

4B-5

形態素解析ツールによるかな漢字プログラミングの実現

*間瀬久雄 *辻 洋 *絹川博之 **川村隆雄

*(株)日立製作所 システム開発研究所 ***(株)日立製作所 ソフトウェア開発本部

1. はじめに

プログラム開発の効率向上の一方法である、かな漢字プログラミングの研究を進めている。これまでに、単語分かち書きのカナ文字列で記述可能なCOBOLプログラミング用簡易言語CORAL(図1)を開発した。CORALは大型計算機VOSシリーズ上で稼働している。

我々は、開発効率をさらに向上させるべく、より可読性に優れた非単語分かち書きかな漢字文による記述(図1)を検討し、形態素解析によってかな漢字プログラムを既存のCORALに変換するプリコンパイラのプロトタイプを開発した。これによりデバッグ効率が向上するほか、プログラムを仕様書として利用できる。なお、本プロトタイプにおけるかな漢字プログラムの構文は、従来のCORALの構文をほぼ継承している。本稿では、本プロトタイプの構成および機能について述べ、また、テストプログラムを用いた評価結果について考察する。

2. プロトタイプの構成

ここでは、本プロトタイプによりかな漢字プログラムをCORALに変換し、既存のコンパイラによりCORALをCOBOLに変換するという2段階方式を採用する。

本プロトタイプの構成を図2に示す。本プロトタイプはワークステーション上で稼働する。

(1) 入力

①かな漢字プログラム(図1)

ワープロやエディタ等を用いて作成したかな漢字プログラムを入力とする。かな漢字プログラムは、Shift-JISの2バイト文字で構成する。ただし、インデント等をしやすくするために、タブ・半角スペースの使用を許している。

②利用者定義データ

CORALでは、DBのスキーマ情報や画面名称等のグローバルデータを定義するためのテーブルがいくつか用意されている。利用者はプログラム中で参照するデータを各テーブルに記述する。その際、かな漢字文字列とカナ文字列とを対応させて記述する。

(2) 出力

①CORALプログラム

カナ文字列に変換されたCORALプログラムを出力する。単語辞書に未登録の単語(未登録語)はかな漢字のまま出力する。

②エラーメッセージ

日本語形態素解析処理の結果、未登録語が存在する場合、その旨を表すエラーメッセージを出力する。本プリコンパイ

<かな漢字プログラム>	<CORALプログラム>
1 #20.	1 #20.
2 画面1を読む。	2 ガン1 オ ヌム。
3 終わりなら終了。	3 くり 行 シヨリ的。
4 10回ループ。	4 10 カイ ル-プ。
5 #30を処理。	5 #30 オ シヨリ。
6 #30-10へ行く。	6 #30-10 ハ イク。

図1 かな漢字プログラムと従来のCORALプログラム

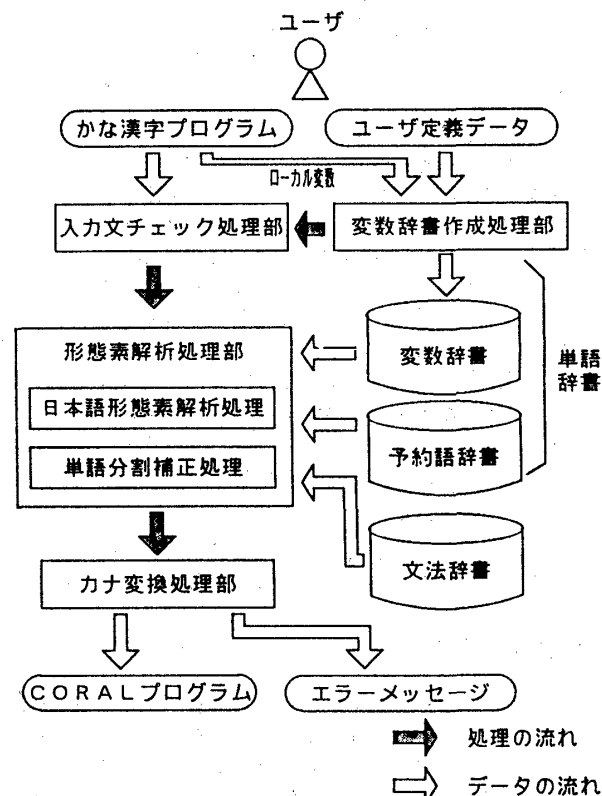


図2 本プロトタイプの構成

ラでは、語彙エラー、すなわちタイプミスや未宣言の変数の使用等によるエラーのみを扱っている。その他の細かいエラー検出については、既存のコンパイラに任せる。

(3) 予約語辞書

CORALで定義されている基本的な予約語は約120語である。予約語辞書にはこれらの予約語に関する形態素情報(品詞、活用等)、およびカナ文字列が格納されている。

(4) 変数辞書

変数辞書には、グローバル変数に関するデータとローカル変数に関するデータが格納されている。辞書アクセスの効率化のため、変数辞書は予約語辞書と同一構造をなす。変数は動的に変化するため、実行時に変数辞書を自動的に更新する

Kana-kanji Character Programming Based on a Japanese Morph Analysis Toolkit.

HIISAO MASE, HIROSHI TSUJI, HIROSHI KINUKAWA, TAKAO KAWAMURA

* Systems Development Laboratory, Hitachi, Ltd.

** Software Development Department, Hitachi, Ltd.

必要がある（（5）参照）。

（5）変数辞書作成処理部

ここでは、プログラムおよびユーザの定義したデータを用いて変数辞書を自動作成（更新）する。

①グローバル変数について

（1）で述べたように、グローバル変数に関しては、利用者が変数の見出しと対応するカナ文字列を定義するので、新たに定義された変数について、その見出しに関するデータを変数辞書に登録する。なお、グローバル変数の品詞はすべて名詞として登録する。

②ローカル変数について

CORALでは、ローカル変数を宣言する場合には行頭に識別子（"T", "V"等）をつけるので、ローカル変数はプログラムから自動抽出できる。また、ローカル変数に対応するカナ文字列を利用者に定義させることは使い勝手を悪くするので、システムが適当な変数コードを割り当てる。

（6）入力文チェック処理部

入力文の中に不正な文字が存在しないかを判別し、存在する場合には、エラーとしてその旨を報知する。また、1バイト文字のタブ、半角スペースを全角スペースに置き換える。

（7）形態素解析処理部

①日本語形態素解析処理

日英機械翻訳システムHICATS/JE (Hitachi Computer Aided Translation System / Japanese to English) の形態素解析処理を切り出してモジュール化した日本語形態素解析ツールを利用して入力文を単語分割し（最長一致法）、各単語に対応するカナ文字列を取得する。

②単語分割補正処理

形態素解析ツールによる単語分割結果がCORALの単語分割仕様と異なる部分について、単語分割結果を補正する。

（例1）" # 20 - 10" → "# 20-10"

（例2）未登録語 "A" + 未登録語 "B" → 未登録語 "AB"

（8）カナ変換処理部

予約語辞書および変数辞書から取得したカナデータに基づいて各単語をカナ文字列に変換して出力ファイルに格納する。

3. 評価と考察

（1）評価方法

テストプログラムを用いて本プロトタイプの性能を評価した。本評価で用いたかな漢字プログラムは、CORALの操作マニュアルに掲載されているサンプルプログラムを、ある人に頼んでかな漢字に変換して頂いたプログラム（176行）である。また、評価内容としては、①変換能力、②エラー検出能力（エラーメッセージの妥当性）、③処理速度の3点とした。

（2）評価結果

テストプログラムの一部とその変換結果を図3に示す。

処理結果として、次の二つのエラーを検出した。その他の部分については、すべて正しくカナ文字列に変換できた。

①語彙エラー：43行目「子」「位置付け」

CORALでは、「ピリオド+スペース」以降の部分はコメントを表す。図3の43行目の「子レコード位置付け」はコメントであるが、ピリオドの直後にスペースを挿入しなかったために、コメントとみなされず、その結果、「子レコー

<かな漢字プログラム>

```
41 C プロジェクト構成員一覧I-PFK = ?
42 -           |PF9| |.
43 A ?を処理 |#40|#20|. 子レコード位置付け
...
140           16回ループする。
141           | ---> 添え字I。
142           #60を処理する。
...
156 V 添え字1 S(4)
```

<変換後のCORALプログラム>

```
41 C プロジェクト構成員一覧I-PFK = ?
42 -           |PF9| |.
43 A ? オ シヨリ|#40|#20|. 子レコード位置付け
...
140           16 カイ ルブ スル。
141           | ---> 添え字 I。
142           #60 オ シヨリ スル。
...
156 V シヨリ S(4)
```

図3 評価に用いたテストプログラムの処理結果の一部

ド」が単語分割され、予約語「レコード」がカナに変換され、「子」「位置付け」が未登録語と判断されエラーとして表示された。これらがともに予約語であるとした場合、この箇所のエラーは検出されず、COBOLへのコンパイル時まで検出が遅れることになる。

②語彙エラー：141行目「添え字」

156行目でローカル変数「添え字1」が定義されているので、変数辞書に登録される。しかし、141行目には「添え字I」と記述されている。「I」はループ変数を表すシステム変数であるので、「添え字I」は「添え字」と「I」に単語分割される。そして「添え字」が未登録語と判断され、エラーとして表示された。

全体としては処理精度は良好であった。ただ、上述したように、未登録語が予約語や変数データを含む場合、エラー検出能力が落ちる恐れがある。語彙エラーをすべて正しく検出するためには、構文・意味的な制約を利用する必要がある。

なお、テストプログラムの処理時間は約3秒であり、妥当な速度範囲内にあると言える。ただし、今回の評価では、変数辞書作成処理に要した時間を含んでいない。本プロトタイプの処理時間の大半は日本語形態素解析ツールが占めている。

4. おわりに

プログラム開発効率の向上を目的として、可読性に優れた非単語分かち書きかな漢字によるプログラミングを検討し、かな漢字プログラムをCOBOLプログラミング用簡易言語CORALに変換するプリコンパイラのプロトタイプを開発した。今後は、COBOLコンパイルの前までに語彙エラーをすべて解消できるように、エラー検出能力を向上させていくとともに、プロトタイプの評価をさらに進めていく。

参考文献

- 1) 西森, 木山, 絹川: 汎用日本語形態素解析ツールの開発, 情報処理学会第44回全国大会(1992)。
- 2) (株)日立製作所: EAGLE2 CORAL入門編, VOS3 6180-3-822(1991)。