

TR-367

構文解析システムSAXのデバッグ環境

山崎重一郎、杉村領一
赤坂宏二、松本裕治

April, 1988

©1988, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03) 456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

構文解析システム SAX のデバッグ環境

1. はじめに

ICOTでは、談話理解システムの研究のベースとなるソフトウェアツール群を汎用日本語処理系(LTB)として研究開発を進めている。SAX [1]は、その中に含まれるツールの一つで、DCG [2]の文法記述を対象とし、並列論理型言語による上昇型の並列構文解析を目的とする構文解析システムAXの逐次論理型言語版システムである。SAXは並列解析の擬似システムとしてのみでなく、逐次論理型言語による効率の良い全探索構文解析システムとしても意味を持つ [3]。しかしながら、SAXによる解析は、一つの解析木を構成するための処理が、各段階で分散して処理されるために実行過程のトレースによるデバッグは極めて困難になる。このため、SAXの文法デバッグ環境は、解析実行による解析の履歴を、利用者の要求にしたがって取り出し、文法規則の適用過程として表示する環境として実現した。本論文では、まずSAXによる解析の実行を概説し、次にSAXのデバッグ環境の実現方法について述べる。

2. SAXによる解析の実行

SAXによる解析は、ボトムアップにある要素が構成されると、それを基にしてより大きな解析木を構成することによって進行する。このときの文法規則の選択は、規則の右辺の左隣要素によって決められる。そして新しい要素がボトムアップに構成されるたびに、選択された規則の2番目以降の要素が埋め合わせられて、完全な木が構成されていく。SAXでは、文法に現れる全ての単語および統語範疇が論理型言語の述語として表現されており、解析の進行による要素の構成は、述語の呼び出しに対応している。呼び出された述語には、新しい規則の選択と、すでに予想されている規則の充足の検査の2種類の役割がある。文法に規則の右辺の左隣要素として現れる要素は新しい規則の選択の役割を担うので、TYPE 1要素と呼び、規則の右辺の2番目以降の要素は、規則の充足の検査の役割を担うので、TYPE 2要素と呼ぶ。文法の複数の規則に現れる要素の中には、TYPE 1でもTYPE 2でもあるような要素も存在する。また、規則の右辺の右隣の要素が充足されることは、予想されていた規則による解析木が完成したことを意味するので、その規則の左辺要素の呼び出しが行われる。したがって、規則の右辺右隣の要素は規則の左辺要素の呼び出しの役割を担うことになる。TYPE 1で右隣要素でもある要素をTYPE 10と呼び、TYPE 2で右隣の要素をTYPE 2-LASTと呼ぶ。またTYPE 2で右隣要素でない要素をTYPE 2-MIDDLEと呼び、右隣要素でないTYPE 1要素を狭義のTYPE 1と呼ぶ。SAXでは、文法の各規則の右辺にそれぞれの位置を表す識別子を付与し、各述語は、この識別子を使用して規則の選択や充足の検査を行う。

〈例1〉 次の文法には、識別子が付与されている。id_1 ~ id_5 が識別子である。

```
規則通番
1: n    --> [cars].
2: aux --> [can].
3: v    --> [be].
4: a    --> [useful].
5: s    --> np, id_1 aux, id_2 vp.
6: np   --> n.
7: vp   --> v, id_3 a.
8: s    --> np, id_4 vp.
9: pp   --> prep, id_5 np.
```

TYPE 1

上の文法のvは狭義のTYPE 1要素であり、次のような述語になる。

```
v(X, [id_3(X) | Y], Y).
```

上の文法のnはTYPE 10であり、次のような述語になる。

```
n(X, Y1, Y2) :- np(X, Y1, Y2).
```

TYPE 2

上の文法のauxはTYPE 2-MIDDLEであり、次のような述語になる。

```
aux([], Y, Y).
aux([id_1(X) | Z], [id_2(X) | Y1], Y2) :- aux(Z, Y1, Y2).
aux([_ | Z], Y1, Y2) :- aux(Z, Y1, Y2).
```

TYPE1かつTYPE2

上の文法のnpは狭義のTYPE1かつTYPE2-LASTであり、次のような述語になる。

```
np(X, Y1, Y3) :- np1(X, Y1, Y2), np2(X, Y2, Y3).
np1(X, [id_1(X), id_4(X) | Y2], Y2).
np2([], Y, Y).
np2([id_5(X) | Z], Y2, Y4) :- pp(X, Y2, Y3), np2(Z, Y3, Y4).
np2([_ | Z], Y2, Y3) :- np2(Z, Y2, Y3).
```

SAXによる解析の実行は、入力文を構成するそれぞれの単語に対応する述語を次のように共有変数で結合して一斉に起動することによって開始される。

〈例2〉 入力文 cars can be useful.

に対する解析の実行は次のようにして行う。

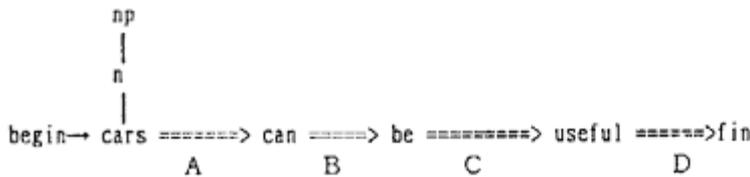
```
cars(begin, A, []),
can(A, B, []),
be(B, C, []),
useful(C, D, []),
fin(D).
```

ここで、beginは入力文の左端を表す識別子であり、finは解析結果を検査する述語である。

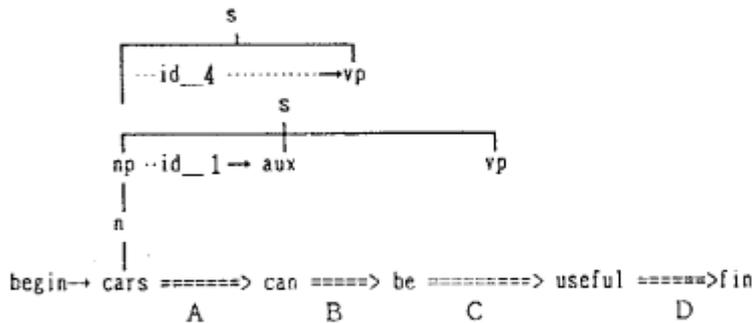
そして、解析の進行とともに、より上位の要素の述語が呼び出され、述語を整えている共有変数は、それらの述語が出力した識別子によって具体化されていく。述語が共有変数を通じて入出力する識別子は、その引数として識別子のリストを持つような再帰的構造を持っており、解析の進行によって、より上位の木構造構成されていくことは、識別子とリストのより深い構造が具体化されていくことに対応している。

〈例3〉 SAXの解析において一つの解析木が構成される過程と共有変数の具体化の過程を例2の入力文の解析によって例示する。

① 述語 cars によつて述語 n が起動され、述語 n によつて述語 np が起動される。



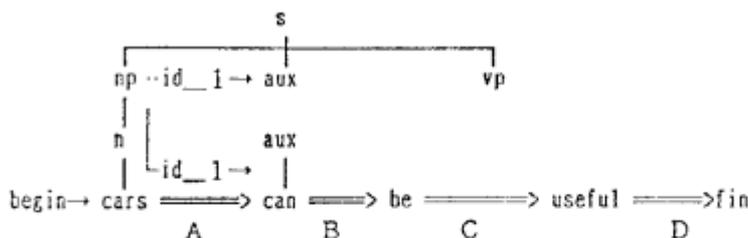
② TYPE1 述語 np が id_1 と id_4 を出力する。これは、次のような二つの木構造を予想することに相当する。



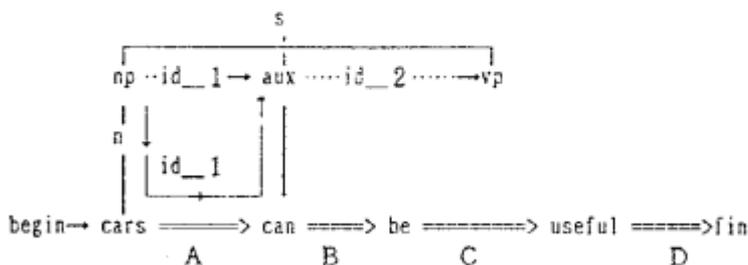
このとき共有変数 A は [id_1(begin), id_4(begin)] に具体化される。

以下では、id_1 によって予想された木構造にのみ注目する。

③ 共有変数 A が述語 cars の出力識別子によって具体化されることにより、述語 can が述語 aux を起動する。

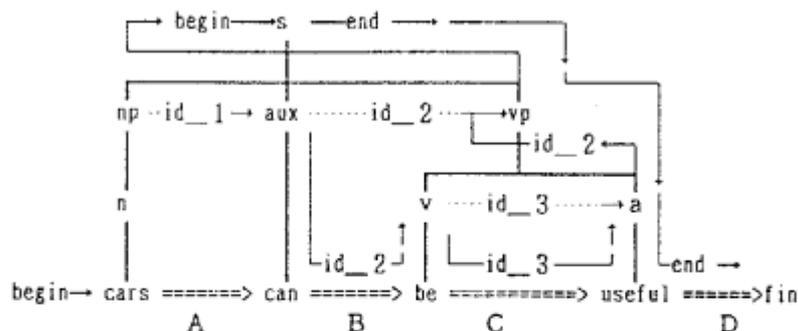


④ TYPE 2—MIDDLE 述語 aux は共有変数 A を通じて渡ってきた識別子 id_1 をうけとると、id_2 を出力する。これは、can から構成された aux が予想された文法規則の aux に適合したことを意味しているので、can から構成された aux の部分木が、その上の不完全な木の aux に接続されたことになる。



このとき共有変数 B は [id_2 (begin)] に具体化される。

⑤ 同様にして be、によって v が起動されると、述語 v は id_3 を出力する、このとき be の入力の id_2 は、id_3 の引数として運ばれる。そして useful が a を起動して、a が id_3 の入力により vp を起動するとき、vp の入力として、運んでいた id_2 を渡す。述語 vp は id_2 を受け取ると、述語 s を起動する。このときも同様に、id_2 の引数として運んでいた識別子 begin を s の入力として渡す。述語 s は、解析の目標となる要素なので、特別に予め begin を受け取ると end を出力する TYPE 2—MIDDLE 要素として定義されている。begin は s を構成する要素の左端を表す識別子であり、end は右端を表す識別子になっている。s の出力の end は vp、a、useful の出力に順に具体化され、fin との共有変数 D に具体化される。このために、cars の入力に begin が具体化されていて、useful の出力変数に end が具体化されたことになるが、この状況は、cars から useful までの間が s を構成することを表している。



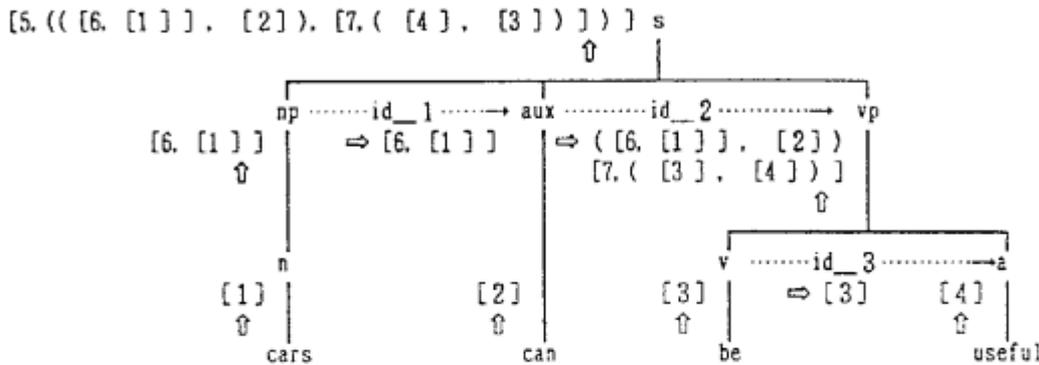
このとき共有変数 C は [id_3 ([id_2 (begin)])] に具体化され、共有変数 D は [end] に具体化される。

3. 解析結果の累積方法

解析の成功は、fin が構成目標要素の右端を表す識別子 (上の例での end) を検出することによって知ることができるが、実際の解析結果についての情報は、これだけではわからない。解析結果を知るためには、識別子の引数に変数を追加し、識別子によって文法規則の適用履歴を累積することが必要である。識別子による文法規則適用履歴の累積は、TYPE 10 及び TYPE 2—LAST で成功した規則の規則番号を追加し、TYPE 2—MIDDLE 及び TYPE 2—LAST で識別子によって運ばれてきた規則適用履歴と配下からきた規則

適用履歴をマージし、狭義のTYPE 1では配下からきた規則適用履歴を出力識別子にのせて送ることによって実現できる。

〈例4〉 識別子による規則通番の累積例（数字は例1の文法の規則通番である）。



ここで、 $[C, (A, B)]$ は、規則Aによる結果と規則Bによる結果が規則Cによってまとめられたことを意味する。ただし、規則Cが3項以上の規則によって構成される場合は、 $[C, ((A, X), Y), B]$ となる。また、終端記号から構成されるときは $[C]$ と記す。

このように、識別子を用いて解析履歴を累積すると、結果として文末の識別子 (end) にその解析の過程で使用された全ての文法規則の履歴が残されていることになる。

4. 文法デバッグ環境

文法デバッグ環境は、識別子とDCG規則の位置の対応づけを行う機能を持ち、解析終了後に識別子の集まりとして共有変数に残されている解析の履歴を利用者の要求にしたがって取り出しそれをDCG規則の適用過程として表示する環境を与える。

4. 1 識別子とDCG規則の対応づけ

SAXの実行結果として得られるのは識別子のリストであるが、利用者にとって必要なのはこれらの識別子が対応している文法規則の位置であり、識別子自体は必要ではない。このためにデバッグ環境では識別子に対応する文法規則の位置を求める機能を持ち、利用者が識別子を意識せずにデバッグできるようにする。これは、識別子名を規則通番と要素位置の対によって構成し、この情報を利用することによってDCG規則の中の対応する位置を知ることによって実現できる。

〈例5〉 識別子名の規約とDCG規則の対応

識別子名 : id_1 2_2 (id_規則通番_要素位置)

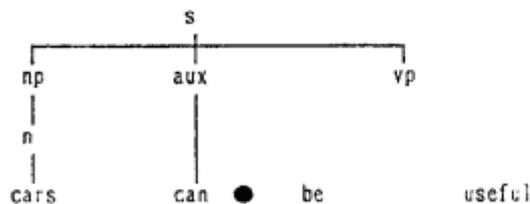
1 2 番目の規則 : a --> b, c, d, f.

この場合、識別子名の情報により (識別子の位置を●で表示すると)

a --> b, c, ●d, f. が得られる。

4. 2 解析過程の表示方法

文法デバッグ環境は、基本的に、構文木の構成過程を表示する環境として実現する。しかしながら、実際の解析過程で現れる木構造は莫大な数になり、適切な選択手段を持たなければ実用にはならない。この選択手段の一つとして、入力文の指定した単語の区切りの位置より左側に構成される不完全な木構造を全て表示することを考える。これは、例3で be と useful の間を区切り位置とすると、次のような不完全な木を表示することに対応している。

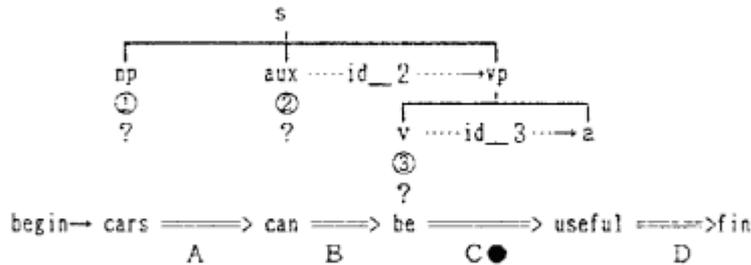


この木構造は cars によって予想された規則に can までが適合していることを示している。この木と同時に表示される他の木構造も can 迄が成功している木構造のみなので、この基準は表示される木構造の選択手段としてはたらく。

4. 3 区切り位置の左側の不完全な木の表示方法

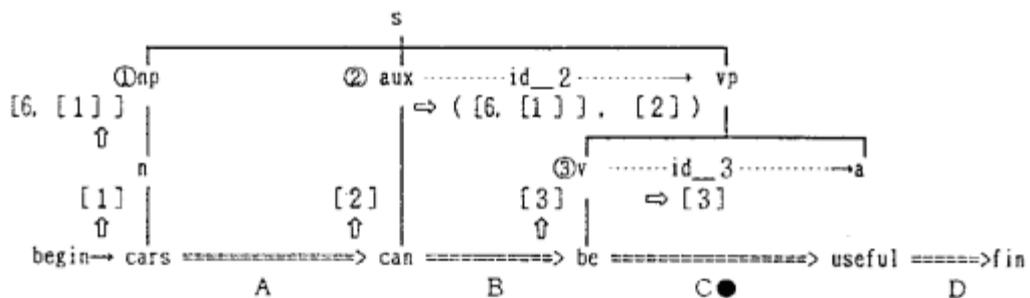
区切り位置の共有変数をモニタして得られる情報を、例3の解析木が構成されたときの共有変数Cをモニタして得られる情報によって説明する。

Cの中ではid_3が選ばれ、id_3の中でid_2が選ばれるのでCをモニタすることによって次のような木構造を得ることができる。



しかしながら、識別子からの情報からは、①や②や③に支配される部分木の情報を得ることはできない。これらの情報を得るために、識別子による文法規則の規則通番の累積を利用する。これによって、id_3で③配下の規則適用履歴がわかり、id_2によって①配下と②配下の規則適用履歴がわかるので、結果として、区切りの左側で構成される不完全な部分木が得られる。

$C = [id_3([id_2(begin, ([6, [1]]), [2]))], [3]]$



4. 4 指定範囲から構成される完全な部分木の表示方法

上の不完全な木構造は useful によって a が構成されることを期待していることも示している。もし、このsの解析が失敗していた場合、利用者は、解析の失敗の原因が a を構成する文法にあることを知ることができる。一般的に、区切り位置の左側に利用者が期待している解析木の構造が現れているにもかかわらず、解析に失敗している場合、その右側の要素を構成する文法に誤りが含まれていることがわかるので、誤り発生位置の特定は、区切り位置を少しずつ右に変更しながら調べていくことによって発見できる。しかし、この方法は、いかにも作業効率が悪そうである。このかわりに、利用者が指定した範囲で構成されている完全な部分木を全て表示できれば、その部分で利用者が期待している統語範疇が完成しているか否かを知ることができるので、誤り発生位置の範囲をもっと簡単に絞りこめる。このために、SAXのデバッグ環境に、指定範囲で構成された完全な部分木を得る機能をもたせることにした。指定範囲で構成された完全な部分木を得る方法は、基本的に解析目標の統語範疇まで完成した解析木を表示する方法と同じであり、指定範囲の左からbeginに対応する識別子を入力して、指定範囲の右側の要素の出力側の共有変数で、その識別子と対をなすendに対応する識別子を検出することによって実現できる。このために、すべての非終端記号に対して、次のようなTYPE 2 MIDDLEの定義を追加する。

要素名 ([left(N, LC) | Z], [right(N, LC) | Y1], Y2, LCC) :- 要素名 (Z, Y1, Y2, LCC).

ただし、N は、入力文の単語の位置を表す番号であり、LCおよびLCC は規則の適用履歴を累積するための変数である。

そして、解析の起動を次のように変更する。

