

第7回産業日本語研究会・シンポジウム  
(20160229) 於東京・丸ビルホール

# 日本語の全体像を知るために

—国立国語研究所による言語資源整備—

前川 喜久雄

国立国語研究所 言語資源研究系・コーパス開発センター



# 国立国語研究所(東京都立川市)



国立国語研究所は1948年に設置された国立の研究機関。独立行政法人を経て、現在は大学共同利用機関法人人間文化研究機構が設置する大学共同利用機関の一つ。「国語及び国民の言語生活並びに外国人に対する日本語教育に関する科学的な調査研究並びにこれに基づく資料の作成及びその公表」を目的とする。

# 自己紹介

- 本来の専門は音声学・音声科学
- 90年代末から「言語資源」学に深入り
- 設計・開発・分析に携わってきた言語資源
  - 『日本語話し言葉コーパス』
  - 『現代日本語書き言葉均衡コーパス』
  - 『国語研日本語ウェブコーパス』

# なぜ今日私がここにいるか？

- 単なる人選ミス
- 「人工知能」で日本語(特許文書)を扱う
  - ⇒ 自然言語処理
    - ⇒ 現在の自然言語処理は統計ベース
      - ⇒ 言語資源
        - ⇒ 前川

# 自然言語処理と言語資源の関係

当初の自然言語処理は「ことばとは、このようなものだ」という規則を書き連ねることで実現していた。しかし、ことばは極めて多様で、常に変化し、人や文脈によって解釈が異なりうる。そのすべてを規則として記すのは、まったく現実的ではなかった。

中谷秀洋「特集統計的自然言語処理—ことばを扱う機会」  
IWANAMI DATA SCIENCE, Vol.2, p. 4, 2016.

# 言語内多様性

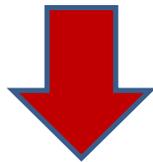
- 多様性の源
  - 歴史的多様性
    - 内発的変化
    - 外国語の影響
  - 地理的多様性
  - 創造的使用
  - 確率的変動(言語変異)
  - その他
    - 非母語話者
    - 機械翻訳

# 言語内多様性

- 多様性の源
  - 歴史的多様性
    - 内発的変化
    - 外国語の影響
  - 地理的多様性
  - 創造的使用
  - 確率的変動(言語変異)
  - その他
    - 非母語話者
    - 機械翻訳

# 言語内多様性を把握する手段

- 直観だけでは把握できない(例は後で)
- 客観的なデータが必要
- 対象を偏りなく代表するデータ
- できるだけ大量に
- 検索用の情報をつけて
- コンピュータで利用できる形式で
- 誰でも利用可能な形で



コーパス ~ 言語資源(language resources)の整備

# コーパス開発センターHP

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

日本語歴史コーパス

現代日本語  
書き言葉均衡コーパス

近代語のコーパス

日本語話し言葉コーパス

国語研日本語ウェブコーパス

中古和文UniDic

近代文語UniDic

UniDic

近代以前

現代語

コーパス利用申込

Subscription

最新情報

> 最新情報リスト

2016/01/04

お知らせ

BCCWJ, CSJの2015年度内の受付は2月12日(金)必着分までとなります

2015/11/09

お知らせ

共同研究 コーパスアノテーションの基礎研究の成果物

2015/10/27

日本語歴史コーパス

「洒落本コーパス」「人情本コーパス」の試作版を公開しました

2015/07/14

お知らせ

「中納言」のログイン方法が変更になりました

コーパス日本語学  
ワークショップ



UniDic - 形態素解析辞書 -

登録不要

少納言



Web茶まめ - 形態素解析支援ツール -

登録制

中納言



# 国語研による言語資源開発の経緯

- 1950年代から各種「語彙調査」を実施してきたがデータは公開しなかった
- 1990年代末にコーパス開発始動(KOTONOHA計画)
  - 『日本語話し言葉コーパス(CSJ)』(構築 1999~2003,公開2004)
  - 『太陽コーパス』(構築 1995~2005, 公開2005)
  - 『現代日本語書き言葉均衡コーパス(BCCWJ)』(構築 2006-2011, 公開2011)
  - 『日本語歴史コーパス(CHJ)』(構築2010~, 段階的に公開)
  - 『国語研日本語ウェブコーパス』(構築2011~2015, 公開2016予定)
  - 『多言語母語の日本語学習者横断コーパス(I-JAS)』(構築2012~, 部分試験公開2016)

詳しくは原稿記載のURLを参照してください

# 日本語コーパス整備の経緯 I : 1990年代

それ以前

近 代

現 代

新聞

青空文庫

新潮百冊

# 日本語コーパス整備の経緯Ⅱ：現状

それ以前

平安

CHJ

狂言

近代

青空文庫

新潮百冊

明  
六

国民  
之友

太陽

女性  
雑誌

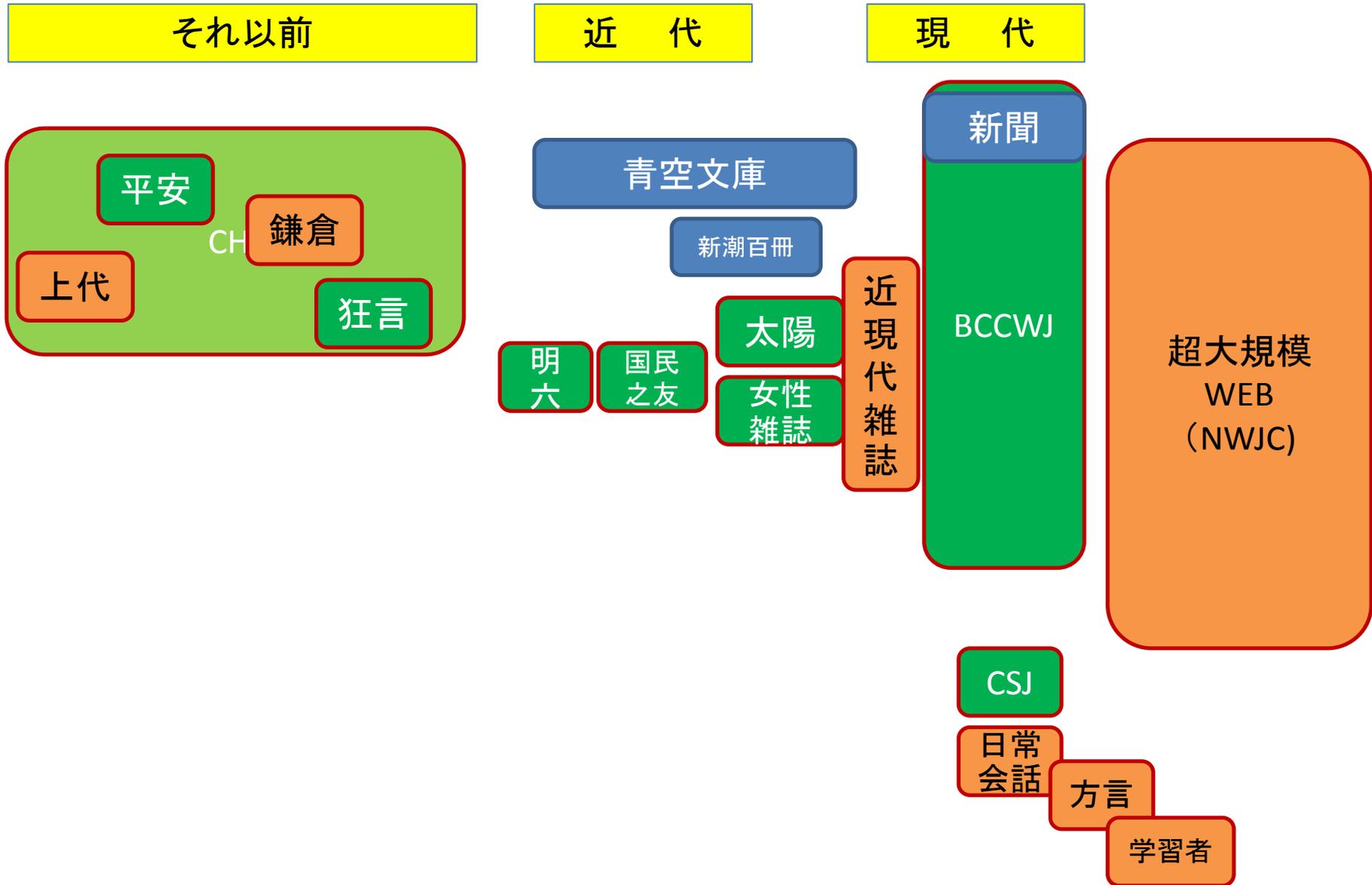
現代

新聞

BCCWJ

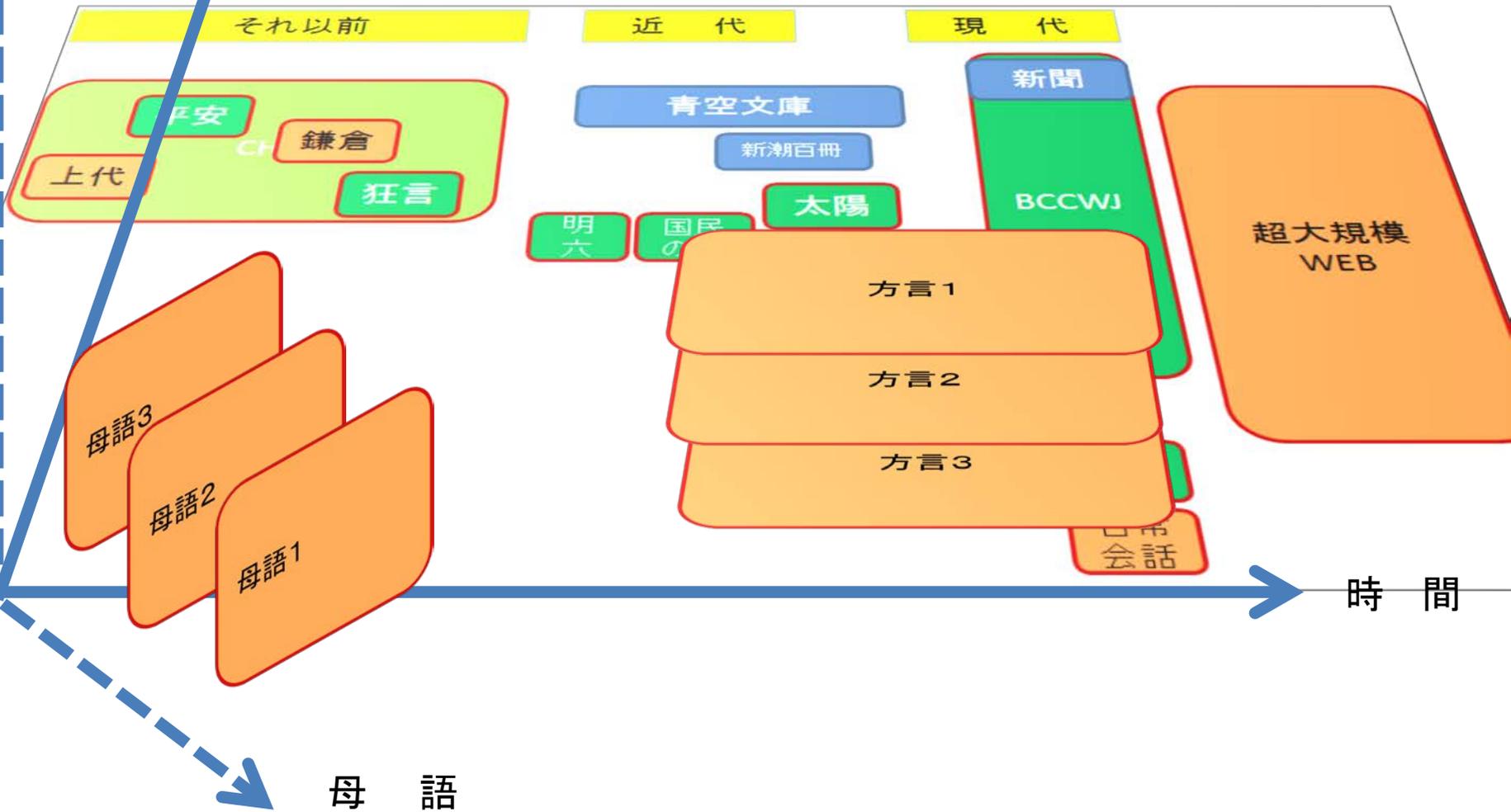
CSJ

# 日本語コーパス整備の経緯Ⅲ：2021年度末



地理（方言）

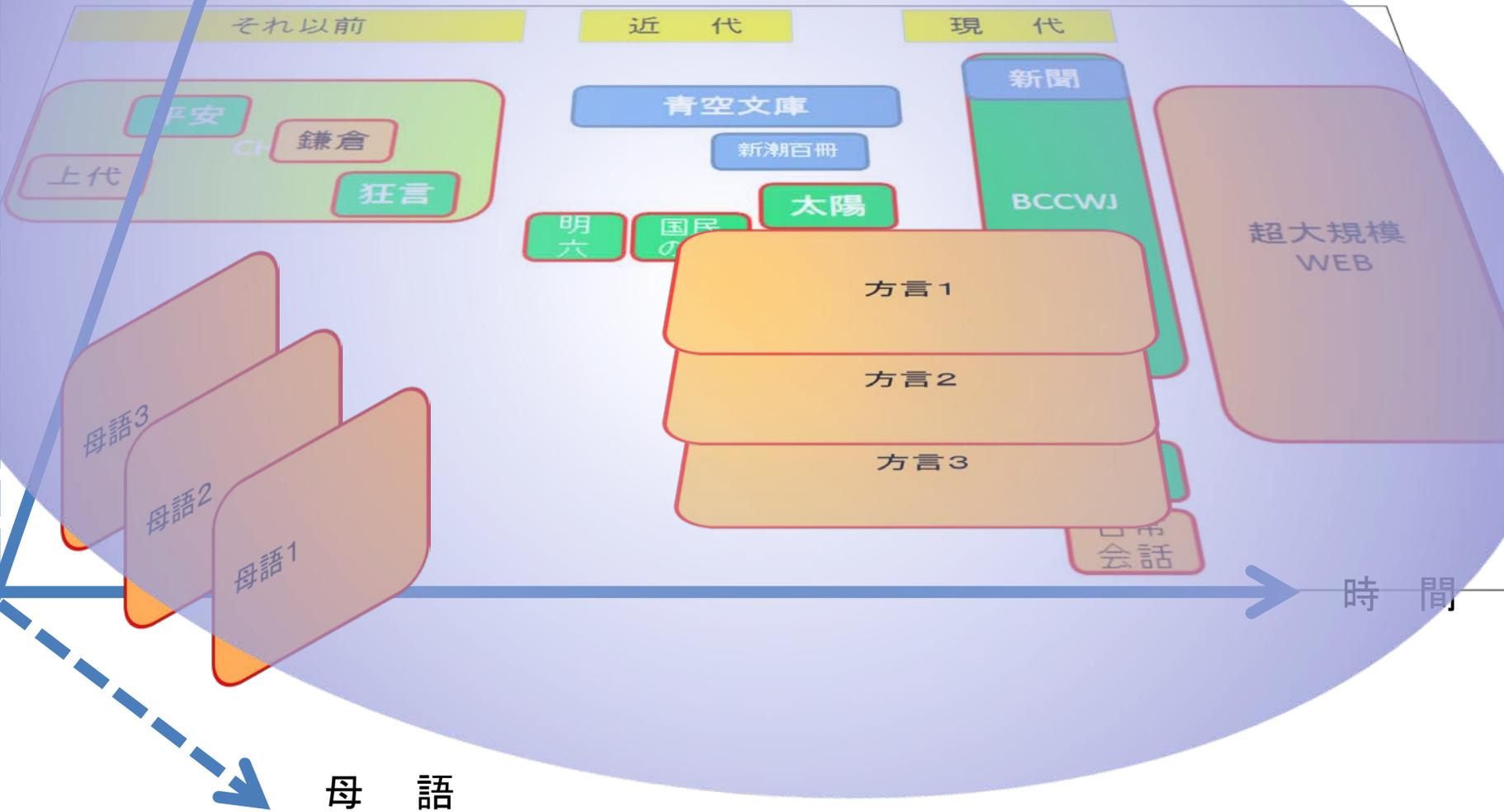
位相（書き言葉と話し言葉、それぞれのスタイル等）



地理（方言）

位相（書き言葉と話し言葉、それぞれの

言語内多様性



# コーパスが捉えた日本語の変異

# 内省できますか？

- 「NHK」はどのように発音されているか？
- 「日本」は？
- 「来られる」と「来れる」は話し言葉でどちらが多い？
- 「読むです」「行くです」等と書く人はいるか？
- 「読めれる」「行けれる」は？
- 「すらさえ」と「さえすら」はどちらが正用？
- Etc.

発音	頻度
エヌエチケー	132
エネーチケー	24
エヌエッチケー	9
エヌエイチケー	7
エヌエチケ	3
エネーチケ	3
エネエチケー	2
エヌエスケー	1
エヌチケー	1
エネーシケー	1

← 発音辞書の見出し

← 発音辞書の見出し

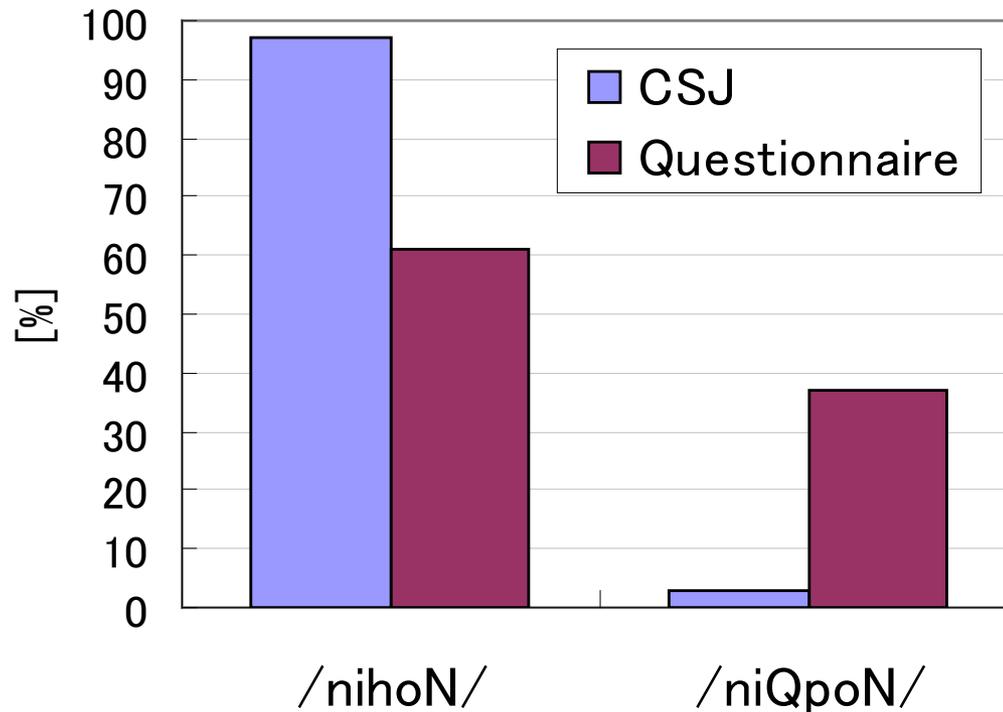
『日本語話し言葉コーパス』(講演などのモノローグ中心、752万語)  
の検索結果

LUW	ニッポン	ニホン	総計	%Nippon
日本一	9	31	40	22.5
日本代表	7	29	36	19.4
日本列島	1	24	25	4.0
日本	122	3108	3230	3.8
西日本	1	30	31	3.2
日本語教育	2	64	66	3.0
日本人	19	1019	1038	1.8
日本語	8	1591	1599	0.5
現代日本語		20	20	0.0
中世末期日本語		25	25	0.0
日本円		20	20	0.0
日本海		26	26	0.0
日本海側		30	30	0.0
日本語らしい		26	26	0.0
日本語らしさ		26	26	0.0
日本語学習者		23	23	0.0
日本語教師		21	21	0.0
日本語話者		148	148	0.0
日本酒		43	43	0.0

- CSJに関する限りニホンの圧勝。ニッポンは全体として3%以下
- 「日本一」「日本代表」では比較的ニッポンになりやすい

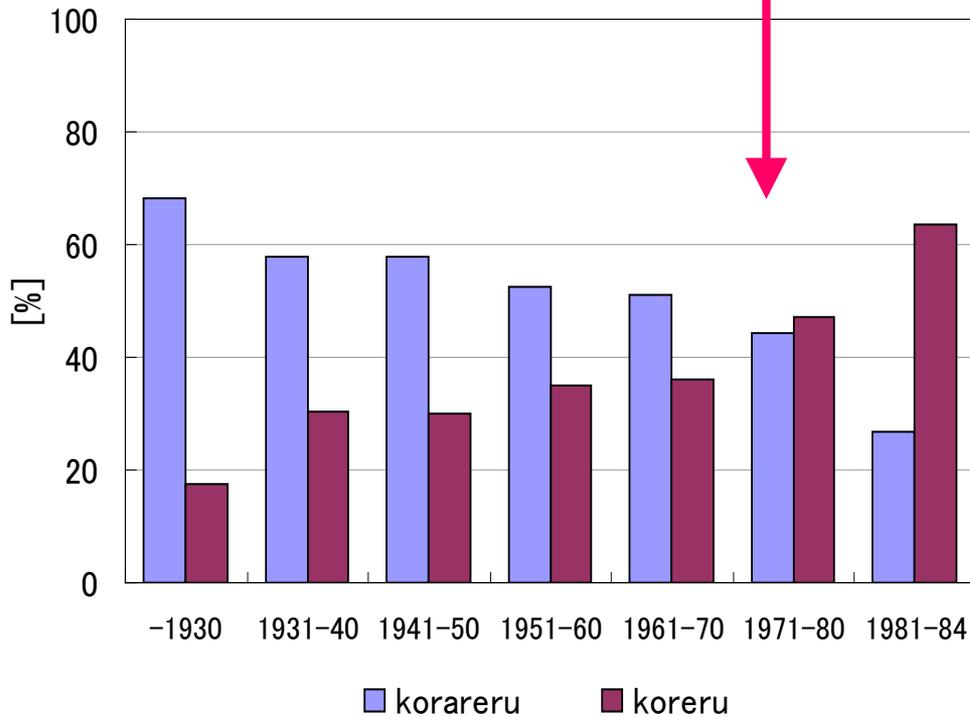
日本中		23	23	0.0
日本的		27	27	0.0

- 『日本語話し言葉コーパス』に記録された行動  
「ニホン」(97%)「ニッポン」(3%)
- NHK放送文化研究所によるアンケート(2004)  
「ニホン」(61%)「ニッポン」(37%)

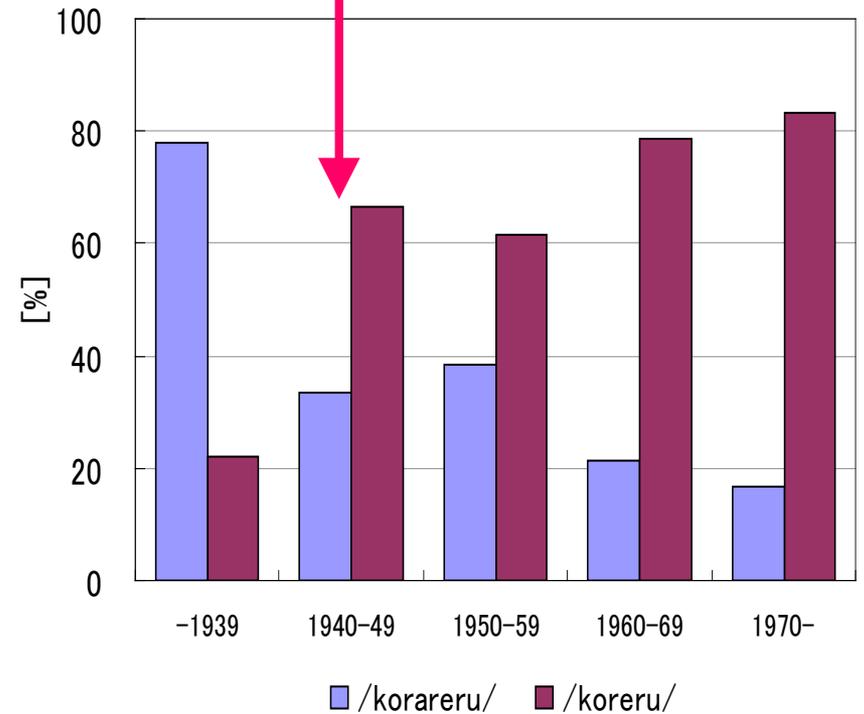


# 「コラレル」か「コレル」か

逆転のタイミングに30年のずれ



文化庁国語課による世論調査  
2001



『日本語話し言葉コーパス』における行動

# 観察結果

表現	BCCWJ(1億語)	NWJC(200億語)
形容詞+デス	11,000	393,619
動詞+デス	232	7,172
～シナ(サ)ソウ	75	312
～シナイベキ	11	205
～サセラレ／ラレサセ	479/0	7450/6
～スラサエ／サエスラ	0/0	1/2
～サエホド／ホドサエ	0/0	0/2

拙稿「コーパスの存在意義」(2013, 朝倉講座1巻)などでとりあげた誤用と正用の境界線上の例をBCCWJと開発中の『国語研日本語ウェブコーパス』で検索した。BCCWJは1億語、『国語研日本語ウェブコーパス』は現状で200億語以上。



# 「動詞＋デス」

## 『国語研日本語ウェブコーパス』

1. [デュエルマスターズ<sup>名 詞</sup> や<sup>助 詞</sup>] [ポケモン<sup>名 詞</sup> の<sup>助 詞</sup>] [カード<sup>名 詞</sup> も<sup>助 詞</sup>] [描か<sup>動 詞</sup> れ<sup>助 詞</sup> **てる** **です** <sup>助 詞</sup> ね<sup>名 詞</sup>]
2. [デモ<sup>名 詞</sup> が<sup>助 詞</sup>] [頻発<sup>動 詞</sup> して<sup>助 詞</sup> て<sup>助 詞</sup>] [危ない<sup>形 容 詞</sup>、<sup>記 号</sup>] [大手<sup>名 詞</sup> メディア<sup>名 詞</sup> の<sup>助 詞</sup>] [嘘<sup>名 詞</sup> も<sup>助 詞</sup>] [感<sup>名 詞</sup> 覚<sup>名 詞</sup> 的<sup>助 詞</sup> に<sup>助 詞</sup> か<sup>助 詞</sup> も<sup>助 詞</sup> し<sup>助 詞</sup> れ<sup>助 詞</sup> ない<sup>助 詞</sup> が<sup>助 詞</sup>、<sup>記 号</sup>] [中<sup>名 詞</sup> 学<sup>名 詞</sup> 生<sup>名 詞</sup>、<sup>記 号</sup>]  
[見破<sup>動 詞</sup> つ<sup>助 詞</sup> てる<sup>助 詞</sup> ん<sup>名 詞</sup> じ<sup>助 詞</sup> ゃ<sup>助 詞</sup> ない<sup>助 詞</sup> か<sup>助 詞</sup> と<sup>助 詞</sup>] [**思** **う** **で** **す** <sup>助 詞</sup> よ<sup>助 詞</sup>]
3. [デメリット<sup>名 詞</sup> を<sup>助 詞</sup>] [**教** **え** **る** **で** **す** <sup>助 詞</sup> か<sup>助 詞</sup> ・<sup>記 号</sup> ・<sup>記 号</sup> ・<sup>記 号</sup> 良<sup>形 容 詞</sup> い<sup>助 詞</sup> か<sup>助 詞</sup> も<sup>助 詞</sup> し<sup>助 詞</sup> れ<sup>助 詞</sup> ま<sup>助 詞</sup> せ<sup>助 詞</sup> ん<sup>助 詞</sup> ね<sup>助 詞</sup> ★<sup>記 号</sup>]
4. [デメリット<sup>名 詞</sup> は<sup>助 詞</sup>] [やり<sup>動 詞</sup> す<sup>助 詞</sup> ぎ<sup>助 詞</sup> る<sup>助 詞</sup> と<sup>助 詞</sup>、<sup>記 号</sup>] [く<sup>名 詞</sup> ら<sup>名 詞</sup> く<sup>名 詞</sup> ら<sup>名 詞</sup> **す** **る** **で** **す** <sup>助 詞</sup>]
5. [デメリット<sup>名 詞</sup> は<sup>助 詞</sup>、<sup>記 号</sup>] [両<sup>名 詞</sup> 方<sup>名 詞</sup> お<sup>名 詞</sup> 金<sup>名 詞</sup> が<sup>助 詞</sup>] [**い** **る** **で** **す** <sup>助 詞</sup>]
6. [デメリット<sup>名 詞</sup>、<sup>記 号</sup>] [日<sup>名 詞</sup> 本<sup>名 詞</sup> の<sup>助 詞</sup>] [裕<sup>名 詞</sup> 福<sup>名 詞</sup> で<sup>助 詞</sup>] [安<sup>名 詞</sup> 定<sup>名 詞</sup> し<sup>助 詞</sup> た<sup>助 詞</sup>] [生<sup>名 詞</sup> 活<sup>名 詞</sup> に<sup>助 詞</sup>] [変<sup>名 詞</sup> 化<sup>名 詞</sup> が<sup>助 詞</sup>] [**起** **こ** **る** **で** **す** <sup>助 詞</sup>]
7. [デミオ<sup>名 詞</sup> だ<sup>助 詞</sup> と<sup>助 詞</sup>] [イン<sup>名 詞</sup> パ<sup>名 詞</sup> ク<sup>名 詞</sup> 不<sup>名 詞</sup> 足<sup>名 詞</sup> か<sup>助 詞</sup> な<sup>助 詞</sup>、<sup>記 号</sup> っ<sup>助 詞</sup> て<sup>助 詞</sup>] [感<sup>名 詞</sup> じ<sup>名 詞</sup> です<sup>助 詞</sup> が<sup>助 詞</sup>] [扱<sup>名 詞</sup> い<sup>名 詞</sup>] [安<sup>形 容 詞</sup> い<sup>助 詞</sup> と<sup>助 詞</sup> は<sup>助 詞</sup>] [**思** **う** **で** **す** <sup>助 詞</sup> が<sup>助 詞</sup> ね<sup>助 詞</sup>] [w<sup>名 詞</sup>]
8. [デボ<sup>名 詞</sup> の<sup>助 詞</sup>] [周<sup>名 詞</sup> り<sup>名 詞</sup> に<sup>助 詞</sup>] [ま<sup>副 詞</sup> た<sup>副 詞</sup>] [ショッ<sup>名 詞</sup> ピン<sup>名 詞</sup> グ<sup>名 詞</sup> セン<sup>名 詞</sup> ター<sup>名 詞</sup> が<sup>助 詞</sup>] [**建** **つ** **で** **す** <sup>助 詞</sup> か<sup>助 詞</sup>]
9. [デブ<sup>名 詞</sup> ち<sup>名 詞</sup> ん<sup>名 詞</sup> は<sup>助 詞</sup>] [何<sup>名 詞</sup> を<sup>助 詞</sup>] [さ<sup>動 詞</sup> れ<sup>助 詞</sup> て<sup>助 詞</sup> も<sup>助 詞</sup>] [ジ<sup>名 詞</sup> ャ<sup>名 詞</sup> ッ<sup>名 詞</sup> と<sup>助 詞</sup>] [し<sup>動 詞</sup> **て** **る** **で** **す** <sup>助 詞</sup>]
10. [デフラグ<sup>名 詞</sup> 統<sup>名 詞</sup> 合<sup>名 詞</sup> 版<sup>名 詞</sup> Win<sup>名 詞</sup> の<sup>助 詞</sup>] [チェ<sup>名 詞</sup> ック<sup>名 詞</sup> ディ<sup>名 詞</sup> スク<sup>名 詞</sup> と<sup>助 詞</sup> デ<sup>名 詞</sup>] [フラ<sup>名 詞</sup> グ<sup>名 詞</sup> を<sup>助 詞</sup>] [有<sup>名 詞</sup> 効<sup>名 詞</sup> に<sup>助 詞</sup>] [**使** **う** **で** **す** <sup>助 詞</sup>]
11. [デフラグ<sup>名 詞</sup> に<sup>助 詞</sup> つ<sup>助 詞</sup> い<sup>助 詞</sup> て<sup>助 詞</sup> は<sup>助 詞</sup>、<sup>記 号</sup>] [諸<sup>名 詞</sup> 般<sup>名 詞</sup> 諸<sup>名 詞</sup> 々<sup>名 詞</sup> の<sup>助 詞</sup>] [件<sup>名 詞</sup> が<sup>助 詞</sup>] [**あ** **る** **で** **す** <sup>助 詞</sup> が<sup>助 詞</sup>]
12. [デフラグ<sup>名 詞</sup> なる<sup>動 詞</sup> を<sup>助 詞</sup>] [知<sup>動 詞</sup> ら<sup>助 詞</sup> ん<sup>助 詞</sup> で<sup>助 詞</sup> も<sup>助 詞</sup>、<sup>記 号</sup>] [HDD<sup>名 詞</sup> 内<sup>名 詞</sup> の<sup>助 詞</sup>] [ファ<sup>名 詞</sup> イ<sup>名 詞</sup> ル<sup>名 詞</sup> を<sup>助 詞</sup>] [整<sup>名 詞</sup> 理<sup>名 詞</sup> 整<sup>名 詞</sup> 頓<sup>名 詞</sup> し<sup>助 詞</sup> て<sup>助 詞</sup> **く** **れ** **る** **で** **す** <sup>助 詞</sup>]

# 「二重可能」(レ足す言葉)

BCCWJ

動詞	レル	レナイ	レタ
行ケ	3	2	1
聞ケ	0	0	0
書ケ	0	0	0
遊ベ	0	0	0
歩ケ	0	0	0
出来	0	0	0
描ケ	0	0	0
飛ベ	0	0	0
聴ケ	0	0	0
読メ	0	0	0

NWJC

動詞	レル	レナイ	レタ
行ケ	32	28	2
聞ケ	5	1	2
書ケ	4	1	0
遊ベ	4	0	0
歩ケ	3	0	0
出来	2	0	0
描ケ	2	2	1
飛ベ	1	0	1
聴ケ	1	1	0
読メ	1	1	0

100億語規模のコーパスではじめて把握できるようになる低確率の  
(しかし実際に生起する)言語現象が少なくない。

# 観察結果

表現	BCCWJ(1億語)	NWJC(200億語)
形容詞+デス	11,000	393,619
動詞+デス	232	7,172
~シナ(サ)ソウ	75	312
~シナイベキ	11	205
~サセラレ/ラレサセ	479/0	7450/6
~スラサエ/サエスラ	0/0	1/2
~サエホド/ホドサエ	0/0	0/2

# 「～シナイベキ」

1. [デフラグ について は、] [早く] [完了 さ せる には] [スクリーンセ이버 は] [同じ よう に] [起動 さ せ **ない べき** です]
2. [タイトレ の] [場合、] [ザラ 場 で] [「エントリー する べき が、」] [し **ない べき** が、] [それが] [問題 た]
3. [テレビ の] [前で] [「光 ちゃん は」] [可愛い けど、] [録画 する べき が] [し **ない べき** が ・ ・ ・ ] と [葛藤 する]
4. [テレビ と] [現実 を] [一緒 に して しまう] [人 は] [見 **ない べき** です]
5. [テレビ っ て、] [何 を] [伝える べき で] [何 を] [伝え **ない べき** な の か を] [わかま えて いる と] [思 い ます か]
6. [チョコレート を] [あげる べき が、] [あげ **ない べき** が で] [悩ん で いる] [人 が] [多い の で は ない で し ょ う か]
7. [チャンス は] [逃さ **ない べき** た ろ う な ー と] [思 い ま し た]
8. [チャリ で] [たばこ 吸わ **ない べき** です よ ね ・ ・ ・ ]
9. [チャイナ エアライン は] [乗ら **ない べき**]
10. [子ビ 太 に] [見せる べき が] [見せ **ない べき** が] [悩ん た ま ま、]
11. [ダッシュ する] [時 も] [ダッシュ す べき が] [し **ない べき** が の] [迷い も] [見て 取れ た し、]
12. [タル は、] [合理的 には] [3つ を] [越え **ない べき** で ある]
13. [タク : タクシー に] [乗ら **ない べき**] [日]

# 観察結果

表現	BCCWJ(1億語)	NWJC(200億語)
形容詞+デス	11,000	393,619
動詞+デス	232	7,172
～シナ(サ)ソウ	75	312
～シナイベキ	11	205
～サセラレ／ラレサセ	479/0	7450/6
～スラサエ／サエスラ	0/0	1/2
～サエホド／ホドサエ	0/0	0/2

拙稿「コーパスの存在意義」(2013, 朝倉講座1巻)などでとりあげた誤用と正用の境界線上の例をBCCWJと開発中の日本語ウェブコーパスで検索した。BCCWJは1億語、超大規模Webコーパスは現状で200億語以上。

# 「～ラレサセ」

(surface:られ)[0,0] (surface:させ)[1,1]

1 - 6 / 6 (0.969801902771 sec)

1. [コイツ だけ は] [ヤバ い、 と] [サンジ は] [自分 の] [身 を] [切り裂く] [巨 き 過ぎる] [男根 に、] [脳 まで]  
[突き通さ れる よう な] [感 覚 に] [一瞬] [息 を] [止め られ させ ながら、] [そう] [思う]
2. [クラブ の] [存続 の] [た め に] [何 か] [手 伝 える] [こ と は] [な い か、 と] [考 え られ させ まし た]
3. [オザンナ] [(主 の] [名 に よ り] [来 た れ る] [者 は] [祝 せ られ させ] [給 え、] [最 高 き] [と ころ まで])
4. [イエズス・キリスト、] [祝 せ られ させ] [賜 え]
5. [アルツハイマー の] [人 は、] [人 の] [生 き 方 を] [考 え られ させ ら れ まし た ね] (\*工\*)
6. [それに、] [す っ ご く] [謎 っ て] [か ん じ で] [読 ん で い た ら] [よ く] [考 え られ させ ら れ ま す] W W

# レジスターの重要性

表現	BCCWJ	NWJC
～ヲ泣ク	1	0
～ヲ死ヌ	3	0

BCCWJで得られた用例のレジスターは、PV(書籍)、LB(図書館書籍)、OB(ベストセラー)、OV(韻文)。先の検索結果は、Webのデータには、相対的に見てこの種の文書が少ないことを示唆しているように思える。

⇒ BCCWJのような均衡コーパスの必要性

# 『特許ライティングマニュアル』より

ルール第F条の1：並立表現の並立要素が同じ表現になるように整える。

改善例：

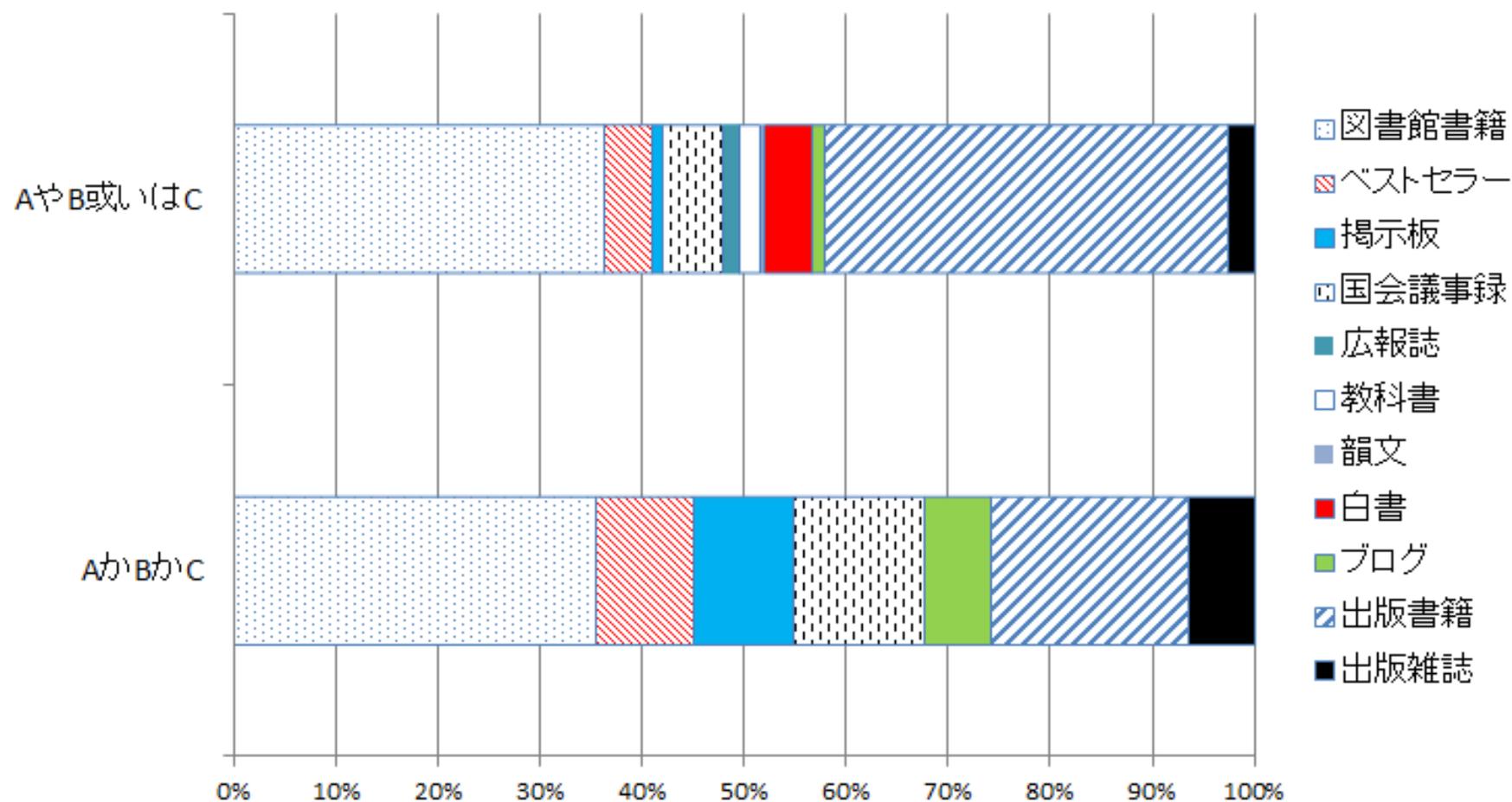
「水やスチーム、或いは、薬品など」  
⇒「水かスチームか薬品など」

BCCWJの検索結果

AやB(、)或いはC ⇒ 288例

AかBかC ⇒ 31例

# レジスター差



おわりに

# 言語研究と情報科学の連携

- 現代の言語資源開発は、情報科学なくして考えられない。機械学習(人工知能)の研究とも関係が深い
- 言語研究からの情報科学への貢献は、日本語の全体像を伝える良質な言語資源の提供
- 今後、複雑多様な言語を過度に単純化するのではなく、複雑さを保ってモデル化する研究を進めるにあたっては、両者の協力が不可欠