

並列形態素解析システム L A X の実現

1T-5

杉村 頌一、赤坂 宏二、松本 裕治

(財) 新世代コンピュータ技術開発機構

1. はじめに

並列実行可能な日本語の形態素解析手法[6]についてのべる。本手法では、形態素辞書をGHC[7]等の並列論理型言語で記述されたプログラムへ変換する。形態素辞書には、形態素の①見出し、②カテゴリー、形態素間の接続可能性を調べるための③接続素性情報[1]、を記述する。接続素性情報は、④接続マトリクス、と共に曖昧さ削減に用いる。解釈に曖昧さがあれば、解析終了時に結果を縮退して一度に出力する。解析結果は、構文解析プログラムSAX[5]への入力に用いる。本方式は未登録語処理をリニアオーダで行なう。また、逐次論理型言語上でも、高速な解析が可能である。

2. 形態素解析アルゴリズム

本アルゴリズムはEarleyのアルゴリズム[2]や、チャートパーズング[4]などで用いられているアルゴリズムに基づいている。

以下、文字列検索プロセスの起動方法とプロセス間の通信方式について、その概要を述べる。

2. 1 プロセスの起動

以下図1を例に説明を行なう。

入力: く る ま で 待 っ .

	1	2	3	4	5	6	7
F 1	dt	pl3	pl4				
	F 2	dt					
		F 3	dt	p35			
			F 4	dt			
				F 5	dt		
					F 6	dt	
						F 7	dt

-図1-

'dt'は初期状態で起動される検索プロセスで、入力文字各々の位置に対応して起動される。初

期状態で、隣接する'dt'はストリームによって一次元に連結される。

'dt'は対応した文字を先頭とする部分文字列についての形態素の検索を次の様に行う。

まず'dt'は対応する1文字を形態素辞書から検索する。これと共に、対応する1文字に始まる長さ2以上の形態素を検索するサブプロセスを起動する。例えば図1で行F1の'dt'は「く」を検索すると共に、「くる」を検索する'pl3'と「くるま」を検索する'pl4'を起動する。起動された検索処理は全て並列に実行が可能である。この処理は次節に説明するように、図4のようなPrologまたは、GHCのプログラムによって実現されている。

'dt'とそのサブプロセスの出力は、d-listによって連結され、1本のストリームとしてまとめられる。

以上の方法により、'dt'の対応する文字を先頭とする部分文字列の検索結果は、'dt'の出力として、その右隣のプロセスへ渡る。

2. 2. 形態素辞書の検索プロセスへの変換

上記の検索プロセス'dt'とそのサブプロセスは以下の方法で形態素辞書から変換して得られる。例えば、図2のような形態素辞書を考える。

見出し	形態素カテゴリー	接続素性
く	力変動詞語幹	1
くる	わ行動詞語幹	2
くるま	名詞	3
くろ	形容名詞	4

-図2-

上記の辞書はTRIE構造[3]を用いて図3の様に表示出来る。

```
[く 力変動詞語幹 1
 [る わ行動詞語幹 2 [ま 名詞 3]]
 [ろ 形容名詞 4]]
```

-図3-

更に図3の構造は図4のようなPrologのプログラムに変換される。

Parallel Lexical Analyzer LAX
Ryoichi Sugimura, Kohji Akasaka, Yuji Matsumoto
Institute for New Generation Computer Technology

