

## 軽量・高速な日本語解析ツール『簡易日本語解析系Q\_JP』

亀田 雅之

(株)リコー・研究開発本部・情報通信研究所

kameda@ic.rdc.ricoh.co.jp

### 概要

従来の大規模な日本語解析系に対し、多くの応用系に手軽に利用できる実用レベルの軽量・高速な日本語解析ツール『簡易日本語解析系Q\_JP』を開発した。

Q\_JPは、機能語と字種に着目した小規模辞書ベースの形態素解析系と、品詞情報ベースのヒューリスティック規則により係り受け選択を行う構文解析系からなる。1. 軽量性 [辞書: 100KB 足らず, メモリ: 約 250KB (DOS 版)], 2. 高速性 [PC (80486/25MHz): 80~150 語/秒, WS (Sun-SS20): 700~800 語/秒], 3. 頑強性 [長文解析, 未知語登録不要] を備え、パソコンでも余裕をもって動作する。

Q\_JPの機能はライブラリ化し、現在、文書検索、読解支援、自然言語 I/F 等への応用を図っている。

### 1 はじめに

日本語解析技術は、今後の高度な日本語文書処理のキーとなる重要な技術である。

しかし、日本語解析のためには、数万語以上の大規模な辞書や曖昧さや意味を扱う重い処理系が必要である。また、未知語やその登録作業、曖昧さの組合せ爆発のために長い文が扱えない、といった問題がある。

こうした問題を踏まえ、多くの応用系に手軽に利用できる実用レベルの日本語解析系を目指し、AWK言語上で試作した『簡易日本語解析系Q\_JP』を報告した [1]。試作システムでは、小規模辞書で形態素解析から構文解析 (係り受け解析) まで行ない、かつ、十分な解析性能と頑強性をもつことを示した。その後、さらにC言語で実装することにより、少ないメモリで極めて高速に解析する系を実現し、所期の目的であった応用系への実用レベルでの展開を可能とした。

Q\_JPでは、日本語文における字種の特徴に着目し、未知語の問題を避けながら、機能語主体に辞書を小規模化するとともに、構文解析での曖昧さの扱いを簡易化し、意味に踏み込まない品詞情報レベルで優先選択することで、処理を軽くしている。

日本語解析系としては、形態素解析系 JUMAN [2] や構文解析系 KNパーサ [3] 等が公開されており、また、各所で汎用化を図っていると思われるが、Q\_JPのように応用系への組み込みを指向し、軽量性、高速性、頑強性を備えた日本語解析系の提示はない。

本稿では、Q\_JPの解析方式を概観し、現状のC言語版の諸元について示す。

### 2 Q\_JPの解析方式

Q\_JPは、漢字かな混じり日本語文を形態素列、単語列に分割する形態素解析系、さらに係り受け構造に解析する構文解析系からなり、各々次の方針を採った。

- 字種の特徴を利用した文節分割 [4] や予備的な単語分割 [5] の方式を、小規模辞書ベースの形態素解析系として拡張・精緻化する
- 構造の曖昧さを簡易的に扱うとともに、小規模辞書を生かすため、意味情報に依らず品詞/構文レベルの情報に基づく構文解析系とする

#### 2.1 形態素解析系

形態素解析系では、品詞ごとの表記の字種 (漢字/ひらがな/カタカナ等)、語彙の数、接続上の制約の各特徴、及び和語系複合語等の形態素構成の特徴に着目し、形態素候補を次の辞書、字種・品詞割り当て、単語合成の分担で扱い、辞書を小規模化する。

辞書: 辞書には、語彙的に閉じた系をなす機能語 (助詞, 助動詞, 接続詞, 連体詞, 副詞等) や活用語尾等を中心に、少数品詞語、次の字種・品詞割り当てから外れる例外語等を登録する。

字種・品詞割り当て: 辞書引きできなかった同一文字種の文字列に対し、字種と文字列長に応じて、語彙数の多い概念語品詞 (名詞, サ変名詞や動詞, 形容動詞の語幹) を候補として割り当てる。語彙数の多いこれらの品詞の単語を字種・品詞割り当てで扱うことで、登録単語数を大幅に削減する。

単語合成: 日本語では、和語系の複合語(「組み-込む」, 「切れ-味」)、派生語(「書-か-す」)、品詞転性(「歩-み」)等の造語力が顕著であり、辞書エントリを増やす要因でもある。これらに対して、短単位の形態素解析と、活用語の語幹と語尾の合成を行う単語合成の枠組みの中で、複合語の形態素の並びを検出し、単語として認識することで、辞書への登録を不要とする。

図1に、短単位の形態素解析結果と単語合成の例を示す。図中の「日本語」他の漢字表記の名詞と「切(る)」/「置(く)」の漢字の動詞語幹、和語複合語「切れ目」は辞書にない。

上記の字種-品詞割り当ては仮想辞書引き処理であり、形態素解析手法は、最長一致と品詞接続検査による方法と変わらない。尚、接続検査で品詞の曖昧さが未解消の形態素に対しては、前後の品詞や単語を検査して、曖昧さを解消したり、優先順位を変えるヒューリスティックな品詞の曖昧さ解消規則がある。

辞書は4種類ある(表1)が、通常は、拡張辞書(非ひらがなの機能語、例外語等)とひらがな辞書(ひらがな機能語、活用語尾、主要なひらがな単語)を参照する。ひらがな拡張辞書は、ひらがな文字列部分の解析失敗でのバックトラック時に、形容動詞語幹辞書は、形容動詞を含む品詞の曖昧さ解消時に参照される。

【形態素リスト】		【単語リスト】	
• [1] (8) 日本語	[J] (523) 名詞	• [1] (8) 日本語 (日本語)	(41) 名詞
• [2] (14) の	[zh] (37) ノ-格助	• [2] (14) の (の)	(51) ノ-格助
• [3] (16) よう	[zh] (283) ヨウダゲ-格助	• [3] (16) ように <ようた>	(60, 22) ヨウダゲ-格助
• [4] (20) に	[zh] (250) ダ-格助		
• [5] (22) 単語	[J] (523) 名詞	• [4] (22) 単語 (単語)	(41) 名詞
• [6] (28) 間	[X] (356) 形式名詞	• [5] (28) 間 (間)	(43) 形式名詞
• [7] (28) に	[zh] (36) ニ-格助	• [6] (28) に (に)	(51) ニ-格助
• [8] (30) 切	[J] (349) 下ラ基	• [7] (30) 切れ+目 (切れ目)	(41) 名詞
• [9] (32) れ	[zh] (351) 下ラ用	• [8] (32) れ (を)	(51) テ-格助
• [10] (34) 目	[X] (507) 単独名詞	• [9] (34) 置+か <置く>	(15, 11) 動: 五カ来
• [11] (36) を	[zh] (41) テ-格助	• [10] (42) な+い <ない>	(60, 4) ナイ-格助
• [12] (38) 置	[J] (369) 五カ幹	• [11] (46) 膠着言語 (膠着言語)	(41) 名詞
• [13] (40) か	[zh] (372) 五カ来	• [12] (54) の (の)	(51) ノ-格助
• [14] (42) な	[zh] (256) ナイ-格助	• [13] (56) 文 (文)	(41) 名詞
• [15] (44) い	[zh] (205) ク-格助	• [14] (58) の (の)	(51) ノ-格助
• [16] (46) 膠着言語	[J] (523) 名詞	• [15] (60) 地理 (地理)	(41) 名詞
• [17] (54) の	[zh] (37) ノ-格助	• [16] (64) に (に)	(51) ニ-格助
• [18] (56) 文	[J] (523) 名詞	• [17] (66) お+い <おく>	(75, 22) オク-五カ来
• [19] (58) の	[zh] (37) ノ-格助	• [18] (70) て (て)	(55) テ-格助
• [20] (60) 地理	[J] (523) 名詞	• [19] (72) の (の)	(92) 読点
• [21] (64) に	[zh] (36) ニ-格助	• [20] (74) 形態素解析 (形態素解析)	(41) 名詞
• [22] (66) お	[zh] (184) オク-五カ来	• [21] (64) は (は)	(52) ハ-格助
• [23] (68) い	[zh] (375) 五カ来	• [22] (66) 第一 (第一)	(46) 数名詞
• [24] (70) て	[zh] (72) テ-格助	• [23] (60) の (の)	(51) ノ-格助
• [25] (72) の	[zh] (512) 読点	• [24] (82) 間門 (間門)	(41) 名詞
• [26] (74) 形態素解析	[J] (523) 名詞	• [25] (86) で (で)	(60, 23) ダ-格助
• [27] (84) は	[zh] (47) ハ-格助	• [26] (88) ある <ある>	(75, 3) アル-五カ来
• [28] (86) 第	[X] (502) 数名詞+頭	• [27] (102) の (の)	(91) 句点
• [29] (88) 一	[zh] (501) 数名詞		
• [30] (90) の	[zh] (37) ノ-格助		
• [31] (92) 間門	[J] (523) 名詞		
• [32] (96) で	[zh] (243) ダ-格助		
• [33] (98) ある	[zh] (170) アル-五カ来		
• [34] (100) る	[zh] (434) 五カ来		
• [35] (102) の	[zh] (353) 句点		

図1. 形態素解析結果例  
【形態素リスト/単語リスト】

(例文出典) 情報処理学会 第42回全国大会 予稿集, 1991

## 2.2 構文解析系

構文解析では、文節間の係り受け解析をベースとし、係り受け文節対の組合せの枠組み [6] で、曖昧さを簡易的に扱う。これは、構造的な全解生成 (積算的組合せ) を文節ごとの係り先の可能性保持 (加算的組合せ; 図5) で代替するもので、組合せ爆発によるメモリ不足/処理過多を回避し、長文も解析可能とする。

日本語構文解析では、一般に意味情報による優先選択が用いられるが、効果は必ずしも大きくない。また、Q\_JP では、意味情報を持つことは、形態素解析系の小規模辞書の特徴と反することになる。そこで、係り先優先選択では、意味情報を用いず、係り受け可能な最近接文節を優先しながら、読点や副詞性単語の有無、表層的/構造的な類似性 [7] といった品詞/構文レベルの情報 [8] を基にしたヒューリスティック規則で、必要に応じて係り先を付け変えるという方式をとる。

構文解析系の全体は、文節単位に、品詞構成から文節属性を設定し、文節属性の組合せから係り受け可能文節対を検出し、さらに係り先文節を選択するという流れである。(図3~5) 各処理は、対応する規則に基づいて行われる。尚、係り先が認定できなかった場合に、形態素解析まで戻って、品詞の再選択を行い、構文解析をやり直す係り受け失敗回復機構をもつ。

## 3 Q\_JP ライブラリ

### 3.1 ライブラリ構成

Q\_JP は、C言語で実装し、応用系に容易に組み込めるようライブラリ化した (DOS 版, UNIX 版)。ライブラリは、形態素解析系 qjp.s 関数、構文解析系 qjp.k 関数を始めとする主要関数群と各種解析結果の情報等を参照/表示する補助関数やマクロ群からなる。

表1に、現版の辞書等のサイズを示す。辞書の合計は、約4500語/50KB 足らずである (今後、例外語 [拡張辞書] や和語のひらがな表記語 [ひらがな拡張辞書] 等を追加する必要があるが、その場合でも全体で100KB

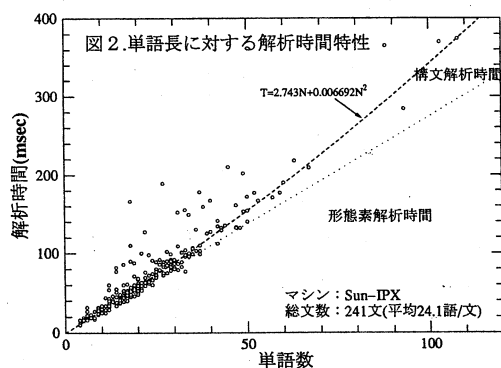
表1. Q\_J P (DOS 版) のファイルサイズ

ワークベンチ実行ファイル	151.2KB	} 合計 228.4KB
制御テーブルファイル	31.8KB	
実行時辞書ファイル (小計)	45.4KB	
		4476語/3839バイト
<hr/>		
拡張辞書	18.7KB	1836語/1469バイト
ひらがな辞書	5.6KB	740語/ 374バイト
ひらがな拡張辞書	10.9KB	869語/ 865バイト
形容動詞語幹辞書	10.2KB	1131語/1131バイト

程度であると見込んでいる)。辞書が小さいため、辞書データをメモリ上に展開する高速モードが可能となる。約30KBの制御テーブルには、形態素解析の533形態素品詞間の接続表、字種・品詞割り当て規則、接辞テーブル、辞書インデクス等が格納されている。

この他、形態素解析系は、文節内接続型規則・品詞曖昧さ解消規則・単語合成規則、構文解析系は、文節属性設定規則・係り受け可能文節対検出規則・係り先優先選択規則・係り受け失敗回復規則をもつ。これらは、IF-THEN形式のC言語関数として実装している。

上記の解析系系に対話/表示系を加え、対話型/バッチ型で利用できる日本語分析ツール「Q\_JP ワークベンチ」を作成した。ワークベンチの実行モジュールサイズは約150KBであり(DOS版)、辞書、制御テーブルを含めた全体でも230KBと小さい(表1)。尚、対話系には、ユーザが品詞の曖昧さ解消、係り先の変更を指示する機構を備える。



【係り受け木構造】

```

1[50](13):機械翻訳方式
2[49](12): する
3[48](11): 特徴と
4[47](11): ことを
5[46](10)>: 生成する
6[45]( 9)>: 〔他言語の〕文を
7[43]( 9)? : 〔その〕概念構造から
8[41]( 9)>: 定め、
9[40]( 8)>: 〔[[[翻訳すべき]言語から]なる]文の〕概念構造を
10[35]( 8)? : 〔[[[既に]選択された]概念記号と]意味関係記号]に基づき、
11[30]( 8)? : 設け、
12[29]( 7)? : 〔[[[利用者が]選択する]概念記号及び]意味関係記号を]受け入れる]手段を
13[23]( 7)? : 翻訳システムにおいて、
14[22]( 6)>: 翻訳する
15[21]( 5)>: 目標言語に
16[20]( 5)? : 〔[[[その]概念構造に]基づいて]翻訳すべき]言語を
17[15]( 5)? : 定め、
18[14]( 4)>: 概念構造を
19[13]( 4)? : 〔[[[各単語に]対応する]概念記号と、] [[各概念記号間の]関係を]示す]意味関係記号から
20[ 6]( 4)? : 分割し、
21[ 5]( 3)>: 単語に
22[ 4]( 3)? : 〔[[[翻訳すべき]言語から]なる]文を

```

【原文】

翻訳すべき言語からなる文を単語に分割し、各単語に対応する概念記号と、各概念記号間の関係を示す意味関係記号から概念構造を定め、その概念構造に基づいて翻訳すべき言語を目標言語に翻訳する翻訳システムにおいて、利用者が選択する概念記号及び意味関係記号を受け入れる手段を設け、既に選択された概念記号と意味関係記号に基づき、翻訳すべき言語からなる文の概念構造を定め、その概念構造から他言語の文を生成することを特徴とする機械翻訳方式

図3. 50 文節文の係り受け縮退木構造

【例文出典】公開特許公報(昭60), 1985

## 3.2 解析実験

「Q\_JP ワークベンチ」により、試作システムと同様に2コーパス[1:平均24.1語,241文/2:平均29.5語,210文]で解析実験を行った[1]。

### 実行性能

実行時メモリは、DOS版で約245KB(高速モード: +25KB)、UNIX版で約500KBであった。このメモリで約60文節までの長文の解析を行うことができる。

解析速度は、パソコン(80486/25MHz)で80~150\*語/秒、ワークステーション(Sun-SS20)で700~800\*語/秒であり(\*:高速モード)、50文節文(図3~5)をパソコン(80486)で1秒未満で解析できる(入出力時間含まず)。図2に単語長に対する解析時間の特性を示す。構文解析時間は、文長の二乗オーダとなるが、係数が小さいため、100単語程度の文長の範囲では形態素解析の線型の時間が大部分を占め、実用の範囲では解析処理全体としてもほぼ線型の特性になる。

比較できる構文解析系の実行性能データはないが、形態素解析系JUMAN 2.0の1MBテキストの解析時間76~93分、メモリ5MB[2](マシン不明)と比較すると、Q\_JPは、構文解析まで行なって、1MBテキスト換算で解析時間約6分(9秒/約200文25KB)、メモリ0.5MB(SS20)であり、各々1桁小さい。

### 解析性能

解析正解率は、若干改善しているが、試作システムの正解率[1][コーパス1/2に対し、形態素正解率: 99%/95%, 文節係り先正解率: 95%/90%, 1文全体の係り受け正解率: 71%/41%]とほぼ同程度である。

