

## Web を母集団とした超大規模コーパスの開発 —収集と組織化—

浅原 正幸 (国立国語研究所コーパス開発センター)  
今田 水穂 (国立国語研究所コーパス開発センター)  
保田 祥 (国立国語研究所コーパス開発センター)  
小西 光 (国立国語研究所コーパス開発センター)  
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)

### Early Results of Page Collection and Linguistics Annotation in Ultra Large Scale Web Corpus Construction

Masayuki ASAHARA (Center for Corpus Development, NINJAL)  
Mizuho IMADA (Center for Corpus Development, NINJAL)  
Sachi YASUDA (Center for Corpus Development, NINJAL)  
Hikari KONISHI (Center for Corpus Development, NINJAL)  
Kikuo MAEKAWA (Dept. of Corpus Studies, Center for Corpus Development, NINJAL)

#### 1. はじめに

国立国語研究所では 2006~2010 年度の期間に 1 億語規模の書き言葉コーパス『現代日本語書き言葉均衡コーパス』(以下“BCCWJ”) (前川 2007, 前川・山崎 2008) を構築し, 2011 年より一般公開した。BCCWJ は種々の母集団に沿った無作為抽出を実施することによって, 高度な均衡性・代表性を備えた均衡コーパスとなっている。しかし, その規模は, 現代のコーパス言語学の趨勢からすれば十分とはいいがたく, 生起頻度の低い言語現象の被覆に問題がある。そのためより大規模な日本語コーパスの構築が望まれている。この問題を解消するため, 国立国語研究所コーパス開発センターでは 2011 年度から 6 年の期間で, Web を母集団とした 100 億語規模の超大規模コーパスを構築する計画に着手した。

現在, 6 年のプロジェクトの 3 年が過ぎ, 収集と組織化について進捗している。またプロジェクトの初期において収集と組織化のための環境構築に工数をかけた。現時点の環境整備・収集・組織化についてはある程度目途がつきつつある。2012 年 10 月より本収集を開始し 2013 年 9 月まで 1 年間収集したデータについて組織化が進んでおり, 基礎統計量が得られている。

本稿ではプロジェクトの進捗状況を過去 1 年間に収集されたデータの基礎統計量とともに示し, 本言語資源を利用した言語研究の可能性について論じる。2 節では収集の進捗状況について述べる。3 節では組織化の進捗状況について述べる。4 節では本言語資源と Google による「Web 日本語 N グラム第 1 版」(1) との頻度上位 10 語の比較を定性的に分析する。5 節で言語資源の大規模化による理論的・応用的研究の可能性について論じ, 6 節にまとめと今後の予定について述べる。

#### 2. 収集の進捗状況

本節では収集の経過と現在までに得られている統計量について示す。

2012 年 6 月に収集のための環境整備が完了し, 環境整備に関して所内のセキュリティ委員会に諮り, クローラ運用の許諾を得た。許諾後すぐにクローラの試験運用 1 か月前よりクローラに関する情報提供・問い合わせ窓口としての Web ページ/メールアドレス/電話を設置した。1 か月の周知期間を経て, 2012 年 7~8 月より試験収集を開始し, Heritrix (2) の各種パラメータを調整した。9 月に最終試験収集となる 1000 万 URL 規模の収集を行い,

その際に得られた収集速度の情報からクローラ運用方針を年間4回収集, 1回あたりの収集量は1億URL規模で3か月間に設定した。URLは1年間通して4回収集し, 季節に偏らないように配慮する一方, URLの更新頻度およびリンク先の情報に基づき, 収集するURLを1年毎に変更する方針にした。2012年10月より3か月ごとに収集を行い, 2013年12月末現在, 5回目の収集の最中で次の1年間に収集すべきURLサンプルを決定する状況にある。

以下収集の統計量について示す。全体については2012年第4四半期(2012-4Q)~2013年第3四半期(2013-3Q)についての統計量を, 各論については2012年第4四半期のみの統計量を示す。

表1: 2012年第4四半期から2013年第3四半期の収集ページ数

	2012-4Q	2013-1Q	2013-2Q	2013-3Q
ページ数 (1期)	61,668,805	58,844,092	61,479,268	57,892,917
内容の重複なしページ数	45,933,605	42,932,982	45,111,527	42,192,931
<b>4期通しての統計</b>				
総異なり URL 数 (4期)	64,539,233			
(内) 内容の重複ありページ数	27,604,915			
(内) 内容の重複なしページ数	36,934,706			

表1に収集したページ数の統計量を示す。1億URLを収集してもrobots.txtの順守や各種HTTPエラーにより, ページとして収集できたものが約六割にすぎない。重複検出はURL毎に各ページのハッシュ値を計算し同一性を認定する。各期において内容の重複なし(異なり)ページ数は4000万強になる。4期通しての総異なりURL数は約6400万URLと1億URLに至らない。4期中2期以上収集できたページ数の内, 内容の重複があるページ数は約四割の2700万ページ, 反対に内容の重複がないページ数は3600万ページになる。

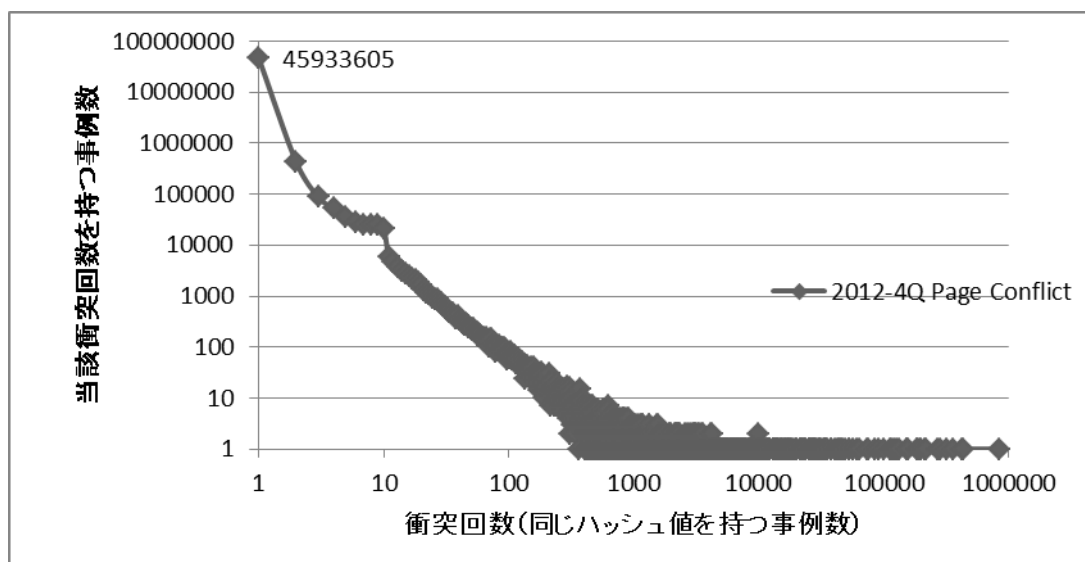


図1: 2012年第4四半期収集ページの重複

図1に2012年第4四半期収集Webページの重複検出結果について示す。同じハッシュ値を持つURLが複数存在することを「衝突」と呼ぶ。グラフ横軸は同じハッシュ値を持つURL数を示し, 「衝突回数」と呼ぶ。グラフ縦軸は横軸の「衝突回数」を持つ衝突事例数(URL数ではなくハッシュ値の異なり数で計算)を示す。グラフは両軸とも対数で表示している。グラフ中の左上の点が表の「内容の重複なしページ数」(他のURLと内容が重複しないページ数)に相当する。衝突回数10以下のものは同一内容の異なるURL表示もしくはいわゆるコピーサイトであると考えられる。それ以上の衝突については, robots.txtや, 「ソフト404」

と呼ばれる当該 URL はサーバ上にないということを 404 HTTP ステータスコードでは返さず 200 HTTP ステータスコードで返し当該ページがないことを示すコンテンツを返すページである。

表 2 : 2012 年第 4 四半期から 2013 年第 3 四半期の収集リンク数

	2012-4Q	2013-1Q	2013-2Q	2013-3Q
リンク先 (のべ)	6,905,805,383	6,610,763,700	7,064,611,259	7,222,958,033
	69 億 URL	66 億 URL	70 億 URL	72 億 URL
リンク先 (異なり)	892,135,930	843,166,672	865,694,816	855,684,918
	8.9 億 URL	8.4 億 URL	8.6 億 URL	8.5 億 URL
<b>4 期通しての統計</b>				
リンク先 (異なり)	1,642,699,579			
	16 億 URL			

表 2 に 2012 年第 4 四半期 (2012-4Q) ~2013 年第 3 四半期 (2013-3Q) の収集リンク数を示す。おおよそ 6000 万 URL の収集に対し, のべ 70 億前後, 異なり 9 億弱の URL が収集できている。4 期を通した集計によるリンク先数が異なり 16 億 URL であることから 1 年間通して同じ URL を 4 期収集することにより 1 期のみクロールするのみ比べてリンクが約 1.8 倍 (8.5 億~8.9 億→16 億 URL) に成長していることがわかる。

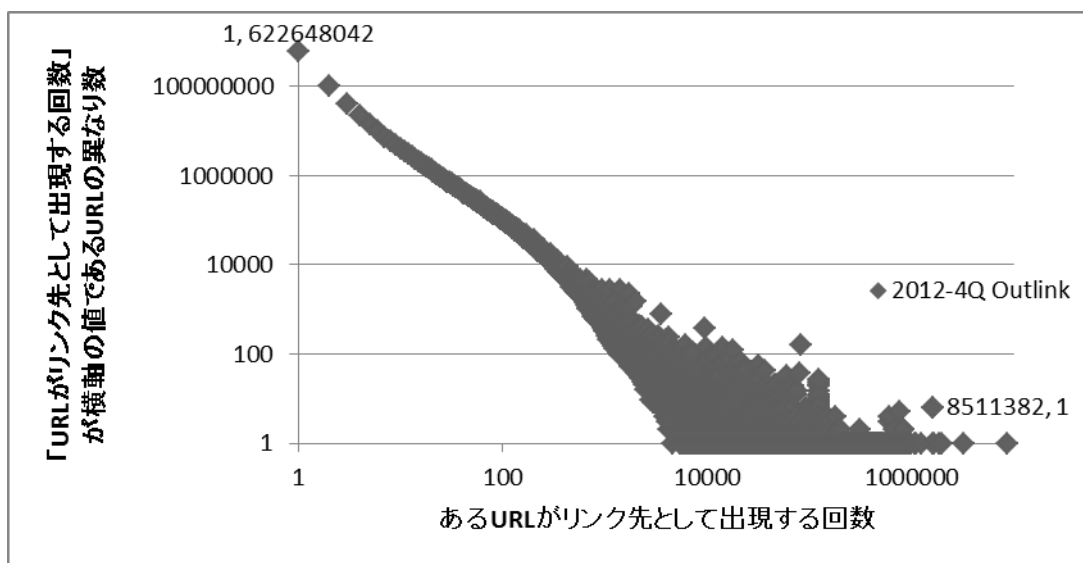


図 2 : 2012 年第 4 四半期収集リンク先の統計

図 2 に 2012 年第 4 四半期の収集リンク先統計を示す。グラフ横軸はある URL がリンク先として出現する回数 (被リンク数) で, グラフ縦軸が横軸相当の被リンク数を持つ URL の異なり数である。グラフの両軸は対数である。グラフ左上が 1 回しかリンクされていない URL 数で約 6 億件である。グラフ右下が最も多い被リンク数約 850 万を持つページで 1 件ある。これは有名ブログサイトのトップページであり, ブログの個々のページからリンクされている。

これらの統計情報から 2014 年以降の収集対象 URL を決定する。収集対象 URL は, 4 期にクロールした同一 URL のうち内容が毎回変わっていた URL と, 収集した Web ページのリンク先 URL の二種類を想定している。現在までに収集した Web ページに対するレジスタ分析は進んでいないが, レジスタ分析が進み次第, レジスタ分析結果の分散を見ながら収集対象 URL を決定したい。

### 3. 組織化の進捗状況

本節では組織化の経過と現在までに得られている統計量について示す。現在までのところ、正規化・形態素解析・係り受け解析までが部分的に進捗している。正規化においては nwc-toolkit による文字コード統制・文抽出を行った。利用するライブラリの関係でサーバ群ではなく安価な 8-core のワークステーション上で正規化作業を行った。8-core の CPU の並列処理により約 1 週間で 1 年分の収集データの正規化作業が完了する。正規化されたデータは MeCab/IPADIC, MeCab/UniDic, JUMAN により形態素解析を行う。32-core の計算サーバ上の並列処理により、それぞれ 1 日弱で 1 年分の収集データの形態素解析作業が完了する。さらに IPA 品詞体系に基づく CaboCha と益岡・田窪品詞体系に基づく KNP により係り受け解析を行う。32-core の計算サーバ上の並列処理により 1~2 週間で係り受け解析作業が完了する。

表 3 に組織化したデータの基礎統計量を示す。Heritrix は収集 Web ページを圧縮 1GB サイズの WARC データに分割して出力する。展開すると約 3 倍程度になるため、表中の収集 WARC ファイル数に 3GB をかけた値が収集 Web ページ容量 (メタデータを含む) と概算することができる。URL 数は前節の収集における URL 数である。正規化処理は nwc-toolkit による。正規化処理の際に文抽出なしに形態素解析 (MeCab/IPADIC) を行うと各期のべ約 620~647 億形態素になる。文抽出を行うと形態素数は各期約 300 億強になることから大体半分の形態素が日本語の文中の形態素ではないとして排除されている。抽出された文数はのべ数で各期 25 億文前後、文単位の同一性を認定すると文の異なり数は各期 10 億文になる。

表 3 : 組織化したデータの基礎統計量

	2012-4Q	2013-1Q	2013-2Q	2013-3Q
収集 WARC ファイル	814	870	910	905
URL 数	61,668,805	58,844,092	61,479,268	57,892,917
形態素数 (文抽出なし)	64,714,650,129 647 億形態素	62,077,520,745 620 億形態素	63,414,252,638 634 億形態素	65,736,027,334 647 億形態素
形態素数 (文抽出あり)	33,767,409,441 337 億形態素	32,651,138,004 326 億形態素	33,073,991,355 330 億形態素	30,923,912,566 309 億形態素
文数 (のべ数)	2,678,315,774 26 億文	2,600,122,908 26 億文	2,659,617,620 26 億文	2,478,309,312 24 億文
文数 (異なり数)	1,097,011,506 10 億文	1,048,772,913 10 億文	1,063,649,324 10 億文	1,007,771,383 10 億文

図 3 に 2012 年第 4 四半期の収集文の重複を示す。横軸が同一文の出現回数で縦軸が当該出現回数の文の異なり数を表す。両軸とも対数で表現する。10 億文のうち約 9 割の 8.9 億文が 1 回しか出現しない文である。以下、ページの重複も含めて同一文が異なりで 1 億文規模存在する。これらは定型的な表現やリンクの見出し語であることが多く、一番多く出現した文は 2,885,654 回出現する「職業とキャリア」(Yahoo!知恵袋のカテゴリ名)であった。ここで文抽出処理の妥当性について検証する。検証データとして BCCWJ に nwc-toolkit を適用して文数とバイト数がどの程度変化するかについて示す。

表 5 に nwc-toolkit のうち文抽出を行う nwc-toolkit-text-filter を BCCWJ に対して適用した場合のジャンル毎の文数とバイト数の変化を示す。表中左の文数はそれぞれの文境界基準により認定しており、OC (Yahoo!知恵袋) のように句点・感嘆符・疑問符が文境界として利用されない場合が多いジャンルについては nwc-toolkit により文数が増える場合があり、単純にどの程度削減されるかは評価できない。表中右のバイト数が純粋にデータとして削

減されたテキスト量を表す。バイト数に関して変化率が低い(多く削除されている)ジャンルは OP (広報紙) と PM (雑誌) であった。これらに対し、何が削除されているのかを確認したところ、人口などの統計情報・イベントの日時・窓口やお店の営業時間・連絡先(電話番号・住所)などが削除されていることがわかった。情報抽出の分野においては重要な情報ではあるが、言語研究においては削除しても多くの研究で影響は少ないと考える。

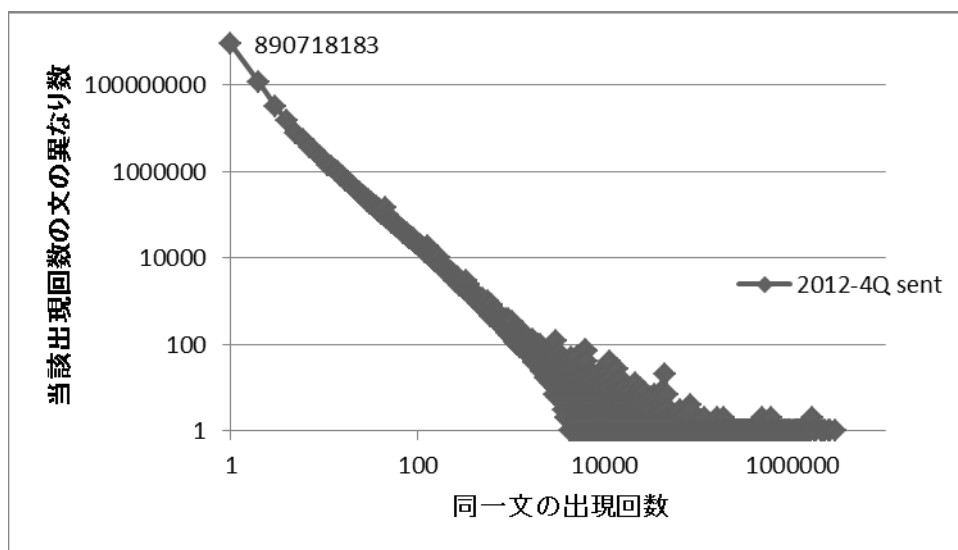


図3 : 2012年 第4 四半期収集文の同一性

表4 : BCCWJ の nwc-toolkit-text-filter の適用による変化

ジャンル名	文数			バイト数		
	処理前	処理後	変化率	処理前	処理後	変化率
LB	451,273	394,782	0.87	186,908	176,792	0.95
OB	222,437	203,467	0.91	23,236	22,242	0.96
OC	42,506	45,082	1.06	369,004	368,572	1.00
OL	38,827	34,761	0.90	6,004	5,864	0.98
OM	140,422	116,295	0.83	26,948	26,052	0.97
OP	257,796	142,660	0.55	22,372	14,636	0.65
OT	64,100	47,680	0.74	5,992	5,340	0.89
OV	18,977	17,807	0.94	1,656	1,632	0.99
OW	146,402	126,140	0.86	28,072	24,876	0.89
OY	117,816	76,205	0.65	224,064	213,480	0.95
PB	392,301	308,827	0.79	176,812	163,476	0.92
PM	301,399	227,089	0.75	30,248	25,812	0.85
PN	80,563	65,055	0.81	10,164	9,152	0.90
合計	2,274,819	1,805,850	0.79	1,111,480	1,057,926	0.95

現在収集しているコーパスデータについて語彙表の作成や n-gram データを進めている。unix コマンド sort, uniq, wc などに基づく単純な集計手法では sort アルゴリズムや途中で出力される中間ファイルにより時間計算量・空間計算量が想定よりも大きく時間がかかることがわかった。現在様々なソフトウェアを組み合わせることで効率化をはかっており、2012 年第 4 四半期の n-gram 統計のみ得つつある。

## 4. 本言語資源とグーグル「Web 日本語 N グラム」との比較

最後に 2012 年第 4 四半期収集データとグーグル「Web 日本語 N グラム」との比較を行う。

表 5 : 2012 年第 4 四半期収集データとグーグル「Web 日本語 N グラム」との比較

	本言語資源(2012-4Q) MeCab/IPADIC 頻度 3 以上 (文抽出有) 文単位での重複性排除有	本言語資源(2012-4Q) MeCab/IPADIC 頻度 3 以上 (文抽出有) 文単位での重複性排除無	「Web 日本語 N グラム」 MeCab/IPADIC 頻度 20 以上 [工藤・賀沢 2007]
総形態素数 (のべ)	18,075,529,620 180 億形態素	33,767,409,441 337 億形態素	255,198,240,937 2550 億形態素
総文数	1,097,011,506 (異なり) 10 億文	2,678,315,774 (のべ) 26 億文	20,036,793,177 200 億文
1-gram	0.039 億	0.050 億	0.025 億
2-gram	0.47 億	0.85 億	0.80 億
3-gram	1.6 億	4.4 億	3.9 億
4-gram	2.1 億	8.7 億	7.0 億
5-gram	1.7 億	10.3 億	7.7 億
6-gram	1.2 億	9.7 億	6.8 億
7-gram	0.84 億	8.5 億	5.7 億

表 5 に 2012 年第 4 四半期収集データと 2007 年公開のグーグル「Web 日本語 N グラム」の比較結果を示す。「Web 日本語 N グラム」では頻度 20 以上の n-gram データのみを配布している。本研究ではできるだけ低頻度のデータを被覆するべく適切な頻度を検討中であるが、速報値として頻度 3 以上の概算値を報告する。総文数においては「Web 日本語 N グラム」の規模の 20 分の 1 から 10 分の 1 くらいの規模である。しかし低頻度のものも組織化することにより n-gram データの被覆では遜色ないレベルにできると考える。

表 6・表 7 は 2012 年第 4 四半期収集データとグーグル「Web 日本語 N グラム」の頻度順位上位 10 件を比較するものである。頻度順位において、文境界メタ記号や句読点・記号類を含む n-gram は排除している。これらより各 n-gram の定性的な分析を行う。

まず本言語資源と「Web 日本語 N グラム」とで標本の量を比較すると 10 倍 (文単位での重複性排除無) ~20 倍 (文単位での重複性排除有) くらい「Web 日本語 N グラム」の方が規模が大きい。単純に収集規模が小さい本言語資源において文単位の重複性排除無で n-gram を取ると 3-gram や 4-gram のレベルで複数のページで出現する定型的な表現が頻度順位上位にきてしまう。2-gram のレベルでも「利用規約」という定型的な表現が確認された。一方、文単位での重複性排除有の本言語資源では、規模が小さくても「Web 日本語 N グラム」と同様に機能語中心の n-gram が得られる。このように文単位での重複性排除を行うことにより Web 特有の定型的な表現の頻度の偏りをなくすことができ、定型的な表現の影響が少なくなる、より規模の大きな言語資源と同じ傾向のデータが得られていることがわかる。

5-gram 以上になると「Web 日本語 N グラム」においても定型的な表現の影響が避けられず、多くみられるようになる。一方、文単位での重複性を排除したものについては、文の一部を変えるような定型的な表現(「タグが付けられた質問」「に関するウェブ上の情報を探す」など)は見られるが、基本的には機能表現中心のコロケーションが得られていることがわかる。

特殊な事例として一文中に同じ文字が連続している表現がある。「ああ ああ ああ ああ

あああああ」や「えええええええ」などは、一文中に同じ文字が連続している表現の形態素解析結果である。例えば「え」が p 回出現するような一文に対して、p-6 回の「えええええええ」7-gram が枚挙される。一文中に複数の n-gram が出現する際のように数え上げるかは議論の余地がある。これは構造をもつデータに対する問い合わせにおいて起きる常に問題で、現在は多くのコンコーダンスにおいて正規化は行っていない。

表 6 : 2012 年第 4 四半期収集データとグーグル「Web 日本語 N グラム」との比較  
頻度順位上位 10 件 (1-gram~4-gram)

	順位	1-gram	2-gram	3-gram	4-gram
本言語資源 MeCab/IPADIC  2012-4Q 文単位での重複性排除有	1	の	して	ています	しています
	2	に	ました	ていた	していました
	3	て	てい	してい	されている
	4	が	ている	している	していた
	5	は	した	と思います	されてい
	6	を	では	されて	たのですが
	7	た	には	になって	てきました
	8	で	され	のですが	れています
	9	と	ません	しました	はありません
	10	し	います	された	になりました
本言語資源 MeCab/IPADIC  2012-4Q 文単位での重複性排除無	1	の	ました	記事への	記事へのトラック
	2	に	でしょう	お願いします	専用ページを表示
	3	を	行って	Q&A	利用することが
	4	は	思って	続きを読む	機能を利用する
	5	て	情報を	マークへ投稿	おすすめの知恵ノート
	6	が	利用規約	専用ページを	正確性の保証
	7	た	おすすめの	機能を利用	お客様自身の責任
	8	で	記事へ	済みの質問	回答を指示する
	9	と	追加する	おすすめの知恵	便利に新規取得
	10	し	場合は	エンターテインメントと趣味	はてなブックマークへ
「Web 日本語 N グラム」 MeCab/IPADIC  [工藤・賀沢 2007]	1	の	して	ています	しています
	2	に	ました	してい	されている
	3	を	てい	ていた	されてい
	4	は	ている	している	はありません
	5	て	した	されて	れています
	6	が	ません	になって	ていました
	7	た	され	しました	になりました
	8	で	には	された	しております
	9	と	では	れている	てきました
	10	し	います	ありません	していた

例えば、「N1 が N2 が V1 する N3 を V2」といった文に対して「『N が』が先行する『V』」といった問い合わせを行う際に何件出現すると言えるだろうか？多くのコンコーダンスは <N1 が,V1>, <N1 が,V2>,<N2 が,V1>,<N2 が V2>の全組み合わせを出力し 4 件と出力するだろう。

このように n-gram を含めて 2 語以上の組み合わせを行う際には、重複して枚挙される現象に注意する必要がある。さらに、係り受け構造など構造を持ったデータに対する部分構造の枚挙に際しては、同様な事例が多く出現する。このような部分構造を共有する事例に

について、n-gram データとして正規化することは一切行わないが、利用者においては統計処理を行う場合に必要に応じて対処する必要がある。

本 n-gram データについては近いうちにデータそのものを公開するとともに、簡単な検索環境を Web API つきで対外的に公開し利用可能なようにする。2014 年度以降進めていく利用者系に関する意見などをいただければ幸いである。

表 7 : 2012 年第 4 四半期収集データとグーグル「Web 日本語 N グラム」との比較  
頻度順位上位 10 件 (5-gram~7-gram)

	順位	5-gram	6-gram	7-gram
本言語資源 MeCab/IPADIC 2012-4Q 文単位での重複性排除有	1	されています	ではないでしょうか	のではないのでしょうか
	2	ではありません	ていたのですが	のタグが付けられた質問
	3	と思っています	のではないでしょう	ではないかと思います
	4	していました	のではないかと	に関するウェブ上の情報を探す
	5	ではないでしょう	に行ってきました	ああああああああああああ
	6	のではないか	ような気がします	のではないかと思い
	7	はないのでしょうか	タグが付けられた質問	していたのですが
	8	になっています	のタグが付けられた	思っていたのですが
	9	ていましたが	させていただきました	えええええええ
	10	ていたのです	たいと思っています	と思っていたのです
本言語資源 MeCab/IPADIC 2012-4Q 文単位での重複性排除無	1	記事へのトラックバック	機能を利用することが	機能を利用することができ
	2	機能を利用すること	利用することができませ	利用することができません
	3	利用することができ	正確性を保証して	正確性を保証しており
	4	正確性を保証し	お客様自身の責任と判断	お客様自身の責任と判断で
	5	お客様自身の責任と	すべての機能を利用する	すべての機能を利用すること
	6	はてなブックマークへ投稿	知恵袋のすべての機能を	知恵袋のすべての機能を利用
	7	更新情報が届きます	おすすめの解決済みの質問	ニックネームの My 知恵袋で確認でき
	8	おすすめの解決済みの	記事へのトラックバック URL	質問年月や画像の有無を
	9	すべての機能を利用	ニックネームの My 知恵袋で確認	質問や知恵ノートは選択さ
	10	質問年月や画像の	することができません	以上更新がないブログに表示
「Web 日本語 N グラム」 MeCab/IPADIC [工藤・賀沢 2007]	1	されています	無料でお届けします	料無料でお届けします
	2	ではありません	料無料でお届けし	配送料無料でお届けし
	3	でお届けします	配送料無料でお届け	国内配送料無料でお届け
	4	無料でお届けし	国内配送料無料でお	以上国内配送料無料でお
	5	1500 円以上国内配送	円以上国内配送料無料	円以上国内配送料無料で
	6	料無料でお届け	以上国内配送料無料で	1500 円以上国内配送料無料
	7	配送料無料でお	1500 円以上国内配送料	はインラインフレームを使用して
	8	国内配送料無料で	を使用しています	フレームを使用しています
	9	以上国内配送料無料	インラインフレームを使用して	インラインフレームを使用してい
	10	円以上国内配送料	この記事へのトラックバック	部分はインラインフレームを使用し

### 5. 本言語資源を利用した理論的・応用的研究の可能性

本節では本言語資源を利用した理論的・応用的研究の展望について述べる。

まず語彙研究のために Web を母集団とするコーパスに対してどのような統計処理をすればよいかについて述べる。コーパスに基づく語彙研究は大きく分けて、生コーパスから得られる統計情報のみから語彙の性質を明らかにする教師なし統計学習に基づく手法と、何らかのアノテーションを付与してその情報を未知データに再現するような教師あり統計学習に基づく手法の二つが考えられる。



前者の教師なし統計学習においては、近年ベイズ統計学が注目されトピックモデルに基づく手法が盛んに利用されている。コーパスを用いてトピックモデルを構築する際に「文書」という単位と「単語」という単位が重要になる。本言語資源においては「文書」が個々の Web ページ、「単語」が形態素解析により認定された形態素が対応する単位になる。文書内の単語の共起などにに基づき潜在変数を推定する。単語の与え方によっては、トピックを文書に遍在するレジスタとみなしレジスタ分析にも利用することができる。ここで本研究では「文」という単位を重要視する。Web の場合、コピーサイトや複数のページで共有する部分ページなどがあり、単純に「文書×単語関係」のみにより表現するとこれらの影響を受けてしまう。3節に示した統計量のように約 1 割の 1 億文については複数の文書に共有されている。そこで文書と単語の間に文という単位を仮定し、「文書×文関係」「文×単語関係」と階層的にモデル化することにより、コピーサイトなどの影響を前者の「文書×文関係」で吸収しながらレジスタをモデル化することが可能になると考える。さらに Web の場合には「文書×文書間関係」としてリンカー被リンク関係が規定されている。単リンクが張られる関係、もしくは同一文書にリンクを張る関係など、レジスタに影響を与える関係がいろいろ考えられる。この「文書×文書間関係」「文書×文関係」「文×単語関係」の多層の共起関係を、現実的に処理可能な計算量でどのようにモデル化するかを検討する。

後者の教師あり統計学習においては、古くは二値分類器が中心に用いられたが、最近では、順位・極性・系列などといった構造学習が自然言語処理の分野で広く用いられている。また特徴量としても木構造や有向非循環グラフなどといったものの部分構造の重複などを効率的に枚挙しながら距離(類似度)を設定して識別する学習手法が多く提案されている。本言語資源では形態論情報だけでなく、係り受け構造や述語項構造などといった情報を付与する。これにより言語研究の分野における構造の利用(特徴量としての利用と推定する対象としての利用)が促進されればと考えている。

利用者系において、3節で通常の n-gram データと文単位での重複を排除した n-gram の二つについての統計量を提示した。n-gram 検索ツールについては、この双方を提供して利用者による評価を待ちたい。また KWIC 検索においては後者の文単位での重複を排除したコーパスに対する検索ツールの開発のみを検討している。3節や4節に示すような様々な情報を提供しながら文単位での重複性の影響について議論していきたい。

形態論研究においては、Web 上の多様な語彙を観察することが可能になるだけでなく、保存する予定であるデータを通時的に観察することで経年変化を分析することが可能になる。Web の世界においては多様な形態論変化が生産され、すぐに廃れていく傾向にある。これらを適切にとらえることができるような環境を構築できればと考えている。

コーパス日本語学において、統語論研究に必要な情報が現状のコーパスから得られているとは言い難い状況にある。まず BCCWJ では品詞情報や係り受け構造に基づく問い合わせをおこなったとしても真に分析したい言語現象が数例しかないということが起きうる。例えばガ格が二つ以上出現する文(複文)の従属節境界の左端のあいまい性(pre-head attachment; Kamide and Mitchel 1999)を研究する際に、BCCWJ に様々な構造を問い合わせると「N ガ N ヲ VP N ガ」という語順が 1 例、「N ガ ADV N ガ」という語順が 3 例、「N ガ N ニ N ガ」という語順が 1 例しか出現しないことがわかる。規模を 100 倍にすることでこういった言語現象を発見する可能性は高くなる。また人文系の研究者側の問題として、多くの研究が語彙情報を単純なベクトル表示にする bag-of-words 的な共起以上の情報を使いこなせていないという問題がある。自動解析ながらもあらかじめ係り受け構造等を付与し、使いやすい利用者系を構築することで、語順や部分木構造などを考慮したコーパス分析手法が根付くような土壌を整備していきたい。

意味論研究においても、頻度情報によらない研究にまで踏み込むことが重要であると考える。この分野においては、共起に基づく研究は盛んに行われている一方、アノテーションを行い、さらに構造学習を導入してアノテーションを未知データに復元することで何かを明らかにするような水準に達していないと考える。ここで重要なのは、解明しようとし

ている研究課題が本質的にコーパスを用いる手法が適しているかどうかを見極める必要があるということである。コーパスを用いる手法が何も証拠を提供できない問題は、言語研究において多々あると考える。そういった問題に対しても、コーパスが仮説を立てるための手がかりを研究者に与えることは可能な場合が多いが、手がかりのみをもって証拠とすることのないように注意されたい。

## 6. おわりに

本論文では国立国語研究所コーパス開発センターで進めている Web を母集団とした超大規模コーパス開発プロジェクトの進捗について報告した。6年間のプロジェクトのうち半分の3年が経ち、単純に100億語規模のコーパスを構築するだけでなく、継続的に同規模のスナップショット的なコーパスを1年間に複数回構築可能である環境が構築されつつある。収集においては Heritrix クローラを利用して年間4回のクローラを行い、組織化においては正規化・形態素解析・係り受け解析まで進捗している。組織化における残る課題は述語項構造解析とレジスタ分析と語彙表構築であるが、これについても早急に環境を整えたい。得られた基礎統計量から3か月単位で目標の100億語規模のコーパスが構築できていることがわかった。今後より一層の効率化を進めたい。

残りの期間で人文系の研究者が柔軟に利用可能な利用者系と保存環境の構築を行う。語彙調査や用例検索に留まらない、自然言語処理で培われた構造に基づく問い合わせ環境を構築したい。これにより高い被覆性を持ちながらも柔軟なコーパス調査を可能にし、統語論・意味論研究を前に進める研究環境を提供できると考える。

一方で、本言語資源に基づく調査で解明できない問題が何であるのかを示すが重要だと考える。規模を大きくすることだけでは解明できない問題についても示していきたい。

## 謝 辞

本研究は国語研「超大規模コーパス構築プロジェクト」によるものです。

## 文 献

- Kamide, Yuki and D. C. Mitchell (1999) Incremental Pre-head Attachment in Japanese Parsing, *Language and Cognitive Processes*, 14(5/6), 631-662.
- 工藤拓・賀沢秀人 (2007) 『Web 日本語 N グラム第1版』, 言語資源協会発行 (<http://www.gsk.or.jp/catalog/GSK2007-C/>).
- 前川喜久雄 (2007) 「コーパス日本語学の可能性-大規模均衡コーパスがもたらすもの-」 *日本語科学*, 22, 13-28.
- 前川喜久雄・山崎誠 (2008) 「『現代日本語書き言葉均衡コーパス』」 *国文学解釈と鑑賞*, 932 (74 巻1号), 15-25.
- 矢田晋 (2010). 『日本語ウェブコーパス 2010 (NWC 2010)』, (<http://s-yata.jp/corpus/>).

## 関連 URL

- (1) Web 日本語 N グラム第1版 : (<http://www.gsk.or.jp/catalog/GSK2007-C/>).
- (2) クローラ Heritrix-3.1.1 :  
(<http://webarchive.jira.com/wiki/display/Heritrix/Heritrix>).
- (3) 日本語ウェブコーパス用ツールキット : (<http://code.google.com/p/nwc-toolkit/>).
- (4) 形態素解析器 MeCab-0.996 :  
(<http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html>).