type="checkbox"/> 小 <input type="checkbox"/> 中 <input type="checkbox"/> 高
<input type="checkbox"/> 保健体育:	<input type="checkbox"/> 高
<input type="checkbox"/> 情報:	<input type="checkbox"/> 高
<input type="checkbox"/> 生活:	<input type="checkbox"/> 小

Textbooks (2005-2007)

Japanese language:	Elementary Middle school High school
Mathematics:	Elementary Middle school High school
Science:	Elementary Middle school High school
Society:	Elementary Middle school High school
Foreign language:	Middle school High school
Home economics:	Elementary Middle school High school
Fine arts:	Elementary Middle school High school
Health/physical education:	High school
Information:	High school
Daily life:	Elementary

広報紙 (2008)

[Local government] newsletters (2005-2007)

The subcategories here are regions in Japan; please let me know if you need translations.

<input checked="" type="checkbox"/> <input type="checkbox"/> Yahoo!知恵袋 (2005)	Yahoo! answers (2005) (Please let me know if you need sub-subcategories translated.)
<input type="checkbox"/> エンターテインメントと趣味:	Entertainment and hobbies
<input type="checkbox"/> インターネット、PCと家電:	Internet/PCs/Home electronics
<input type="checkbox"/> ビジネス、経済とお金:	Business, economics, and money
<input type="checkbox"/> 職業とキャリア:	Employment and careers
<input type="checkbox"/> ニュース、政治、国際情勢:	News, politics, international
<input type="checkbox"/> スポーツ、アウトドア、車:	Sports, outdoors, autos
<input type="checkbox"/> 暮らしと生活ガイド:	Lifestyle guides
<input type="checkbox"/> 健康、美容とファッション:	Health, beauty, and fashion
<input type="checkbox"/> 子育てと学校:	Parenting and education
<input type="checkbox"/> マナー、冠婚葬祭:	Etiquette, ceremonial occasions
<input type="checkbox"/> 教養と学問、サイエンス:	Culture, education, science
<input type="checkbox"/> 地域、旅行、お出かけ:	Geography, travel, excursions
<input type="checkbox"/> Yahoo! JAPAN:	Yahoo! JAPAN
<input type="checkbox"/> その他:	Other

<input checked="" type="checkbox"/> <input type="checkbox"/> Yahoo! ブログ (2008)	Yahoo! blogs (2008) (Please let me know if you need sub-subcategories translated.)
<input type="checkbox"/> ビジネスと経済:	Business and economics
<input type="checkbox"/> コンピュータとインターネット:	Computers and internet
<input type="checkbox"/> 生活と文化:	Lifestyle and culture
<input type="checkbox"/> エンターテインメント:	Entertainment
<input type="checkbox"/> 家庭と住まい:	Family and home
<input type="checkbox"/> 政治:	Politics
<input type="checkbox"/> 健康と医学:	Health and medicine
<input type="checkbox"/> 学校と教育:	School and education
<input type="checkbox"/> 科学:	Science
<input type="checkbox"/> 出会い:	Dating/marriage
<input type="checkbox"/> 地域:	Regional
<input type="checkbox"/> 特集:	Specialized — <i>single subcategory is Hobbies and sports</i>
<input type="checkbox"/> 芸術と人文:	Arts and culture
<input type="checkbox"/> Yahoo!サービス:	Yahoo! service
<input type="checkbox"/> 趣味とスポーツ:	Hobbies, sports — <i>subcategories are gambling, sports, leisure, hobbies, vehicles</i>

<input checked="" type="checkbox"/> <input type="checkbox"/> 韻文 (1980~2005)	Poetry (1980-2005) (These subcategories have no further sub-subcategories.)
<input type="checkbox"/> 短歌:	Tanka
<input type="checkbox"/> 俳句:	Haiku
<input type="checkbox"/> 詩:	Poems

法律 (1976~2005)

Laws (1976-2005)

There are many subcategories listed here; please let me know if you need the subcategories translated.

国会会議録 (1976~2005)

<input type="checkbox"/> 衆議院:	<input type="checkbox"/> 本会議	<input type="checkbox"/> 常任委員会	<input type="checkbox"/> 特別委員会	<input type="checkbox"/> その他
<input type="checkbox"/> 参議院:	<input type="checkbox"/> 本会議	<input type="checkbox"/> 常任委員会	<input type="checkbox"/> 特別委員会	<input type="checkbox"/> その他

Records of Diet [national legislature] proceedings (1976-2005)

Lower house:	Regular session	Standing committees	Special committees	Other
Upper house:	Regular session	Standing committees	Special committees	Other

- (8) The final black box at the bottom of the search page allows you to restrict your search by time period.
- You can select entire decades (年代) or individual years
 - To search all years, check the box in the black bar (全期間 = 'all time periods')

全期間

<input type="checkbox"/> 1970年代:	<input type="checkbox"/> 1971	<input type="checkbox"/> 1972	<input type="checkbox"/> 1973	<input type="checkbox"/> 1974	<input type="checkbox"/> 1975	<input type="checkbox"/> 1976	<input type="checkbox"/> 1977	<input type="checkbox"/> 1978	<input type="checkbox"/> 1979	
<input type="checkbox"/> 1980年代:	<input type="checkbox"/> 1980	<input type="checkbox"/> 1981	<input type="checkbox"/> 1982	<input type="checkbox"/> 1983	<input type="checkbox"/> 1984	<input type="checkbox"/> 1985	<input type="checkbox"/> 1986	<input type="checkbox"/> 1987	<input type="checkbox"/> 1988	<input type="checkbox"/> 1989
<input type="checkbox"/> 1990年代:	<input type="checkbox"/> 1990	<input type="checkbox"/> 1991	<input type="checkbox"/> 1992	<input type="checkbox"/> 1993	<input type="checkbox"/> 1994	<input type="checkbox"/> 1995	<input type="checkbox"/> 1996	<input type="checkbox"/> 1997	<input type="checkbox"/> 1998	<input type="checkbox"/> 1999
<input type="checkbox"/> 2000年代:	<input type="checkbox"/> 2000	<input type="checkbox"/> 2001	<input type="checkbox"/> 2002	<input type="checkbox"/> 2003	<input type="checkbox"/> 2004	<input type="checkbox"/> 2005	<input type="checkbox"/> 2006	<input type="checkbox"/> 2007	<input type="checkbox"/> 2008	

(self-explanatory!)

Last update: September 28, 2023