

金田一春彦「国語動詞の一分類」(1950)



動詞として挙げたいものは、(中略)、ある状態を示す動詞と言いたいものであるが、例えば「山が登える」がこれである。

「登る」は「いる」の形で状態を表すのに用い、ただ「登る」だけの単独の形で動作・作用を表すために用いることがないのを特色とする。

寺村秀夫「ムードの形式と否定」(1979)



「思弁の遊びのように思われるかもしれないが必ずしも日本語を外国語として習得しようとする者にとっては、**ダロウ**の前の部分を否定にする、**ダロウ**は否定の形である、**〜べき**のように、その否定は**〜べき(ダ)**イベキダという形はない、一方、**ハズダ**のように、内側から否定して〜べきものもある、等々といったことは必要な文法的知識であるに相違なく、また「内」外」両面の否定があり得るとき、それはどう意味的に違ってくるのか、といったこともまた必要な文法的知識に属するであろう、これらの問いに多少とも統一的に答えるためには、どうしても上のような議論を避けて通るわけにはいかないのである。

さて、実際は

理論をうのみにするのではなく、世の中で実際に用いられている用例を偏りなく、大量に観察してみる

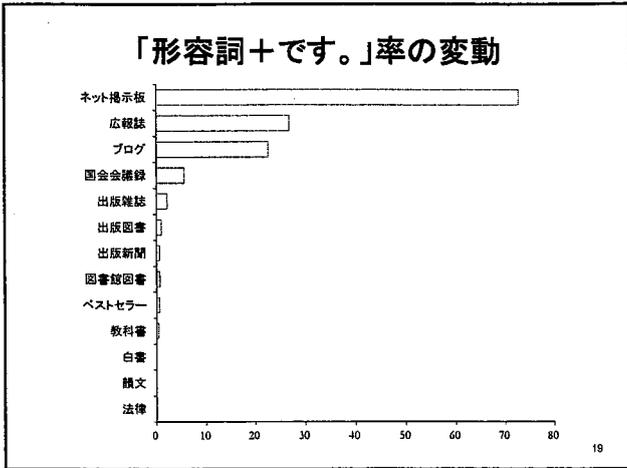
⇒ 国立国語研究所が昨年公開した『現代日本語書き言葉均衡コーパス』(BCCWJ,1億語規模)を検索

まずは「そびえている」から

「そびえている」が多数派

195 年間の使用例が挙げられた。

ラングム	品詞類	品	種別類	ラテン語	漢字	読み仮名	出典
1950-1959	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
1960-1969	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
1970-1979	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
1980-1989	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
1990-1999	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2000-2009	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2010-2019	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2020-2029	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2030-2039	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2040-2049	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2050-2059	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2060-2069	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2070-2079	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2080-2089	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2090-2099	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2100-2109	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2110-2119	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2120-2129	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2130-2139	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2140-2149	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2150-2159	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2160-2169	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2170-2179	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2180-2189	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2190-2199	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2200-2209	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2210-2219	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2220-2229	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2230-2239	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2240-2249	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2250-2259	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2260-2269	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2270-2279	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2280-2289	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2290-2299	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2300-2309	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2310-2319	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2320-2329	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2330-2339	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2340-2349	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2350-2359	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2360-2369	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2370-2379	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2380-2389	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2390-2399	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2400-2409	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2410-2419	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2420-2429	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2430-2439	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2440-2449	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2450-2459	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2460-2469	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2470-2479	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2480-2489	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2490-2499	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2500-2509	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2510-2519	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2520-2529	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2530-2539	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2540-2549	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2550-2559	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2560-2569	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2570-2579	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2580-2589	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2590-2599	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2600-2609	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2610-2619	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2620-2629	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2630-2639	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2640-2649	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2650-2659	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2660-2669	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2670-2679	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2680-2689	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2690-2699	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2700-2709	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2710-2719	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2720-2729	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2730-2739	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2740-2749	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2750-2759	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2760-2769	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2770-2779	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2780-2789	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2790-2799	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2800-2809	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2810-2819	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2820-2829	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2830-2839	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2840-2849	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2850-2859	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2860-2869	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2870-2879	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2880-2889	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2890-2899	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2900-2909	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2910-2919	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2920-2929	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2930-2939	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2940-2949	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2950-2959	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2960-2969	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2970-2979	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2980-2989	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
2990-2999	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3000-3009	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3010-3019	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3020-3029	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3030-3039	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3040-3049	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3050-3059	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3060-3069	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3070-3079	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3080-3089	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3090-3099	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3100-3109	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3110-3119	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3120-3129	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3130-3139	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3140-3149	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3150-3159	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3160-3169	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3170-3179	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3180-3189	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3190-3199	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3200-3209	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3210-3219	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3220-3229	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3230-3239	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3240-3249	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3250-3259	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3260-3269	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3270-3279	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』
3280-3289	動詞	1	他動詞	そびえる	そびえる	そびえる	『現代国語』



ここで、ちょっと注意

- ここまで挙げてきた例は「使われているかもしれないが、誤った日本語だ」という立場
- 規範を示す立場(国語教育)としてはありうる
- しかし「正しい」日本語だけを研究していたのでは、日本語の情報処理はできない。外国人に対する日本語教育にも支障がでる
- 規範を離れた(非情な)立場での研究が必要
肺魚も魚、ダチョウも鳥という立場

ここまできをまとめると

- ある表現を「言うか言わないか」あるいは「どう言っているか」の判断は予想外に難しい。ケースによっては内省に著しい限界がある
- 問題はなにが「ケースによる」かがわからないこと
- 非情に世界を見渡してデータを集める必要がある

非情な言語研究のために何が必要か

- 偏りのない大量の電子化データ構築 コーパス(corpus)
- データへの検索用情報付与 形態素解析、係り受け解析 etc.
- データの分析技術 種々の統計手法等

言語資源学

国立国語研究所による日本語コーパスの開発

言語資源研究系が開発してきたコーパス

言語資源研究系が開発してきたコーパス

- 全訳古語コーパス
- 現代日本語コーパス
- 日本語訳し書庫コーパス
- 現代日本語 教育書コーパス
- 現代日本語 新聞コーパス
- 現代日本語 雑誌コーパス
- 現代日本語 小説コーパス
- 現代日本語 論文コーパス
- 現代日本語 法律コーパス
- 現代日本語 白書コーパス
- 現代日本語 教科書コーパス
- 現代日本語 ベストセラーコーパス
- 現代日本語 ネット掲示板コーパス
- 現代日本語 広報紙コーパス
- 現代日本語 ブログコーパス
- 現代日本語 国会会議録コーパス
- 現代日本語 出版雑誌コーパス
- 現代日本語 出版図書コーパス
- 現代日本語 出版新聞コーパス
- 現代日本語 図書館図書コーパス
- 現代日本語 ベストセラーコーパス
- 現代日本語 教科書コーパス
- 現代日本語 白書コーパス
- 現代日本語 論文コーパス
- 現代日本語 法律コーパス

『日本語話し言葉コーパス』
Corpus of Spontaneous Japanese (CSJ)
2004年公開

25

『日本語話し言葉コーパス』

Corpus of Spontaneous Japanese

利用目的：音声認識研究と言語研究

対象：現代日本語の知的なモノローグ

規模：662時間(752万語)

学会での研究発表と一般的なスピーチが中心(330H)

対照用に対話音声(12H)と朗読音声(20H)も

付加情報：非常に豊富

音声信号、話者情報、転記テキスト、形態論(品詞)情報、節単位情報、係り受け情報、要約・重要文情報、談話構造情報、分節ラベル、韻律ラベル(X-JToBI)、講演が与える印象の評価

⇒ XML文書として統合

開発：国立国語研究所、情報通信研究機構、東京工業大学

世界の主要話し言葉コーパスとの比較

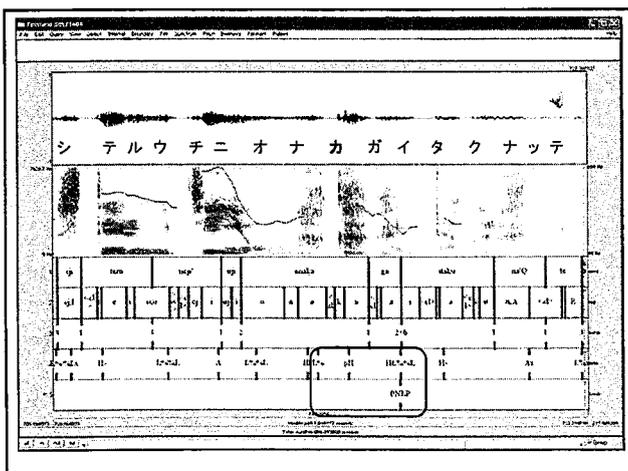
名称	音声タイプ	音声公開	話者	時間	語数(万)	転記	形態	分節	イントネ	言語
CSJ	独話中心	○	1408	661	752	○	○	○	○	日本語
Switchboard Corpus	対話	○	543	250	---	○	○	○	×	米語
CSPAE	対話	○	---	6	5	○	×	×	○	米語
CSPAIE	演説・記者会見、教授会	×	400	---	200	○	○	×	×	米語
Map Task Corpus	対話	○	64	34	---	○	×	×	×	英語
BNC	独話・対話・朗読・放送	×	2700	1220	1000	○	△	×	×	英語
現代漢語口語対話語料庫	対話	○	60	27	---	○	○	○	×	中国語

音声を公開しているコーパスとしては世界最大、付加情報は抜群に豊富

X-JToBI

- 東京語の自発音声を対象とした韻律ラベリングスキーム
- J_ToBI(實際上朗読音声用)を拡張
- 音声信号との対応づけが可能
- Word層、Segment層、Tone層、BI層、Prominence層、Misc層にわけて情報を記述
⇒ BI: Break Index 後半ではただBreakと呼んでいる

⇒ 自発音声の研究には非常に有益



kikuo maekawa
National Institute for Japanese Language and Linguistics
Phonetics, Language Acquisition, etc.
VeeRed email of nupl.ac.jp

■ 日本語

■ CSJを

■ 博士論

■ 代表的論文

- 母音
- /z/
- /b/
- PNI
- 日本語

能の解明

よど

『太陽コーパス』

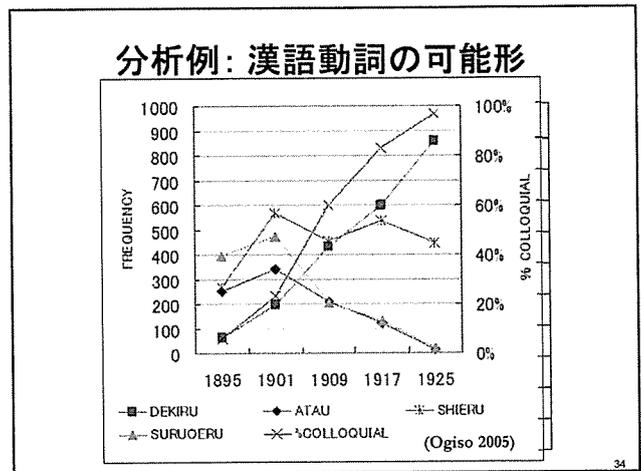
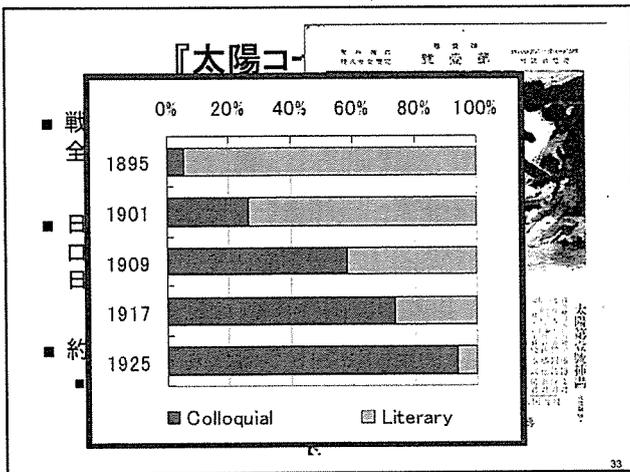
Taiyo Corpus
2005年公開

31

『太陽コーパス』

- 戦前の有名な総合雑誌である『太陽』(1895-1925)の全文テキストコーパス
- 日本語書き言葉の文法が、いわゆる文語文法から口語文法へと急速に変化した時期(言文一致期)の日本語の記録
- 約700万語(形態素解析されていないので推計)
 - 近い将来に形態素解析データを公開予定

32



『現代日本語書き言葉均衡コーパス』

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

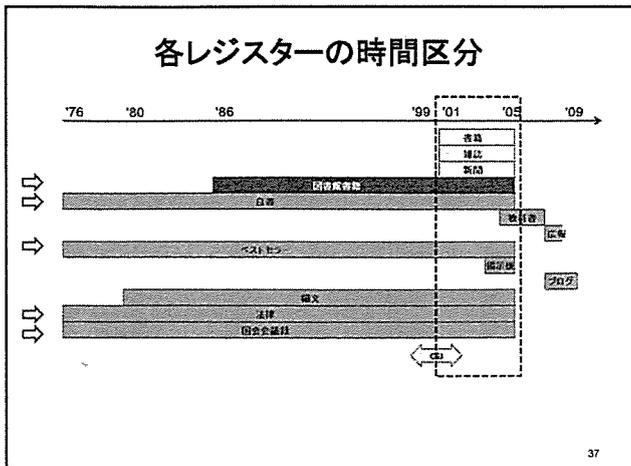
特定目的(非母集団)サブコーパス
ウェブ上の文書、白書、教科書、国会会議録、ベストセラー等
対象期間はさまざま、最長30年。3500万語

35

BCCWJを構成するレジスター

サブコーパス	レジスター	サンプル(個)	語数(万語)
出版サブコーパス	書籍(PB)	10,117	2,855
	雑誌(PM)	1,996	444
	新聞(PN)	1,473	137
図書館サブコーパス	書籍(LB)	10,551	3,038
	白書(OW)	1,500	488
	教科書(OT)	412	93
	広報紙(OP)	354	376
	ベストセラー(OB)	1,390	374
	Yahoo!知恵袋(OC)	91,445	1,026
	Yahoo!ブログ(OY)	52,680	1,019
特定目的サブコーパス	随文(OV)	252	25
	法律(OL)	346	108
	国会会議録(OM)	159	510
	合計	172,675	10,493

36



日本語コーパス構築の問題点

38

- ### 日本語コーパス開発上の困難
- 日本の著作権法
権利者の保護に手厚い。ユーザーとしてはきわめて不便。
 - 日本語表記の特異性
多くの文字を併用(漢字、ひらがな、カタカナ、Roman alphabet 等)
 - いわゆる膠着語(こうちやくご)的性格
語の境界がはっきりしない ⇒ 分かち書きの習慣がない
- 39



- ### 日本語検索の難しさ
- 「国語」 ⇒ 国語、中国語、韓国語、外国語, etc.
 - 「リズム」 ⇒ リズム、アルゴリズム、フォルマリズム, etc.
 - 「明かり」 ⇒ 「明り」「灯り」「燈」「あかり」, etc.
 - 「やはり」 ⇒ 「やはり」「やっぱり」「やはし」, etc.
- 41

感動詞「嗚呼」

- BCCWJには33種類の表記がある

あ～ ああ ああ あゝ あア あー アゝ アア
 アア アー 吁嗟 嗚呼 嗟 嗟呼 噫 於戲
 臆 あ～っ ああっ あーっ アーッ 嗚乎 あああ
 あああ あああ ああー アアア あ～あ あーあ
 あーア アーア 嗚～呼 あああ～

42

人名「ヒロシ」

- BCCWJには75種類の表記がある

ひろし ヒロシ 博 博史 博司 博士 博師 博志 博至
博資 啓 啓史 啓志 坦 大 大志 央 宏 宏史 宏司
宏士 宏志 宏至 容 寛 寛志 広 広史 広士 広志
廣 廣司 廣志 弘 弘史 弘司 弘士 弘志 弘至 弘視
怒 日呂志 昊 普 比呂志 比露思 汎 洋 洋司
洋士 洋志 洪 洽 浩 浩嗣 浩士 浩志 浩至 溥 滉
熙 碩 礼嗣 秦 紘 綽 裕 裕史 裕士 裕志 飛呂志
鴻 熙 Hiroshi

43

複合動詞「沸き起こる」

- 終止形だけで20種

わきおこる わき起こる わき起る 沸きおこる 沸き起こる 沸き起る
沸起こる 沸きおこる 沸き起こる 沸き起る 沸起こる 沸きおこる
沸き起こる 沸き起る 沸起こる 沸起る わきおこれる 沸き起これる
沸き起れる 沸起れる

- すべての活用形を考慮すると324種

沸きおこれ 沸きおこりゃ 沸きおころ 沸きおこれ 沸きおころう
沸きおころっ 沸きおころ 沸きおこら 沸きおこん 沸きおこる 沸きおこん
沸きおこる 沸きおこん 沸きおこり 沸きおこっ 沸きおこん 沸きおこりゃ
沸き起これ 沸き起こりゃ 沸き起ころ 沸き起これ 沸き起ころう
沸き起ころっ 沸き起ころ 沸き起こら 沸き起こん 沸き起こる 沸き起こん
沸き起こる 沸き起こん 沸き起こり 沸き起こっ 沸き起こん 沸き起こりゃ
Etc.

44

日本語検索の難しさについての結論

- 日本語のコーパスは単にテキストを電子化しただけでは使い物にならない
- テキストを単語に区切って、表記の多様性と用言の活用による多様性を吸収し、いわゆる「原形」で検索できるように加工する必要がある
- 形態素解析(形態論情報アノテーション)

45

日本語の形態素解析

- 情報科学の研究者達の努力により、1990年前後に実用化
京都大学 JUMAN
奈良先端大 ChaSen, MeCab
- ただし言語研究に使うには解析用辞書に問題が
辞書の内部で「単語」の定義に一貫性がない

46

Yahoo!	IPA	JUMAN
国立国会図書館	国立	国立
	国会図書館	国会
		図書
		館
国立 公文書館	国立 公文書 館	国立 公文書 館
国立 歴史 民俗 博物館	国立 歴史 民俗 博物館	国立 歴史 民俗 博物 館

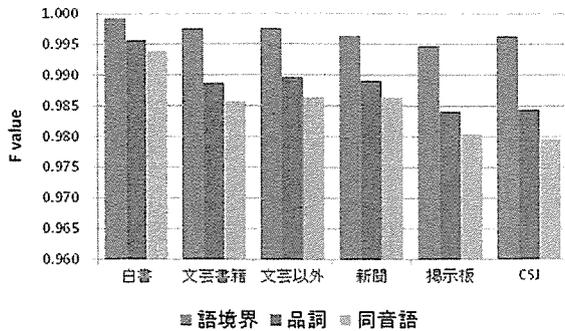
47

Yahoo!	IPA	JUMAN	UniDic
国立国会図書館	国立	国立	国立
	国会図書館	国会	国会
			図書
			館
国立 公文書			国立 公
	館	館	文書 館
国立 歴史 民俗 博物館	国立 歴史 民俗 博物館	国立 歴史 民俗 博物 館	国立 歴史 民俗 博物 館

UniDicはCSJの経験を踏まえて
BCCWJ用に開発した電子化辞書。
同じ粒度で「語」を切り出せる

48

UniDic + MeCabによる解析精度



二重形態素解析

短単位	短単位品詞	長単位	長単位品詞
公害	名詞-普通名詞-一般	公害紛争処理法	名詞-普通名詞-一般
紛争	名詞-普通名詞-サ変可能		
処理	名詞-普通名詞-サ変可能		
法	名詞-普通名詞-一般		
に	助詞-格助詞	における	助詞-格助詞
おけ	動詞-一般		
る	助動詞		
公害	名詞-普通名詞-一般	公害紛争処理	名詞-普通名詞-一般
紛争	名詞-普通名詞-サ変可能		
処理	名詞-普通名詞-サ変可能		
の	助詞-格助詞	の	助詞-格助詞
手続	名詞-普通名詞-サ変可能	手続	名詞-普通名詞-一般
は	助詞-係助詞	は	助詞-係助詞

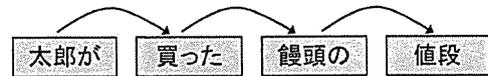
短単位～辞書見出し語 長単位～複合語(複合辞)

形態素解析の次のステップ

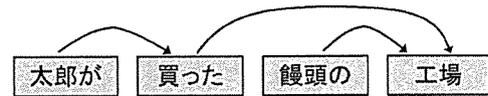
- 係り受け構造のアノテーション
語と語の修飾関係
- 動詞項構造のアノテーション
動詞と目的語などの関係
- 文末・節末のアノテーション
文末が句読点などで明示されているとは限らない
- 意味のアノテーション
多くの語には複数の用法がある
- 拡張モダリティのアノテーション
書き手の意図や判断
- その他、時間表現、長単位内部構造、等々

係り受け構造解析

- 太郎が買った饅頭の値段



- 太郎が買った饅頭の工場



※音声合成にも有用。

節(従属句)の特徴

Cf. 南不二男(1974)『現代日本語の構造』

A ながら つつ	歯を磨かないながら	電話がかかってきた
	歯を磨くだろうながら	歯磨き粉が
		歯磨き粉が
B のに ので ば		
	あんなに勉強するだろうのに	
C けれど から が		

BCCWJの一般公開

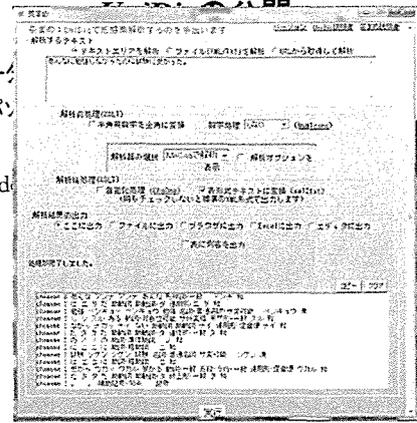
文字列検索専用IF
無償、契約不要

短単位・長単位の
検索IF、無償、要
契約

UniDicの公開

- データ全体を無償公開
パソコンやスマートフォンのOSで既に利用されている。
- Windows PC用インターフェース「茶まめ」

55



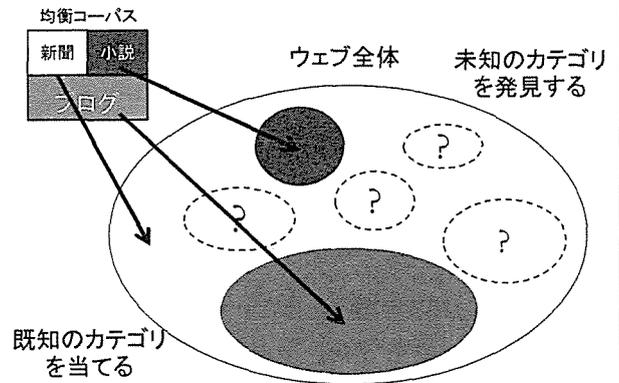
56

いまはどこに向かっているか

- 各種アノテーションの基礎研究
多くの大学と共同研究
- 日本語歴史コーパスの準備研究
平安から江戸時代までの代表的文学作品
形態素解析も
- 超大規模コーパスの構築
Webを母集団として100億語規模
サンプルのレジスター推定

57

ウェブのレジスター推定問題



レジスターの推定実験

		子 潤					
		掲示板	白書	ブログ	書籍	雑誌	新聞
正 解	掲示板	42	0	5	2	1	0
	白書	0	49	0	0	0	1
	ブログ	3	0	40	3	0	4
	書籍	2	0	0	47	1	0
	雑誌	0	1	3	12	27	7
	新聞	0	0	0	3	5	42

BCCWJの6レジスターからランダムに選んだ300サンプルをSVM (support vector machines)という手法で分類。データは文長、品詞比率、語種比率

59

まとめ

コンピュータで日本語をきちんと研究できるようにするためには、

- 各種コーパスの構築
 - 検索用情報(形態素、係り受け等)の付与
 - 検索用インターフェースの開発
 - 著作権処理
- などの壁を超える必要があった。

これらの仕事を一応終えたいま、過去の日本語を対象とした歴史コーパスや、100億語規模の超大規模コーパスの開発を進めている。また遠くない時期に、日常会話を対象とした新しい話し言葉コーパスにも着手したいと考えている。

60