

4E-3

LangLABの辞書2次記憶化

西山研司 坂本浩一

(日立ソフトウェアエンジニアリング株式会社)

1. はじめに

prologにより自然言語処理システムを構築する場合、辞書の保存形式が問題となる。最も単純な方式は、辞書データをprologのホーン節で表現し、それらをすべてprologプログラム中に常駐させる方法である。こうしておけば、prologの特徴であるユニフィケーションによりシステムが勝手に辞書引きをしてくれるため、システム開発者が辞書引きに関していさゝい考慮する必要がなかった。しかしながら、実用システムの構築時には、大規模な辞書の作成が必要であるため、プログラムのメモリ容量の問題から、辞書をすべてメモリに常駐させることができないので、2次記憶化を図る必要が生じる。

本稿では、東京工業大学情報工学科田中穂積研究室で開発されたLangLABシステムの辞書2次記憶化と、その実用性の評価について報告する。

2. システム構成

LangLABは、prologで記述された自然言語処理システムの開発支援システムである。図1にLangLABシステムの構成を示す。

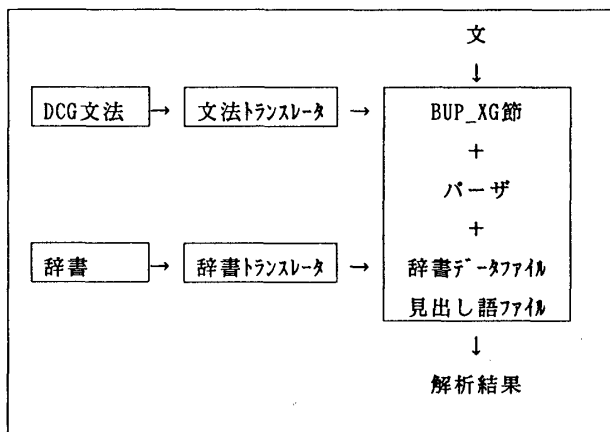


図1. LangLABシステムの構成

LangLABシステムは、ユーザが記述したDCG形式の文法規則をBUP_XGと呼ばれるprolog

プログラムに変換し、このプログラムとprolog記述のパーザおよび辞書に解析すべき文を与えると、構文処理と意味処理とを融合したボトムアップで縦型探索の自然言語処理を行うことができるシステムである。

3. 辞書2次記憶化の方法

3.1 従来の辞書引き

従来のLangLABにおいては、辞書引きはdictaという述語の起動により行っていた。

<辞書>

```
"door" : dicta(door, [[n, [[cnt|'+']], door]]], []).
"key"  : dicta(key, [[n, [[cnt|'+']], key]]], []).
```

<辞書引きの起動>

```
dicta(WORD, ARGs, TRIEs) .....①
```

ここで、例えば、WORDを'door'とユニファイした状態で、①を起動すると、dictaの第1項のWORDという変数と'door'のマッチングにより辞書引きが行われる。

3.2 辞書2次記憶化による辞書引き

(1) 辞書ファイルの構成

辞書ファイルの構成を図2に示す。

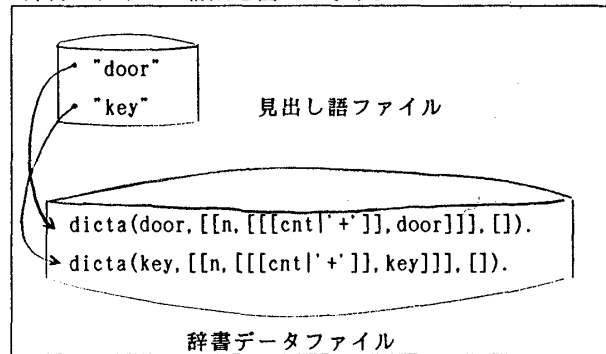


図2. 辞書ファイルの構成

(2) 辞書引き方法

新しい辞書引きのための述語をdictbとする。C言語で記述された辞書サルーチンdict_searchにより、辞書データファイル中の辞書prologのホーン節を得る。得られたホーン節をassertした後、従来の辞書引

き述語"dicta"を起動する。

```
dictb(WORD, ARGs, TRIEs):-
    call(dict_search, [WORD], F), !, assert(F),
    dicta(WORD, ARGs, TRIEs).
```

(3)バッファによる高速化

辞書引き高速化のために、以下の二種類のバッファを設ける。prologの場合、これらのバッファへの格納は、assert述語により行える。

(a)辞書バッファ

一度辞書引きした結果を格納

(b)未登録語バッファ

辞書引きの結果、未登録語であることがわかった単語を格納

以上のメカニズムを組み込んだdictbの内容を図3に示す。

```
dictb(WORD, ARGs, TRIEs):-
    not_found(WORD), !, fail.          .....①
dictb(WORD, ARGs, TRIEs):-
    dicta(WORD, ARGs, TRIEs), !.      .....②
dictb(WORD, ARGs, TRIEs):-
    call(dict_search, [WORD], F), !, assert(F),
    dicta(WORD, ARGs, TRIEs).        .....③
dictb(WORD, ARGs, TRIEs):-
    assert(not_found(WORD)), !, fail.  .....④
```

①対象単語が未登録語バッファ(not_found)にあれば、failする。
 ②対象単語が辞書バッファ(dicta)にあれば、それを辞書引き結果とする。
 ③辞書ファイルのサーチ。
 ④辞書未登録語を未登録語バッファに格納。

図3. dictbの内容

4. 評価

C-prolog上で、LangLABの辞書を2次記憶化する前(辞書メモリ中)と後(辞書ファイル中)の実行結果を表1に示す。

表1の結果より、2次記憶に辞書をおいても、LangLABは十分高速で、実用に耐えうる速度であることがわかる。

5. おわりに

LangLABシステムを題材に、prologで記述された自然言語解析システムの辞書2次記憶化の方法について述べた。

表1. 辞書2次記憶化速度比較

例 文	解析時間*1[msec]	
	メモリ	ファイル
I open the door with a key.	2,888	2,717
Tell me who swims in the river.	3,250	3,017
Tell me what to do.	2,000	1,850
I tell you that it happens.	5,167	4,967
I open not only the door but also the window.	2,217	1,850
I know as if he bought the car.	3,700	3,533
I saw the man in the park with a telescope.	5,817	5,733
It is not tied to a particular domain of applications.	8,133	7,767

*1:全解析木が得られる時間

[謝辞]

日頃御指導いただく東京工業大学田中穂積教授ならびに徳永助手をはじめ田中研究室の方々に深謝いたします。

[参考文献]

[1]上脇正、田中穂積、"辞書のTRIE構造化と熟語処理" PROCEEDING OF LOGIC PROGRAMMING CONFERENCE '85 12.2.
 [2]田中穂積、"自然言語処理システムLangLAB"、情報処理学会第33回予稿集(1986)
 [3]奥村学、田中穂積、"LangLABにおける高水準辞書記述言語SRL/O"、情報処理学会第33回予稿集(1986)
 [4]Hartono、田中穂積、"LangLAB上出のインドネシア語文法試作"、情報処理学会第33回予稿集(1986)
 [5]沼崎浩明、田中穂積、"LangLABにおける文法開発環境"、情報処理学会第33回予稿集(1986)
 [6]上脇正、田中穂積、沼崎浩明、"LangLABの構文処理アルゴリズムの高速化"、情報処理学会第33回予稿集(1986)
 [7]今野聡、田中穂積、"左外置を考慮したボトムアップ構文解析システム"、コンピュータソフトウェアvol.3 No.2(1986)
 [8]QUEK CHEE HUEI、"英語の文法の機能拡張に関する研究"、東京工業大学情報工学科卒業論文(1988)