

Manyoshu and Hachidaishu in the Corpus of Historical Japanese ver. 2019.3

小木曾 智信 / Toshinobu OGISO

国立国語研究所 / NINJAL

2019/09/20 EAJRS in Sofia, Bulgaria

Topics

- NINJAL Corpora & CHJ
- Recent Updates to CHJ
- CHJ Manyoshu & ONCOJ etc.
- CHJ Hachidaishu
- (Preliminary Analysis)

Center for Corpus Development, NINJAL

コーパス開発センター
Center for Corpus Development, NINJAL

国立国語研究所コーパス開発センターでは、日本語の全貌を把握するための言語コーパス (language corpus) を構築しています。

English 国立国語研究所

コーパス ツール 申込方法 KOTONOHA計画 語彙調査データ 報告書 イベント

ご覧になりたいコーパス名をクリックしてください

現代日本語書き言葉均衡コーパス 日本語歴史コーパス

日本語話し言葉コーパス 国語研日本語ウェブコーパス

多言語母語の日本語学習者横断コーパス 名大会話コーパス

現日研・職場談話コーパス 日本語日常会話コーパス

近代語のコーパス コーパスアノテーション

FREE 中納言 利用申込 SUBSCRIPTION

CHARGED BCCWJ(有償版) 利用申込 SUBSCRIPTION

CHARGED CSJ(有償版) 利用申込 SUBSCRIPTION

最新情報 > 最新情報リスト

2019/05/09 お知らせ
2019年7月10日(水)~9月30日(月)の間、BCCWJ(DVD版)・CSJ(USB版)の利用申込受付を中断いたします。

2019/03/26 最新情報
『日本語歴史コーパス』(ver.2019.3)を公開しました。

講義・講習ビデオ

言語資源活用ワークショップ

UniDic - 形態素解析辞書 -

Web茶まめ - 形態素解析支援ツール -

分類語彙表 - データベース -

トライ!コーパス!WEB検索ツール

登録不要 少納言

登録制 中納言

登録制 (一部登録不要) 梵天



NINJAL Corpora

Balanced Corpus of
Contemporary Written Japanese

Corpus of
Spontaneous Japanese

International Corpus of
Japanese As a Second language
-In Japanese Only-

Gen-Nichi-Ken
Corpus of Workplace Conversation

Corpus of
Historical Japanese

NINJAL Web Japanese Corpus

Nagoya University
Conversation Corpus

Corpus of
Everyday Japanese Conversation

NINJAL Corpora

- Spoken 話し言葉 → CSJ
- Written 書き言葉 → BCCWJ
- **History 歴史 → CHJ**
- Web (100 billion words) → NWJC
- Dialects 方言 → COJADS
- Everyday Conversation 日常会話 → CEJC

Corpus as a research infrastructure

BCCWJ (Balanced Corpus of Contemporary
Written Japanese)

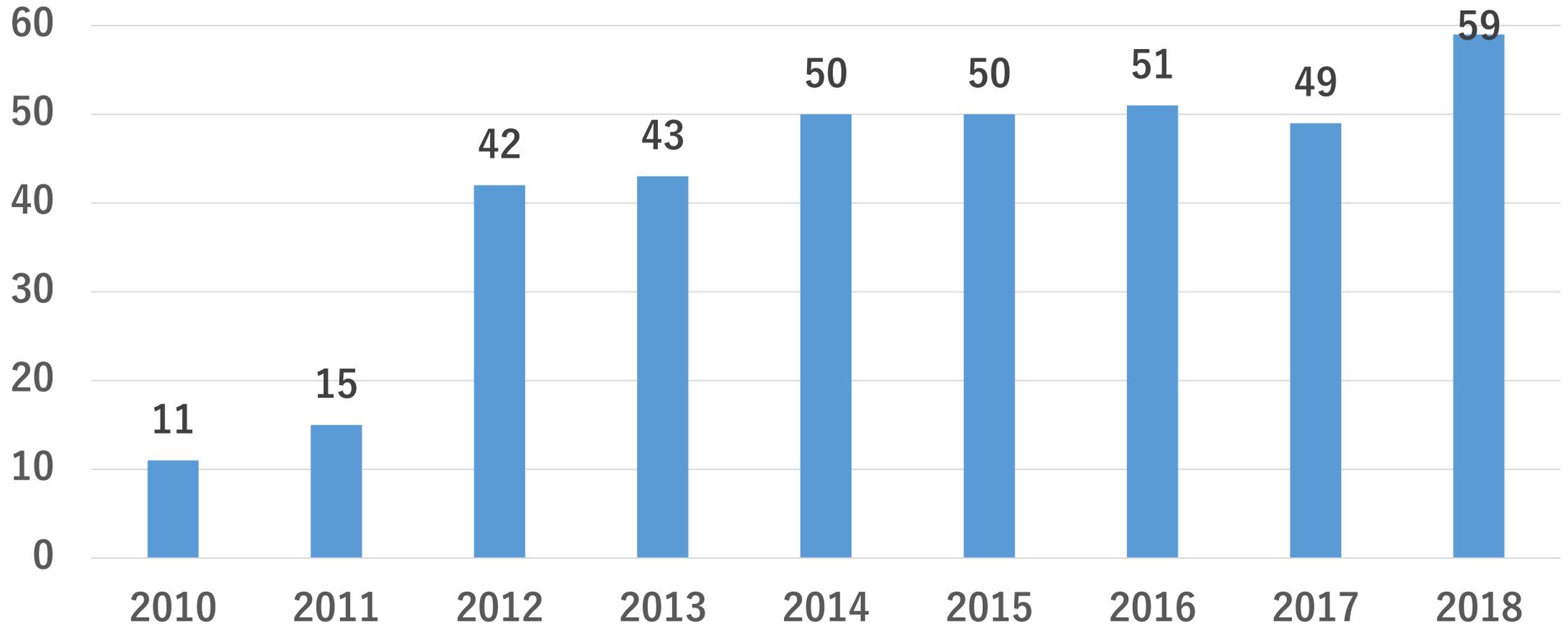
- Registered users : 20,000
- Queries : 500,000/year
- Papers : 70/year

Corpus as a research infrastructure

CHJ (Corpus of Historical Japanese)

- Registered users : 10,000
- Queries : 260,000/year
- Papers : 50/year
(including proceedings)

Number of papers using CHJ



NINJAL Diachronic Corpus Project

The Construction of Diachronic
Corpora and New Developments
in Research on the History of
Japanese (2016-2022)

「通時コーパスの構築と日本語史研究の新展開」プロジェクト

1

1000年をこえる日本語の歴史を収録する通時コーパスを構築

- ・ 上代（奈良時代以前）から明治・大正時代まで

2

通時コーパスを活用した新しい日本語史研究を展開

- ・ マクロな視点
- ・ 新しい手法の導入

3

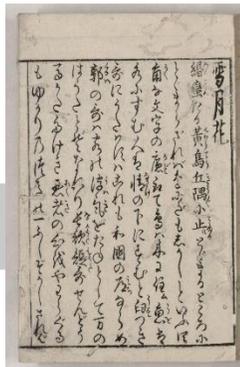
日本語史研究者以外の人々にも

- ・ 歴史的研究への参入障壁をコーパスで排除、海外の研究者や関連分野の研究者にオープンに
- ・ ポータルサイトを構築し言葉に関心を持つ全ての人々に情報提供

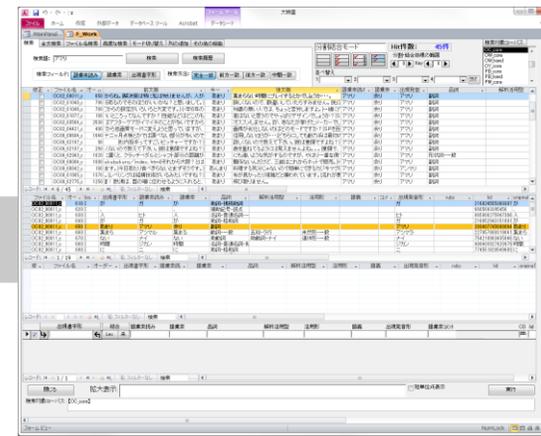
➤ 総合的日本語研究の骨格となる日本語史の研究基盤の構築と活用

Corpus of Historical Japanese

- **Diachronic corpus** from 8th to 20th century
 - **Full text** corpus of historical Japanese materials
 - Lemmatized head words, part-of-speech and morphological information are **annotated**
 - **Linked** to related resources on the web



```
1 <|text|><|id|> 成朝本辞書 039 大名 徳久入 <|series|> 成朝本辞書・大名辞書之類V上R30  
2 <|title|> 徳久入 <|year|> 1642 <|year_w|> 寛永19 >  
3 <|front|>  
4 <|titleBlock|><|block_type|>title <|s|><|pb.n|>177 </|b|></|info.originalPage|>256 </|>  
5 徳久入 </|s|>  
6 </|block|></|titleBlock|></|front|>  
7 <|body|>  
8 <|p|>  
9 <|stage|><|s|></|b|>「名乗<|char_script|>カタカナ</|>は<|char|>ひ<|rend|>傍線</|>入間<|hi|>  
10 </|stage|><|speech|><|speaker_value|>大名 </|>  
11 <|p|>ら <|ruby_resp|>校注者 <|ruby_text|>この <|ruby|>計 <|ruby_resp|>校注者 <|ruby_text|>計  
12 <|p|>ちゅう <|ruby|>某 <|ruby_resp|>校注者 <|ruby_text|>計 <|ruby|>きも <|ruby|>計 <|ruby|>計  
13 <|p|>e=eratum <|resp|>annotator <|cor|>馳 <|cor|>走 <|ruby_resp|>校注者 <|ruby_text|>計  
14 <|p|>計 <|ruby|>計  
15 </|p|>  
16 </|speech|><|speaker_value|>太郎冠者 </|>  
17 </|text|>
```



Recent Updates to CHJ

Sub-corpora of CHJ

奈良時代	<input checked="" type="checkbox"/> 万葉集 <input type="checkbox"/> 宣命
平安時代	<input checked="" type="checkbox"/> 仮名文学
鎌倉時代	<input checked="" type="checkbox"/> 説話・随筆 <input checked="" type="checkbox"/> 日記・紀行 <input type="checkbox"/> 軍記
室町時代	<input checked="" type="checkbox"/> 狂言 <input checked="" type="checkbox"/> キリシタン資料
江戸時代	<input checked="" type="checkbox"/> 洒落本 <input checked="" type="checkbox"/> 人情本 <input type="checkbox"/> 近松
明治・大正	<input checked="" type="checkbox"/> 雑誌 <input checked="" type="checkbox"/> 教科書 <input type="checkbox"/> 文学作品 <input type="checkbox"/> 新聞 <input checked="" type="checkbox"/> 明治初期口語資料

Update of CHJ / 2018

March 2018

- Muromachi Series II: *Kirisitan Shiryo*
(Amakusa Edition of Feiqe, Esopo) thanks to BL

October 2018

- Meiji/Taishō Series II: Textbooks
(国定読本)

Update of CHJ / March 2019 (1)

Meiji/Taishō Series

- Magazines
+ Toyo Gakugei Zasshi 1881, 1882
- Early Meiji Spoken Language Materials
(明治初期口語資料)
+ 10 works

Update of CHJ / March 2019 (2)

- Edo Series

 - + 洒落本 Sharebon (30 books)

 - + 人情本 Ninjobon (7 works=93 books)

- Wakashu Series

 - + **八代集 Hachidaishu**

Web Interface Chunagon 中納言

- Online concordancer for NINJAL corpora
- Enables query of morphological information



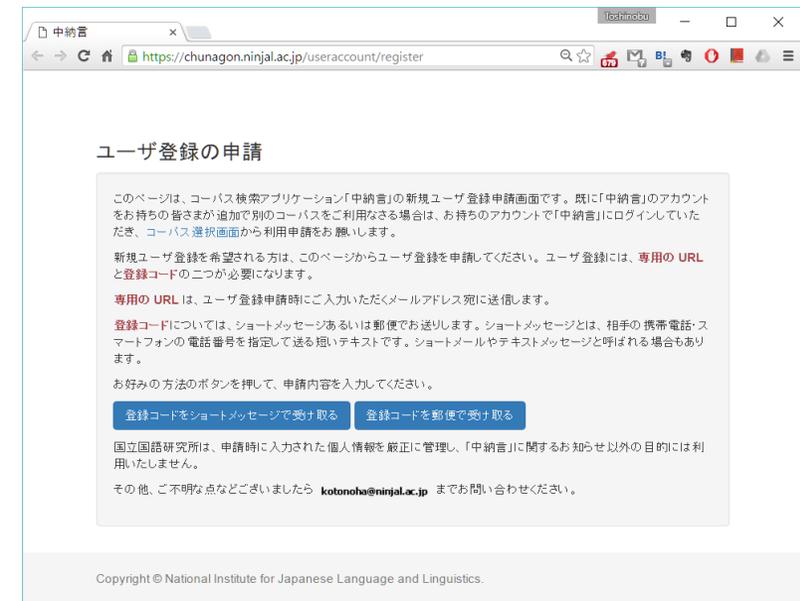
<https://chunagon.ninjal.ac.jp>



Registration

- Free of charge, Registration needed
- Online registration using email and SMS (text)

<https://chunagon.ninjal.ac.jp/useraccount/register>



How to use Chunagon?

Resource providers workshop

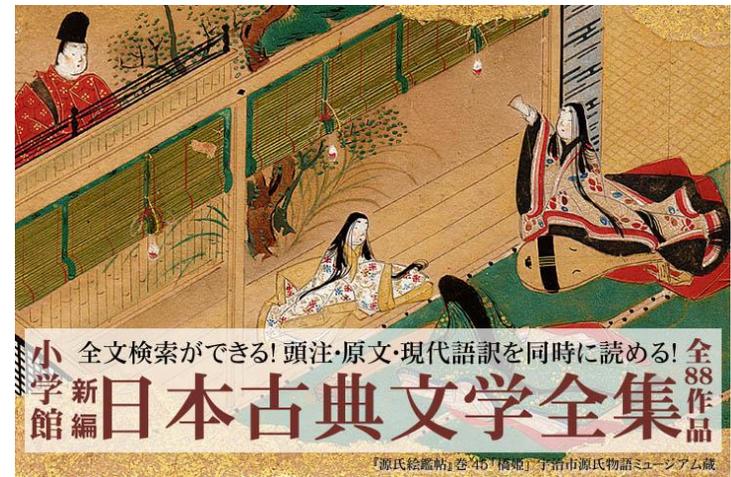
TODAY 15:00-18:00

NINJAL 国立国語研究所

CHJ Manyoshu 万葉集

CHJ奈良時代編 | 万葉集

- 「万葉集」 全20巻 4,516首
- 約98,000語
- 小学館「新編日本古典文学全集」
- 万葉仮名の原文本文と対応付け



万葉仮名

籠毛与 美籠母乳 布久
思毛与 美夫君志持 此
岳尔 菜採須兒 家告閑
名告紗根 虚見津 山跡
乃国者 押奈戸手 吾許
曾居 師吉名倍手 吾己
曾座 我許背齒 告目
家呼毛名雄母

籠もよ み籠持ち ぶくし
もよ みぶくし持ち この
岡に 菜摘ます兒 家告ら
せ 名告らさね そらみつ
大和の国は おしなべて
我こそ居れ しきなべて
我こそいませ 我こそば
告らめ 家をも 名をも

CHJ万葉集の特長

検索結果に万葉仮名の原文を併記

10-万葉 0759_00008	37020	葛花 なでしこ が 花 をみなへし また 藤袴 朝顔 が 花 #	秋	の 日 の 穂田 を 雁がね 暗けく に 夜 の
		花瞿麦之花姫部志又藤袴朝貞之花 #	秋	日乃穂田乎鴈之鳴闇尔夜之穂杼呂
10-万葉 0759_00004	22630	しきたへ の 衣手 交へ て 自妻 と頼め る 今夜	秋	の 夜 の 百夜 の 長 さ あり こせ ぬ
		而敷細乃衣手易而自妻跡憑有今夜	秋	夜之百夜乃長有与宿鴨
10-万葉 0759_00019	54900	万代 に 記し 継が む そ やすみしし 我 が 大君	秋	の 花 し が 色々 に 見し たまひ 明らめ たまひ
		仁万世尔記続牟曾八隅知之吾大皇	秋	花之我色と尔見賜明米多麻比酒見
10-万葉 0759_00010	41020	が 衣 摺ら む に にほひ こそ 島 の 榛原	秋	立た ず とも # 風 に 散る 花橘 を 袖 に
		思ふ之木摺排々保比与嶋之榛原	秋	不立之 # 風散花橘則袖受吾為君御跡

CHJ万葉集の特長

- 万葉仮名の原文（漢字）での検索が可能
- 用例を参照するpermalink
https://chunagon.ninjal.ac.jp/chj/permalink?unit=short&position=10-万葉0759_00005,13660

× 令和

- コーパスのアノテーション対象は歌のみ
- 「令和」の出てくる左注は対象に含まれず
- https://chunagon.ninjal.ac.jp/chj/permalink?unit=short&position=10-万葉0759_00005,13660

CHJ Manyoshu & ONCOJ

ONCOJ

Oxford-NINJAL Corpus of Old Japanese

Oxford-NINJAL上代日本語コーパス



Front page

オックスフォード・NINJAL上代語コーパス (ONCOJ)

Front page

ONCOJ

Prof. Bjarke Frellesvig (Univ. of Oxford) の
プロジェクトが構築した, 統語情報付きの
上代日本語コーパス (OCOJ)



オックスフォードと国語研との協力で整備
拡充して国語研から公開中

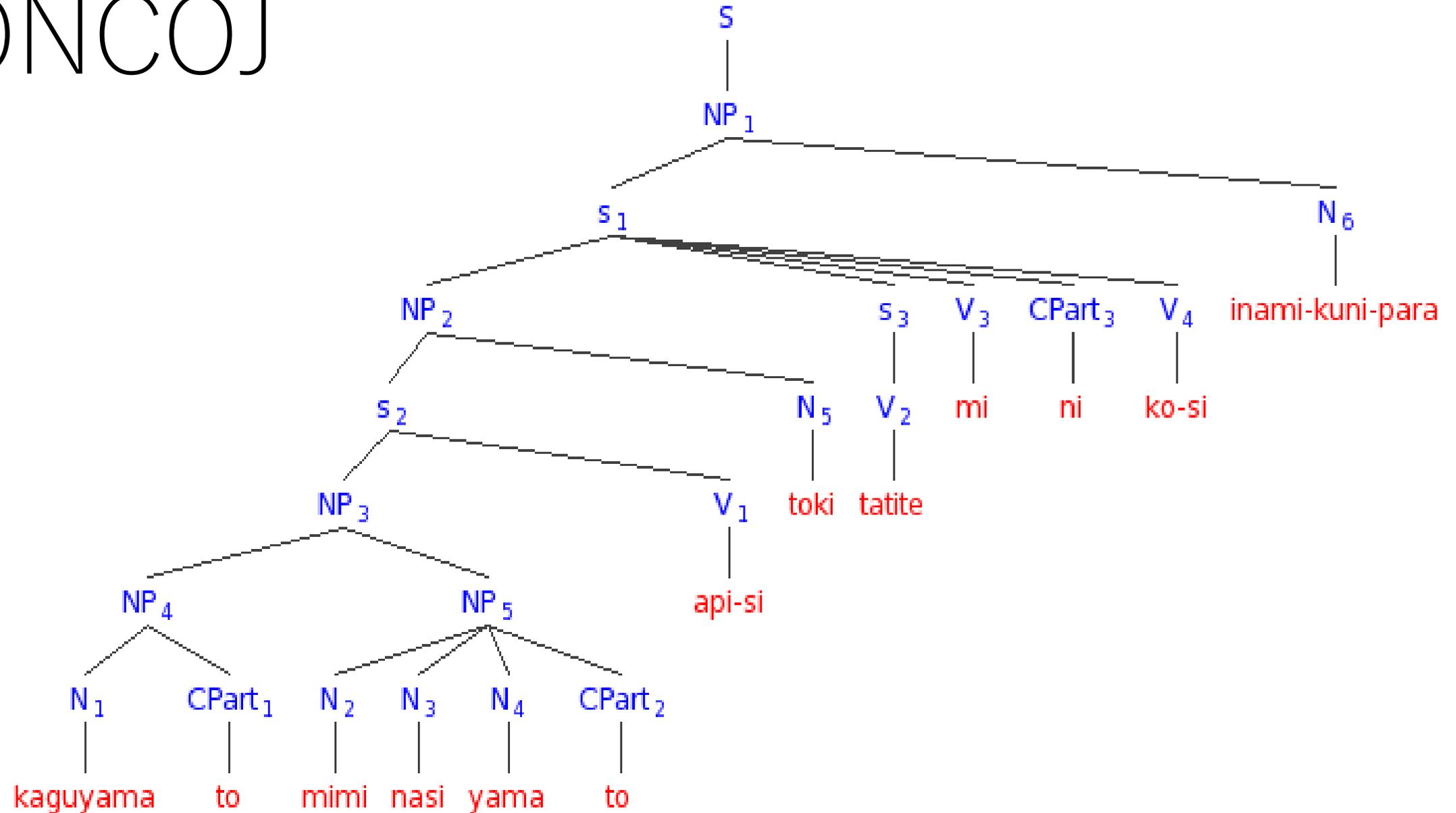
<https://oncoj.ninjal.ac.jp>

ONCOJ

• **MYS.1.14**

- 高山与
- 耳梨山与
- 相之時
- 立見尔来之
- 伊奈美國波良
- kaguyama to
- mimi nasi yama to
- apisi toki
- tatite , mi ni kosi
- inamikunipara .

ONCOJ



ONCOJとCHJ

- ONCOJ

- 上代日本語の共時的なコーパス
- 詳細な文法情報付き

- CHJ

- 上代から近現代までの通時的なコーパス
- 単語情報＋外部リンク

CHJ ↔ ONCOJ

歌番号による相互リンク

ONCOJ → CHJ

CHJ → ONCOJ / coming soon

万葉集の写本

- 万葉集の原本は失われ、平安時代中期以降の写本が残る（最古の写本、桂本は巻第四のみ）
- 全巻揃ったものは鎌倉時代後期（西本願寺本）

万葉集校本DBへのリンク

万葉集校本データベース

(万葉集校本データベース作成委員会)

- https://www.manyou.gr.jp/SMAN_1/

- CHJから歌番号でリンク
coming soon

万葉集校本DBへのリンク



● 歌首表示画面

● 各古写本表示

任意の歌句表示 で囲まれている各部分を選択すると、本文異同を表示。

● [本文表示](#)

各注釈本に採用されている本文一覧が表示されます。

● [訓み表示](#)

各注釈本に採用されている訓み一覧が表示されます。

● 注釈リスト

(※現在データ不備のためご使用いただけません)

注釈本リストが表示されます。

この中から見たい注釈データを選択できます。

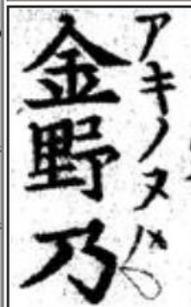
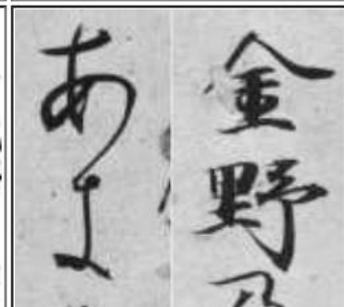
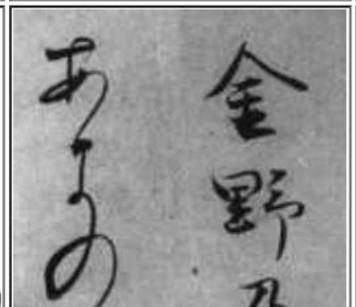
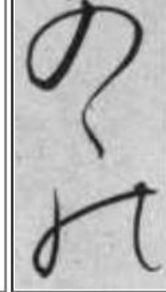
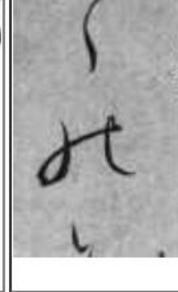
● [口訳表示](#)

各注釈本に採用されている口訳一覧が表示されます。



万葉集校本DBへのリンク

- 注釈本の本文、古写本の画像表示

		寛永版本	元暦校本	広瀬本	類聚古集	紀州	
C03	万葉集拾穂抄	金野乃 美草刈菫 屋杼礼里之 兔道乃宮子能 借五百磯所念					
C04	万葉代匠記 (初稿本) (精撰本)	金野乃 美草刈菫 屋杼礼里之 兔道乃宮子能 借五百●所念					
C05	万葉集僻案抄・童蒙抄・割記	金野乃 美草刈菫 屋杼禮里之 兔道乃宮子能 借五百磯所念					
C06	万葉考	金野乃 美草刈菫 屋杼禮里之 兔道乃宮子能 借五百磯所念					

CHJ Hachidaishu 八代集

コーパス化の対象

- 底本

国文学研究資料館所蔵 正保版本
『二十一代集』 テキストデータ

<http://base1.nijl.ac.jp/~selectionfulltext/21textpagelist.html>

- 二十一代集より八代集のみ

古今、後撰、拾遺、後拾遺、金葉、詞花、
千載、新古今

表 「和歌集編」作品別語数(短単位)

和歌集名	成立年	語数
古今和歌集	905	29,943
後撰和歌集	955	38,881
拾遺和歌集	1005	34,164
後拾和歌集	1087	40,250
金葉和歌集	1128	20,543
詞花和歌集	1151	12,376
千載和歌集	1188	37,095
新古今和歌集	1205	49,160
合計		262,412

本文校訂

- 濁点（゛）の付与
- 踊り字（ゝ / ㇵ）の処理
- 送り仮名（久恋→久しき恋）修正
- 誤り（幡磨→播磨）訂正

NIJL原文画像へのリンク

- NIJL正保版本『二十一代集』影印画像
- CODHのIIF「日本古典籍ビューア」經由リンク

<http://codh.rois.ac.jp/pmjt/book/200007092/>

NIJL原文画像へのリンク

季	後文脈	語彙素読み	語彙素	品詞	原文文字列	本文種別	歌番号	作品名	巻名等	底本	ページ番号	底本リンク	参考リンク
秋	後文脈	アキ	秋	名詞-普通名詞-副詞可能	秋			古今和歌集	仮名序	古典選集本文DB 正保版本「二十一代集」	7	Nijl	Maruzen
秋	の月 の夜 こと に さ ぶら ふ 人 々 を め し	アキ	秋	名詞-普通名詞-副詞可能	秋			古今和歌集	仮名序	古典選集本文DB 正保版本「二十一代集」	8	Nijl	Maruzen
秋	の月 の夜 こと に さ ぶら ふ 人 々 を め	アキ	秋	名詞-普通名詞-副詞可能	秋			古今和歌集	仮名序	古典選集本文DB 正保版本「二十一代集」	9	Nijl	Maruzen
秋	の月 を見る に あ か つ き の 雲 に あ へ	アキ	秋	名詞-普通名詞-副詞可能	秋			古今和歌集	仮名序	古典選集本文DB 正保版本「二十一代集」	10	Nijl	Maruzen



歌人情報の整備とリンク

- 『勅撰集 付新葉集作者索引』和泉書院索引叢書(1986)により歌人名を標準化
- NDL Authoritiesの当該歌人にリンク

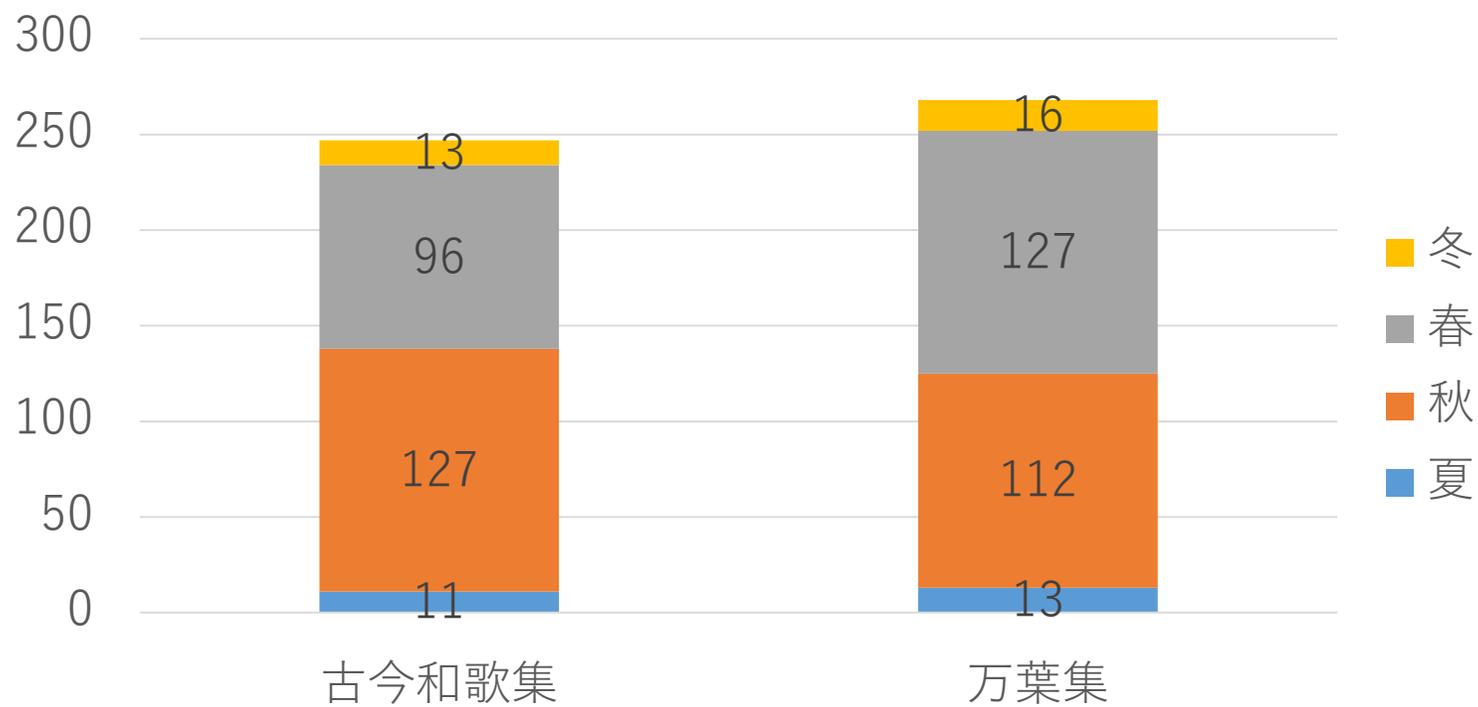
既存の資料と比較した CHJ八代集の長所

国歌大観、岩波「八代集」CD-ROM etc.

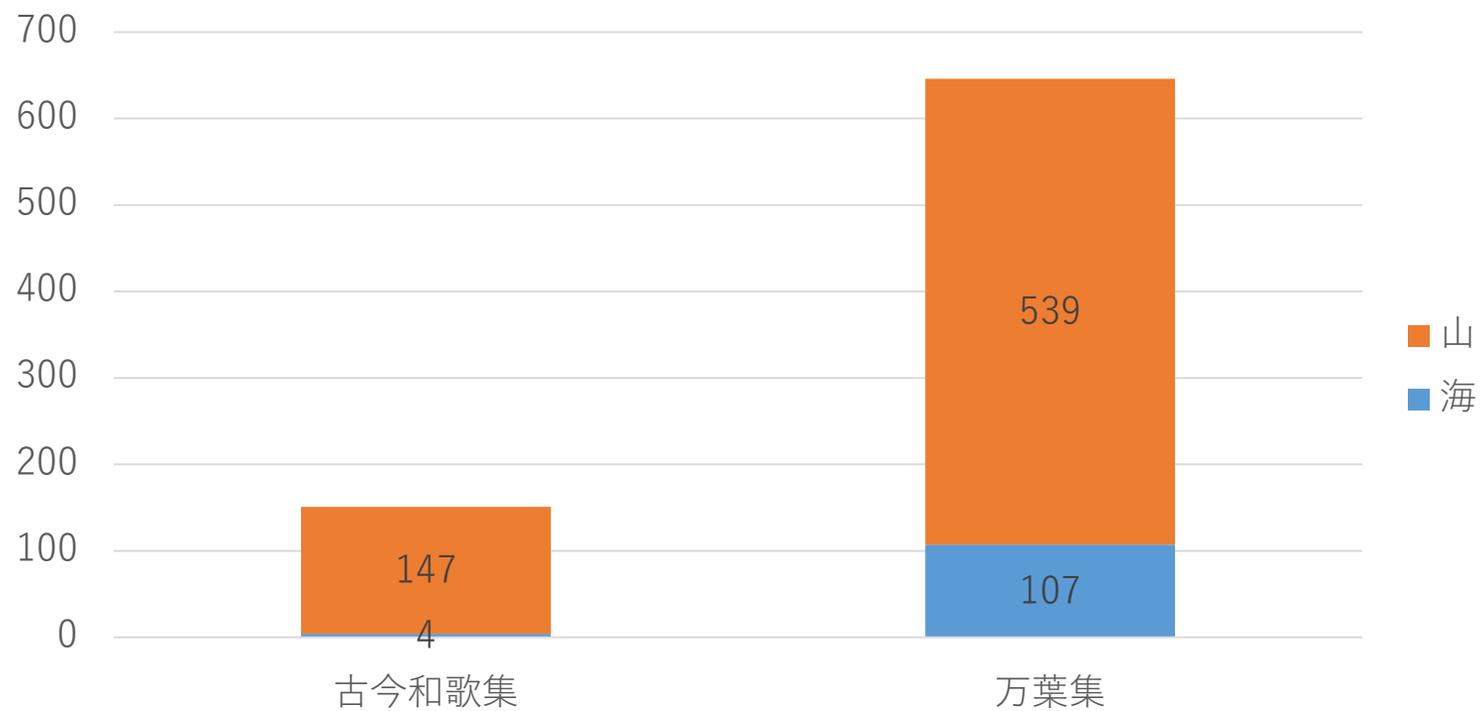
- 同一基準による単語情報付き
 - 句や文字列でなく見出し語で検索可能
 - 散文を含む他の作品と通時的に比較可能

Preliminary Analysis

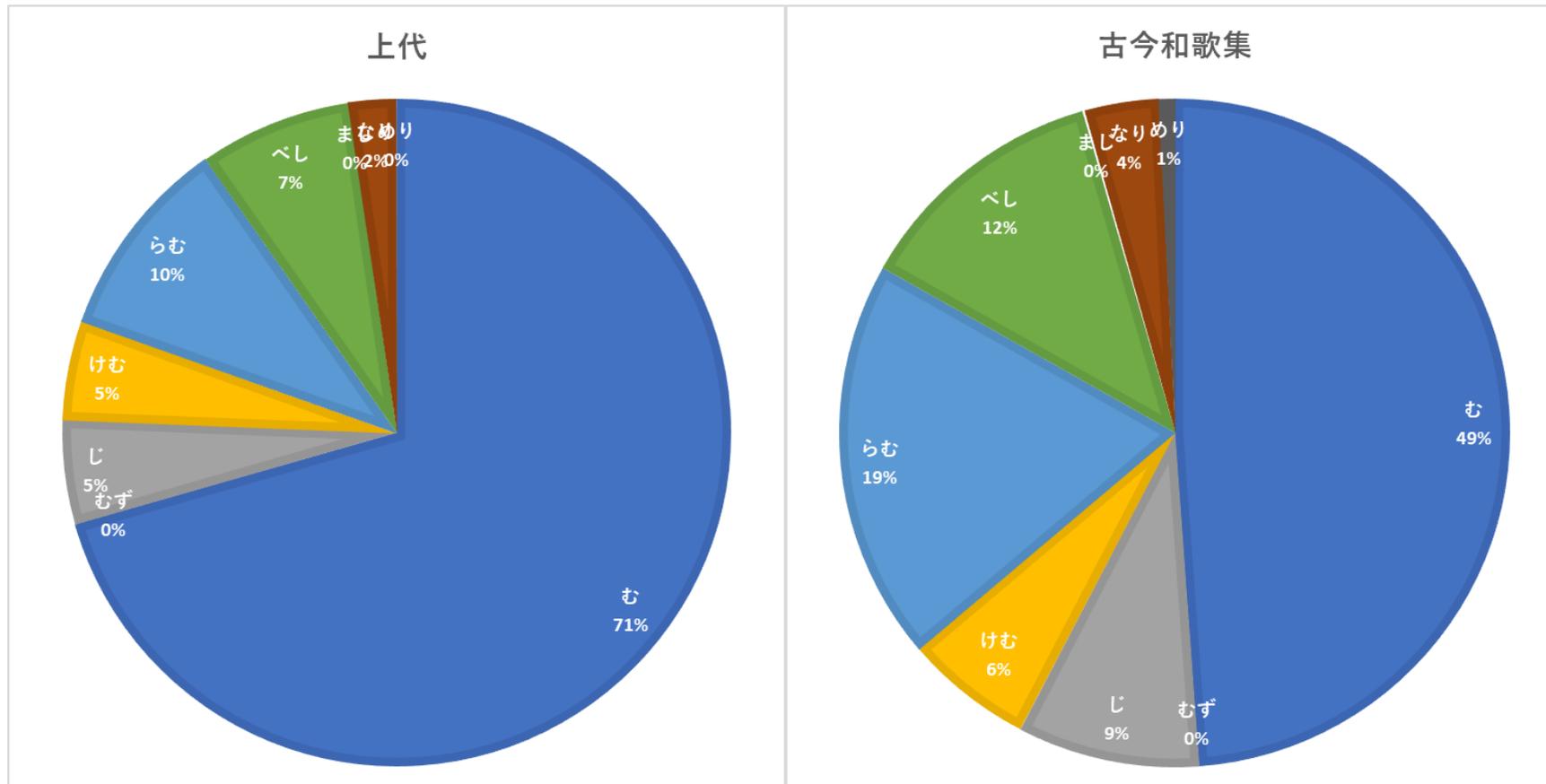
春夏秋冬



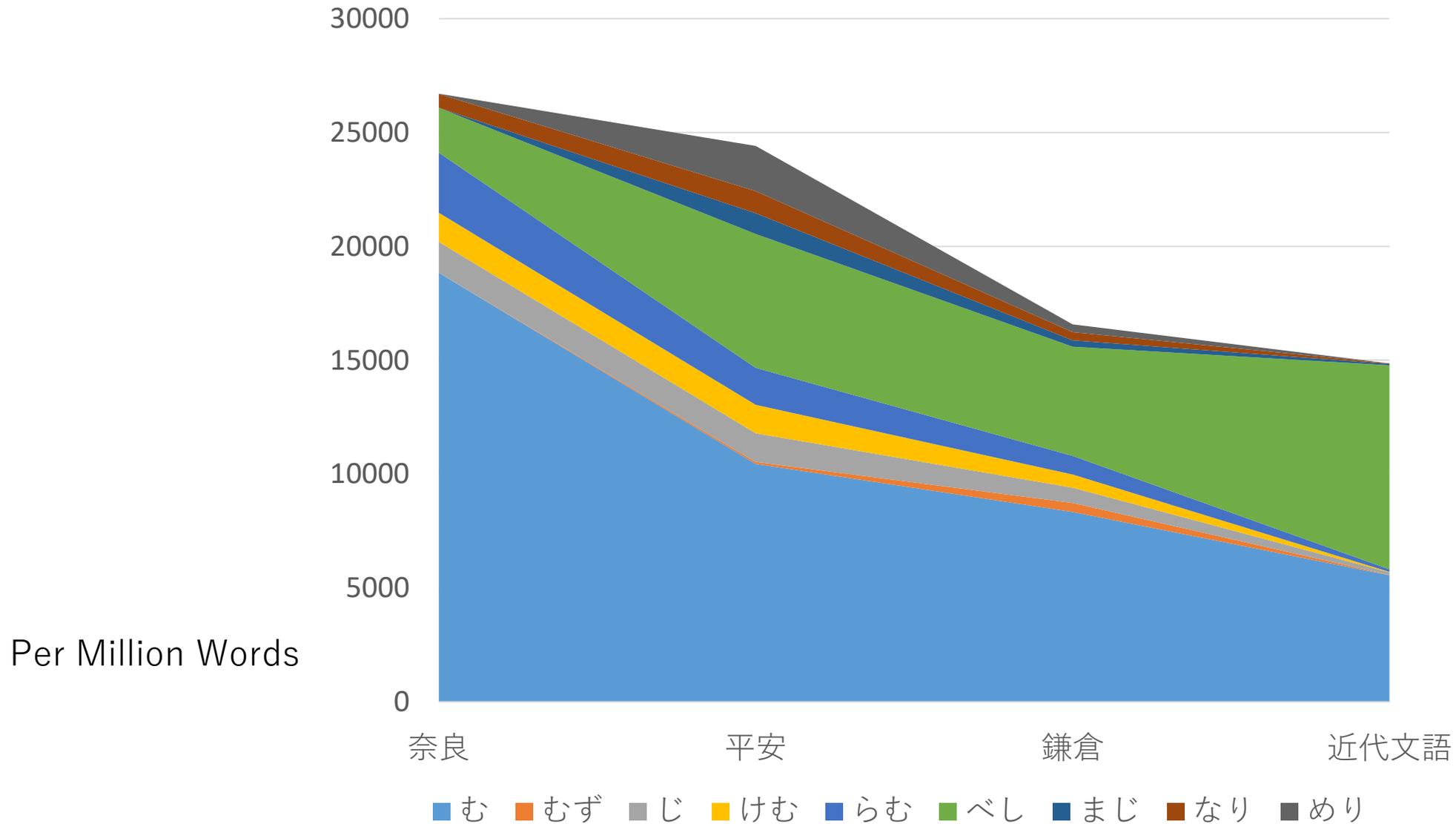
山と海



Modality (Auxiliary verbs)



Modality (Auxiliary verbs)



特徴語の抽出

- 対数尤度比 (Log-Likelihood Ratio, LLR)

$$\begin{aligned} G/2 &= \sum \text{観察値} (\log \text{観察値} - \log \text{期待値}) \\ &= a \log \frac{aN}{(a+b)(a+c)} + b \log \frac{bN}{(a+b)(b+d)} \\ &\quad + c \log \frac{cN}{(a+c)(c+d)} + d \log \frac{dN}{(b+d)(c+d)} \end{aligned}$$

- 他から期待される頻度と比べて、どれだけ多いか少ないか

歌集ごとの特徴語（形容詞）

語	金葉和歌集	古今和歌集	後拾遺和歌集	後撰和歌集	詞花和歌集	拾遺和歌集	新古今和歌集	千載和歌集	万葉集
賢い	-9.94876	-16.7869	-2.11058	-26.3492	-6.33119	-7.64882	-9.95342	-7.71912	131.3103
乏しい	-4.4072	-7.43597	-9.44282	-11.6708	-2.80474	-9.56765	-10.9757	-7.91391	104.0594
術無い	-4.60807	-2.99195	-9.87325	-12.2028	-2.93257	-10.0038	-11.476	-8.27463	97.52756
清い	-0.63127	-4.44909	-18.2387	-12.046	-0.14476	-14.1112	-4.58569	-3.47072	96.8924
良い	-4.5296	0.003582	-1.01714	-18.5622	-6.20255	-3.34575	-9.57577	-7.43956	70.8016
繁い	-0.69435	0.016517	-19.9639	-9.73172	-3.76372	-8.13308	-2.35568	-0.4346	65.67397
欲しい	-6.4181	0.138568	-13.7519	-4.47346	-4.08443	-1.29771	-3.85446	-11.5251	53.22038
日長い	-2.20072	-3.71305	-4.71506	-5.82745	-1.40055	-4.77739	-5.4804	-3.95168	51.90292
痛い	-4.4072	-0.8731	-9.44282	-3.20438	-0.11791	-4.34593	-2.77134	-7.91391	49.17338
尊い	0.15263	-5.23487	-6.64761	-8.21598	-1.97455	-6.73548	-2.92735	-5.57131	48.58252
鬱しい	-2.00042	-3.37509	-4.2859	-5.29703	-1.27308	-4.34255	-4.98157	-3.592	47.17404
麗しい	-2.40107	-4.05109	-1.18714	-6.35801	-1.52806	-5.21234	-5.97935	-4.31145	46.68617
間無い	-3.80486	-0.46106	-8.15217	-0.80118	-2.42142	-3.33663	-4.26108	-6.83225	42.01455

歌集ごとの特徴語（形容詞）

- **万葉集**：賢し、乏し、すべなし、清し、良し、繁し etc.
- **古今集**：恋し、憂し、正し、強し、侘し、疾し、あやなし etc.
- **新古今集**：深し、涼し、儂し、憂し、空し、脆し、凄し etc.

CHJ 今後の公開予定

2020年3月

- 奈良時代編 II **宣命** (続日本紀)
- 江戸時代編 III **近松浄瑠璃** (世話物)