システムソフトウェア特論'19 # 8 構文解析(3)

久野 靖 (電気通信大学)

2019.1.9

上向き解析

シフト還元解析器

前回までで、構文木を上の方(出発記号の側)から組み立てて行く下向き解析を扱いました。今回は構文木を下の方から組み立てていく上向き解析(bottom-up parsing)を扱います。上向き解析を行う構文解析器は通常、シフト還元解析器(shift-reduce parser)の形を取ります。そこで具体的な解析方式の説明に先立ち、シフト還元解析器の枠組について説明しましょう。

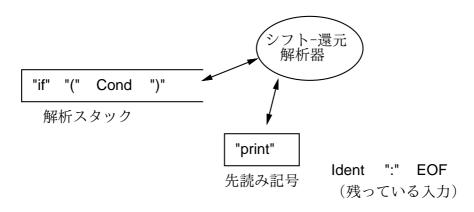


図 1: シフト還元解析器の枠組

シフト還元解析器の基本的な道具立てを図??に示します。解析には構文記号を積むスタック (解析スタック) 1 つと、先読み記号1個を使用します (より一般的には入力記号をn 個先まで見ることも考えられますが、動作の原理としては同じです)。そして、解析器の行う動作は必ず次の2 つのうちいずれかです。

- シフト ― 先読み記号をスタックに積み、入力を進める。
- 還元 生成規則 $A \to \alpha$ に対応し、 α の長さ分スタックを取り降ろして代りに A を積む。

シフト還元解析器 (2)

ただしこれでは終りがないので、文法で出発記号Sを左辺に持つ生成規則が1つだけになるようにし、その規則番号を1番とします。そして還元時に規則番号が1番だったら解析を終了します。以下では例題として次の文法を使用しますが、これも上記の条件を満たしています。

```
prog ::= statlist
statlist ::= statlist \ stat \mid \epsilon
stat ::= ident = expr \ ; \mid \mathbf{read} \ ident \ ; \mid \mathbf{print} \ expr \ ;
\mid \mathbf{if} \ (\ cond \ ) \ stat \mid \mathbf{while} \ (\ cond \ ) \ stat \mid \{ \ statlist \ \}
cond ::= expr < expr \mid expr > expr
expr ::= ident \mid iconst
```

先の言語の最小限のプログラム「read x; print x;」をシフト還元解析器で解析する様子を図??に示します。これからわかるように、シフト還元解析器では入力は次々にスタックに移され、スタック上に生成規則の右辺と一致するものができると左辺の非端記号に置き換える、という形で解析が進んで行きます。ここまでで「なるほど、スタック上に規則の右辺が現れたら還元、それ以外ではシフトをすればいいのか」と思われたかも知れませんが、それほど単純ではありません。上の例でもprint文のidentはスタックに積まれた直後にexprに還元されていますが、一方read文のidentはそのまま残されています。これはもちろん、read文の規則の右辺がread ident;であるからですが、ともかく右辺が現れたらすぐ還元してよい、というものでないわけです。

このシフトと還元の選択を正しく行うところがシフト還元解析器の「きも」です。その手法としてここでは最も広く使われているLR解析器(LR parser)を説明します。

シフト還元解析器 (3)

LRの最初のLはLLと同様「left-to-rightに1回入力を捜査するだけで解析を行う」という意味、2文字目のRは「rightmost derivationの逆順の系列を生成する」という意味です。逆順なのは上向き、つまり出発記号から遠いところから構文木をつくっていくため出発記号が最後になるからで、逆順にする前で考えれば左側の導出が先に出てきます(入力を左から読むのでそれが普通です)。

実際、図??に現れる規則を下から順に並べて導出列を作ると、 $prog \Rightarrow statlist \Rightarrow stat \ statlist$

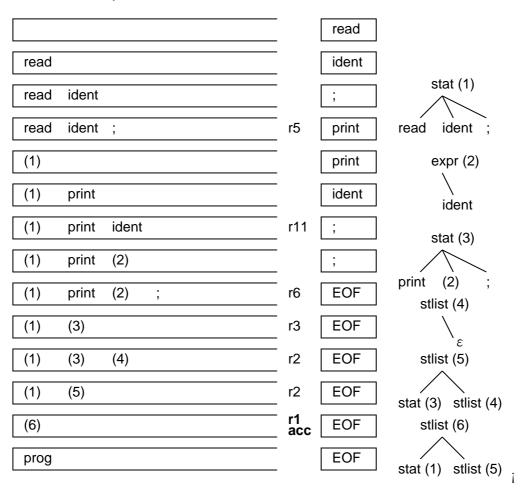


図 2: シフト還元解析器による構文解析

LRオートマトンと項

先の *ident* を還元するかどうかは、先に read をシフトしたか print をシフトしたかによって選択が違っていました。「置かれた状態によって同じ入力に対する動作が変化する」ことが問題ですが、この種の問題はオートマトンで表現するのが適しています。

ここではLRオートマトンとよばれるものを使いますが、それは字句解析のときの有限オートマトンとはだいぶ違います。まず、LRオートマトンでは各状態が「構文規則の右辺のどの位置にいるか」に対応します。これを表すために、構文規則の右辺の任意の位置に \bullet を記入して「現在いる位置」を表します。これを項(item)と呼びます(より厳密にはLR(0)項。0は項の中に先読み記号の情報が含まれていないことを意味する)。

例えば、規則 $stat \rightarrow if$ (cond) stat からは6つの項ができます (項は[]で囲んで表す)。

```
[ stat \rightarrow \bullet if ( cond ) stat ]

[ stat \rightarrow if \bullet ( cond ) stat ]

[ stat \rightarrow if ( \bullet cond ) stat ]

[ stat \rightarrow if ( cond \bullet ) stat ]

[ stat \rightarrow if ( cond ) \bullet stat ]

[ stat \rightarrow if ( cond ) stat \bullet ]
```

LR(0) オートマトンの作成

それでは、上向き構文解析のためのオートマトンの構成方法について説明しましょう。字句解析の決定性オートマトンでは、1つの状態は複数の正規表現のそれぞれ特定の場所に「並行して」対応していました。ここでも同様なことを考える必要があります。例えば、次の状態にいるものとします。

```
[ stat \rightarrow \mathbf{if} ( \bullet cond ) stat ]
```

つまり cond の直前にいるわけですが、そこで次の生成規則を参照します。

```
cond \rightarrow expr ; expr
cond \rightarrow expr ; expr
```

ということは、上の状態は同時に次の2つの状態にも「並行して」存在していることになります。

```
[ cond \rightarrow \bullet \ expr \ ; \ expr \ ]
[ cond \rightarrow \bullet \ expr \ ; \ expr \ ]
```

さらに次の規則も参照します。

```
expr 	o ident
expr 	o ionst
```

ということは、これらの状態にいるということは次の状態にもいることになります。

```
[expr \rightarrow \bullet ident][expr \rightarrow \bullet iconst]
```

LR(0)オートマトンの作成 (2)

このように「並行していることになる」項の集合を (LR(0)項集合の) 閉包 (closure) と呼びます。したがって、LRオートマトンの各状態は項の集合で、なおかつ閉包になっている必要があります。そして、状態間の遷移については、●の直後にある記号(端記号・非端記号の双方)をラベルに持つ遷移により、●がその記号の後ろに移動した項(の閉包)の状態に移ることになります。

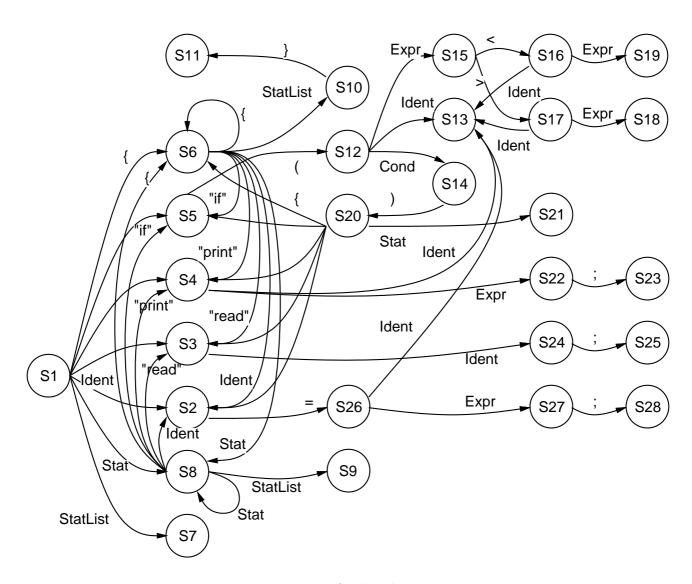


図 3: LR(0) オートマトン

先の文法に対応する LR(0) オートマトンを作成すると、図?? のようになります。中身の項集合は表?? に別に示しています。

LR(0)オートマトンの作成 (3)

表 1: LR(0) オートマトンの各状態

```
S1:
           [prog \rightarrow \bullet statlist]
                                                                    S12:
                                                                               [stat \rightarrow if ( \bullet cond ) stat]
           [statlist \rightarrow \bullet statstatlist]
                                                                                [cond \rightarrow \bullet expr < expr]
           \lceil statlist \rightarrow \bullet \rceil
                                                                                [cond \rightarrow \bullet expr > expr]
           [stat \rightarrow \bullet ident = expr;]
                                                                                [expr \rightarrow \bullet ident]
           [stat \rightarrow \bullet \text{ read } ident ;]
                                                                                [expr \rightarrow \bullet iconst]
            [stat \rightarrow \bullet \text{ print } expr ;]
                                                                                [expr \rightarrow ident \bullet]
                                                                    S13:
           [stat \rightarrow \bullet \text{ if } (cond) stat]
                                                                                [stat \rightarrow if (cond \bullet) stat]
                                                                    S14:
           [stat \rightarrow \bullet \ \{ \ statlist \ \}]
                                                                    S15:
                                                                                [cond \rightarrow expr \bullet < expr]
 S2:
           [stat \rightarrow ident \bullet = expr;]
                                                                                [cond \rightarrow expr \bullet > expr]
 S3:
           [stat \rightarrow read \bullet ident ;]
                                                                    S16:
                                                                               [cond \rightarrow expr < \bullet expr]
 S4: [stat \rightarrow print \bullet expr ;]
                                                                                [expr \rightarrow \bullet ident]
           [expr \rightarrow \bullet ident]
                                                                                [expr \rightarrow \bullet iconst]
           [expr \rightarrow \bullet iconst]
                                                                    S17:
                                                                               [cond \rightarrow expr > \bullet expr]
           [stat \rightarrow if \bullet (cond) stat]
                                                                                [expr \rightarrow \bullet ident]
 S5:
 S6:
           [statlist \rightarrow \bullet statstatlist]
                                                                                [expr \rightarrow \bullet iconst]
           [statlist \rightarrow \bullet]
                                                                    S18:
                                                                                [cond \rightarrow expr > expr \bullet]
           [stat \rightarrow \bullet ident = expr;]
                                                                                [cond \rightarrow expr < expr \bullet]
                                                                    S19:
            [stat \rightarrow \bullet \text{ read } ident ;]
                                                                                [stat \rightarrow \bullet ident = expr;]
                                                                    S20:
           [stat \rightarrow \bullet \text{ print } expr ;]
                                                                                [stat \rightarrow \bullet \text{ read } ident ;]
           [stat \rightarrow ullet \ if \ (cond) \ stat]
                                                                                [stat \rightarrow \bullet \text{ print } expr ;]
           [stat \rightarrow \bullet \ \{ \ statlist \ \}]
                                                                                [stat \rightarrow \bullet \text{ if } (cond) stat]
           [stat \rightarrow \{ \bullet statlist \}]
                                                                                [stat \rightarrow if (cond) \bullet stat]
         [prog \rightarrow statlist \bullet]
                                                                                [stat \rightarrow \bullet \ \{ \ statlist \ \}]
 S7:
           [statlist \rightarrow \bullet statstatlist]
                                                                                [stat \rightarrow \text{if } (cond) stat \bullet]
 S8:
                                                                    S21:
            \lceil statlist \rightarrow stat \bullet statlist \rceil
                                                                    S22:
                                                                                [stat \rightarrow print \ expr \bullet \ ;]
           [statlist \rightarrow \bullet]
                                                                                [stat \rightarrow print \ expr; •]
                                                                    S23:
           [stat \rightarrow \bullet ident = expr;]
                                                                                [stat \rightarrow read \ ident \bullet \ ;]
                                                                    S24:
           [stat \rightarrow \bullet \text{ read } ident ;]
                                                                    S25:
                                                                                [stat \rightarrow \texttt{read} \ ident \ ; \ ullet ]
           [stat \rightarrow \bullet \text{ print } expr ;]
                                                                    S26:
                                                                                [stat \rightarrow ident = \bullet expr ;]
           [stat \rightarrow ullet \ if \ (cond) \ stat]
                                                                                [expr \rightarrow \bullet ident]
           [stat \rightarrow \bullet \ \{ \ statlist \ \}]
                                                                                [expr \rightarrow \bullet iconst]
 S9:
           [statlist \rightarrow statstatlist \bullet]
                                                                    S27:
                                                                                [stat \rightarrow ident = expr \bullet ;]
           [stat \rightarrow \{ statlist \bullet \}]
                                                                                [stat \rightarrow ident = expr ; \bullet ]
S10:
                                                                    S28:
           [stat \rightarrow \{ statlist \} \bullet ]
S11:
```

LR(0)オートマトンの作成 (4)

S1	s3	read
S1 read S3	s24	ident
S1 read S3 ident S24	s25	;
S1 read S3 ident S24; S25	r5	print
S1 stat	s8	print
S1 stat S8	s4	print
S1 stat S8 print S4	s13	ident
S1 stat S8 print S4 ident S13	r11	;
S1 stat S8 print S4 expr	s22	;
S1 stat S8 print S4 expr S22	s23	;
S1 stat S8 print S4 expr S22; S23	r6	EOF
S1 stat S8 stat	s8	EOF
S1 stat S8 stat S8	r3	EOF
S1 stat S8 stat S8 stlist	s9	EOF
S1 stat S8 stat S8 stlist S9	r2	EOF
S1 stat S8 stlist	s9	EOF
S1 stat S8 stlist S9	r2	EOF
S1 stlist	s7	EOF
S1 stlist S7	r1 acc	EOF

図 4: LR 解析器による解析のようす

このLRオートマトンを用いた解析過程を手でシミュレートしてみましょう(図??)。プログラム例は再び「read x; print x;」です。

LR(0)オートマトンの作成 (5)

今度は道具立てとして、解析スタックと先読み記号に加えて「現在の状態」が必要なので、解析スタックには構文記号とLRオートマトンの状態を交互に載せます。まず状態S1から開始し、入力記号をread、ident、;、とシフトしながらS3、S24、S25 と遷移すると同時に、これらの記号と状態をそれぞれスタックに載せていきます。

次はS8でprint なのでS4に行き、次がidentでS13です。すると今度は先程と違い、ここで「行き止り」になってidentがexprに還元さてS4へ戻り、ここから改めてS22へ進み、ここでようやく;がシフトされてS23へ進み、print 文全体が還元されてS8へ戻ります。このように、「シフトすべきか、還元すべきか」という選択はオートマトンの状態を通じて的確に指示されるわけです。

このあとはもう次が EOF なので、 ϵ を statlist に還元し、それと前に認識した stat を併せて statlist に還元することを繰り返していき、S1 まで戻ったところで statlist により S7 へ行きます。次の動作は最初の規則 R1 による還元なので、これで終了となります。

SLR(1)解析器

先の説明では「行き止り」になったとき、●が最後にある項(「行き止り」だから必ずそういう規則がある)を用いて還元していました。

しかし、行き止りではないが●が最後に持つ項も含んでいる状態もあり得ますし、還元に使える項が複数ある可能性もあります。これらの場合の動作はオートマトンだけでは決まらず、別の情報を用いなければなりません。

1つの考え方として、選択をFollow集合に基づいて行うという方針が挙げられます。つまり、ある状態に [$A \to \alpha \bullet$] なる項が含まれ、かつ次の先読み記号t がFollow(A) に含まれている時のみ、この項に対応する規則による還元を行うわけです。

入力が構文的に正しいならばAの後に来る端記号はFollow(A)に含まれるはずなので、もしそうでないなら還元してはならないのは明らかです。この方針に基づく解析器をSLR(1)解析器と呼びます(最初のSはSimlpeの意)。

SLR(1)解析器 (2)

ident	=	;	()	{	}	<	>	eof	if	read	print	stlist	stat	cond	expr
s2					s6	r3			r3	s5	s3	s4	7	8		
	s26															
s24																
s13																22
			s12													
s2					s6	r3			r3	s5	s3	s4	9	8		
									accept							
s2					s6	r3			r3	s5	s3	s4	9	8		
						r2			r2							
						s11										
r8					r8	r8			r8	r8	r8	r8				
s13															14	15
		r11		r11			r11	r11								
				s20												
							s16	s17								
s13																18
s13																18
				r10												
				r9												
s2					s6					s5	s3	s4		21		
r7					r7	r7			r7	r7	r7	r7				
		s23														
r6					r6	r6			r6	r6	r6	r6				
		s25														
r5					r5	r5			r5	r5	r5	r5				
s13																27
		s28														
r4					r4	r4			r4	r4	r4	r4				
<u> </u>																
	s2 s24 s13 s2 s2 r8 s13 s13 s13 s13 s13 s2 r7 r6 r5 s13	s2 s26 s24 s13 s2 s2 r8 s13 s13 s13 s13 s13 s13 s13 s2 r7 r6 r5 s13	\$2	\$2 \$24 \$13 \$12 \$2 \$2 \$2 \$2 \$13 \$13 \$13 \$13 \$2 \$2 \$7 \$7 \$23 \$76 \$25 \$13 \$13 \$28	\$2 \$24 \$13 \$13 \$12 \$2 \$2 \$2 \$2 \$2 \$2 \$2 \$2 \$3 \$13 \$13 \$13 \$13 \$13 \$13 \$13 \$13 \$13	\$2	\$2	\$2	\$2 \$26 \$24 \$13 \$24 \$13 \$3 \$32 \$4 \$313 \$32 \$32 \$34 \$34 \$34 \$34 \$34 \$34 \$34 \$34 \$34 \$34	\$26 \$26 \$24 \$13 \$13 \$12 \$2 \$6 \$73 \$73 \$2 \$6 \$73 \$73 \$2 \$6 \$73 \$73 \$11 \$11 \$11 \$11 \$13 \$13 \$16 \$17 \$13 \$13 \$13 \$13 \$13 \$13 \$16 \$17 \$2 \$6 \$77 \$77 \$77 \$23 \$2 \$6 \$6 \$6 \$7 \$7 \$7 \$7 \$7 \$23 \$16 \$16 \$16 \$16 \$16 \$25 \$75 \$75 \$75 \$15 \$15 \$28 \$28 \$28 \$28 \$28 \$28 \$28	\$26 \$26 \$24 \$13 \$12 \$12 \$2 \$12 \$2 \$6 \$73 \$73 \$5 \$2 \$6 \$73 \$73 \$5 \$2 \$6 \$73 \$73 \$5 \$2 \$6 \$73 \$72 \$72 \$72 \$11	\$2 \$26 \$24 \$13 \$24 \$13 \$12 \$6 \$73 \$73 \$5 \$3 \$2 \$6 \$73 \$73 \$5 \$3 \$2 \$6 \$73 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$73 \$8	\$2 \$6 \$6 \$73 \$73 \$5 \$3 \$4 \$24 \$13 \$12 \$12 \$12 \$12 \$12 \$13 \$5 \$3 \$4 \$2 \$6 \$6 \$73 \$73 \$5 \$3 \$4 \$2 \$6 \$73 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$72 \$73 \$73 \$73 \$73 \$73 \$74 <td< td=""><td>\$2 \$6 r3 r3 \$5 \$3 \$4 7 \$24 \$13 \$12 \$6 r3 \$73 \$5 \$3 \$4 9 \$2 \$11 \$11 \$11 \$11 \$11 \$11 \$18</td><td>\$2 \$6 \$6 \$7 \$7 \$8 \$24 \$13 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$13 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$11</td><td>\$2 \$6 \$6 \$73 \$73 \$5 \$3 \$4 \$7 \$8 \$24 \$13 \$12 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$11 \$11 \$11 \$11 \$11 \$11 \$14<!--</td--></td></td<>	\$2 \$6 r3 r3 \$5 \$3 \$4 7 \$24 \$13 \$12 \$6 r3 \$73 \$5 \$3 \$4 9 \$2 \$11 \$11 \$11 \$11 \$11 \$11 \$18	\$2 \$6 \$6 \$7 \$7 \$8 \$24 \$13 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$12 \$13 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$11	\$2 \$6 \$6 \$73 \$73 \$5 \$3 \$4 \$7 \$8 \$24 \$13 \$12 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$2 \$6 \$73 \$73 \$5 \$3 \$4 \$9 \$8 \$11 \$11 \$11 \$11 \$11 \$11 \$14 </td

図 5: SLR(1) 構文解析表

LR解析器の動作は一般に、状態と先読み記号の対から「シフトしてどの状態へ行く」か「どの規則で還元する」のいずれかを指示する動作表 (action table) と、状態と非端記号の対から「還元した後どの状態に進むか」を指示する行き先表 (goto table) の組からなる LR 構文解析表 LR parsing tableで表せます。先の例に対応する解析表を図??に示します。

先の図??の動作は、この解析表に従って解析した場合を反映し

たものです。どこにFollow集合の情報が使われていたか分かりますか? それは、 $statlist \rightarrow \epsilon$ による還元を行なうのは先読み記号がEOFか $\}$ のいずれかの時だけ、という部分です。この制約がないと、空列はいつでもstatlistに還元できてしまうので、それをどこで適用すべきか分からないことになります。

SLR(1)解析器 (3)

図??についてもう1つ補足です。先にはわかりやすさのためスタック上に構文記号と状態を交互に積むとしていましたが、実際には構文記号は必要としません。その理由は、オートマトンの状態には「現在どのような入力列を読んだところか」という情報が必要なだけ含まれているためです。このため、実際には解析スタック上には状態だけ積めば十分です。実際、動作表ではシフト時には「どの状態に進む」かは直接示されていますし、還元時には還元された記号が分かるのでそれを状態と照らし併せて次の状態が決められます。

正準LR(1)解析器とLALR(1)解析器

SLR(1)構文解析では還元の選択を Follow 集合によっていますが、Follow はもともと構文全体を通してどこかで後続記号になっているかを見るものなので、細かい文脈の区別には不向きです。

例えば次のような文法を考えてみます。

```
prog ::= stat
stat ::= left = right \mid right
left ::= ident \mid * right
right ::= left
```

これはCのように「式は単独でも文になり得、また代入の左辺はポインタ参照でもよい」ような言語の文法を簡略化したものになっています。これから先のプログラムによりSLR(1)構文解析器をつくろうとすると、leftを認識したあとに「=」が先読み記号になった場合、シフトすべきか $right \rightarrow left$ で還元すべきかが判別できません。これは、right =の前に来ることが可能なためFollow(right)に=が含まれているからです。しかし実際には代入記号の直前にいるのだから、leftをrightに還元するのは誤りです。

このような *Follow* による「おおざっぱな判定」の弱点を克服するため、最初から状態の中に先読み記号も含めて計算を行うことを考えます。そこで、今度は項の中に次のようにして先読み記号を含めることにします。

```
[ prog \rightarrow \bullet stat ; EOF ]
[ right \rightarrow left \bullet ; = ]
```

正準LR(1)解析器とLALR(1)解析器 (2)

このような項をLR(1) 項と呼びます(1というのは先読み記号の長さを示す)。そしてあとは、この項どうしの遷移により、これまで通りに解析表を作ります。これを正準(cannonical)LR(1)構文解析表、それに従って動作する解析器を正準LR(1)解析器と呼びます。

正準 LR(1) 文法は SLR(1) よりも広い文法クラスを含んでいます。実は、正準 LR(1) 文法は左から右へ後戻りなしで解析可能な最大の文法クラスとなっていることが知られています。

一方、正準LR(1)解析器の問題点はその状態数が非常に多くなるということです。それは文法における端記号の数がnとして、項の数がSLR(1)のn倍になっていることに起因します。このため、正準LR(1)解析器をそのままコンパイラに採用するのは実用的でないとされます。

実用のコンパイラでは、SLR(1)より受け入れられる文法の範囲が広く、しかもLR(1)ほどには状態数が多くない実用的な構文解析器として、LALR(1)解析器が多く使われています。

正準LR(1)解析器とLALR(1)解析器 (3)

これは、LR(1)でオートマトンを構成したあと、先読み記号部分(;の後ろ)を取り除いて同じになる状態を互いに併合することで構成されます(併合をおこなうときに行き先となる状態が食い違って併合できないなどのことは起きないことが知られています)。

併合をおこなった後の状態数はLR(0) オートマトンと同じになるため、状態数の問題が解消され、なおかつ先読み記号の情報が使えるため、SLR(1) よりは広いクラスの文法が扱えます(ただし併合により正確さが減じるため、正準LR(1) よりは狭くなる)。

なお、解析表の生成時だけとは言え、正準LR(1)オートマトンをいったん作るのでは計算の手間が大きくなりますがLR(1)オートマトンを経ずに直接必要な状態のみを生成しながら合せて先読み記号の情報を計算する手順が知られています。

- 演習 8-1 例題のオートマトンや構文解析表で別のプログラムを与えたとき確かに解析ができることを確認してみなさい。解析のようすを記録すること。
- 演習 8-2 簡単な言語を BNF で記述し、その LR(0) オートマトンを構成してみなさい。 SLR(1) 構文解析表までできるとなおよい。これらを使って簡単なプログラムを解析するようすを示すこと。

曖昧な文法の活用

ここまでは曖昧でない文法のみを想定して来ましたが、場合によっては曖昧な文法を許す方が人間、コンピュータ双方にとって有利なこともあります。その代表的な例としてぶらさがりelse(dangling else)があります。多くの言語ではif文に対して次の構文定義を用いています。

```
stat \rightarrow \mathbf{if} \ cond \ \mathbf{then} \ stat \ \mathbf{else} \ stat
stat \rightarrow \mathbf{if} \ cond \ \mathbf{then} \ stat
```

この場合、if C1 then if C2then A else Bという文は次の 2通りに解釈でき、曖昧です。

```
if C1 then ( if C2 then A else B ) -- (1) if C1 then ( if C2 then A ) else B -- (2)
```

そこで「ただし、おのおのの else はまだ対応する else をもたない最も近い if に対応させる」という規則を (自然言語で)付すのが恒例です (これにより上の解釈 (2) は削除される)。

なお、この例の場合は代りに文法を次のように直しても曖昧さ は解消できますが、明らかに繁雑になります。

```
balstat \rightarrow \mathbf{if} \ cond \ \mathbf{then} \ balstat \ \mathbf{else} \ stat — 各種の文 stat \rightarrow \mathbf{if} \ cond \ \mathbf{then} \ stat — balstat
```

このように文法が曖昧なまま構文解析器を作ろうとすると、LR 構文解析の場合にシフト還元衝突 (shift-reduce conflict) や還 元還元衝突 (reduce-reduce conflict) が発生します。つまり、 構文解析表を構成しているときに、特定の位置の動作表に「シ フト」「還元」の両方の動作が記入できることになったり、複数 の規則による「還元」が記入できるようになることを言います。

曖昧な文法の活用(2)

たとえば、ぶらさがり else の場合、文法の曖昧さに対応して「then 部だけから成る if 文を全て読み終った時点とも、else 部まである if 文の else の直前とも取れる」状態が現れ、そこでシフト還元衝突となります。

そこで、このような状態では常にシフトする—つまり、後者の解釈を取る—ことにして解析表を作っていまうことができます。これによって構成される解析器は上述の「elseを最も近いifに対応させる」解析動作を行います。LL(1)解析器の場合でもこれと同様なことが行えます。

曖昧な文法が役に立つもう1つの代表例は式と演算子の構文です。「 $expr \rightarrow expr + expr$ 」という文法は「1 + 2 + 3」を「(1 + 2) + 3」とも「1 + (2 + 3)」とも解釈できます。

この曖昧さを除くには、これまでに見てきたように、式-項-因子のように演算子の強さごとに別の非端記号を使用すればよいわけですが、これも読みやすくはありません。ここでは+は左結合的なので、前者を選択し、還元を行うように解析表を作ればよいわけです。同様に、演算子間の順位や右結合なども、そのことを予め考慮して解析表を作ることで対応できるわけです。

曖昧な文法を使用することの利点は1つは、構文記号の数が少なく人間にとってわかりやすくなることですが、同時に解析器の側でも状態数が少なくてすみ、また余分な還元がなくなるので効率よく解析が行えるという利点があります。

構文解析器生成系

Cupとその構文定義

解析表とドライバが分かれた形の構文解析器では、コンパイラ 作成者が自ら解析表の計算を行うことはほとんどなく、構文記 述を入力すると解析表を作成してくれるツールを使用するのが 普通です。これを構文解析器生成系 (parser generator) と呼び ます。

たとえば、Unixに古くから備わっているYaccやそのフリーソフト版であるBisonはLALR(1)解析器を生成する生成系で、出力はC言語のコードになっています。ここでは言語としてJavaを使っていて、JFlexとの相性もよいことから、Cupと呼ばれる生成系を使ってみます (細かい書き方はYaccと違っていますが、大筋は同じです)。

まず先頭の import は CUP のライブラリにアクセスするため 必須のもので、そのままコードの先頭の入ります。そのあと、端記号および非端記号として使う名前を宣言します。ここでは端記号を大文字、非端記号を小文字にしています。

その次ですが、上にのべたように曖昧な文法を使うため、優先順位の低いものから演算子の端記号を列挙します。また、nonassoc、left、rightにより結合のしかた (非結合、左結合、右結合)を指定します。左結合は「 $x \diamondsuit y \diamondsuit z \to (x \diamondsuit y) \diamondsuit z$ 」の意味で、非結合は「 $x \diamondsuit y \diamondsuit z$ 」のようにつなげては書けないことを意味します (比較演算子は通常そうですね)。なお、最後のUMINUSという端記号は最も優先順位が高くなりますが、実際には現れません。これについてはすぐ後で述べます。

Cup とその構文定義 (2)

```
import java_cup.runtime.*;
terminal READ, PRINT, IF, WHILE, ASSIGN, SEMI, LPAR, RPAR, LE
terminal GT, LT, PLUS, MINUS, ASTER, SLASH, UMINUS;
terminal ICONST, IDENT;
non terminal prog, statlist, stat, expr;
precedence nonassoc LT, GT;
precedence left PLUS, MINUS;
precedence left ASTER, SLASH;
precedence left UMINUS;
prog ::= statlist
statlist ::= stat statlist
stat ::= IDENT ASSIGN expr SEMI
       | READ IDENT SEMI
       | PRINT expr SEMI
       | IF LPAR expr RPAR stat
       | WHILE LPAR expr RPAR stat
       | LBRA statlist RBRA
expr ::= expr PLUS expr
       | expr MINUS expr
```

```
| expr ASTER expr
| expr SLASH expr
| expr GT expr
| expr LT expr
| MINUS expr %prec UMINUS
| IDENT
| ICONST
| LPAR expr RPAR
```

さて文法ですが、見て分かるとおり普通のBNFで、ただし上述のように式のところは曖昧な文法を使っています(なので簡潔です)。「MINUS exp %prec UMINUS」というところが謎ですが、これは単項のマイナス演算子で、端記号としては MINUS が使われますが、ただし順位は他の演算よりも高い UMINUS の強さにします、という指定です。

演習 8-3 上の例題を Cup で処理してみなさい。うまくいったら、 自分独自の簡単な言語を BNF で定義し、次に Cup で構文記 述を書いて変換してみなさい。エラーがあれば直すこと。

JFlex と Cup の連携

上記の文法指定をcupコマンドで処理すると、parser.javaとsym.javaという2つのファイルが生成されます。前者はパーサのファイルですが、後者はToken.javaのように端記号の番号を定義するためのもので、使い方も同様です。これを用いてこの言語用のJFlexファイルを作成します。

まず冒頭でCupと同じライブラリのimportを行ないます。これは、字句解析器が返すものがSymbolオブジェクトである必要があるためです。「%cup」という指定は字句解析用のメソッド名をCup用の名前にするために必要です。次の3行は、Symbolオブジェクトを生成するためのオブジェクトを用意するものです。%eofval の3行はEOFになったときにEOFシンボルを返す指定です。

JFlexとCupの連携 (2)

```
import java_cup.runtime.*;
%%
%class Lexer
%cup
%{
  SymbolFactory sf = new DefaultSymbolFactory();
%}
%eofval{
  return sf.newSymbol("EOF", sym.EOF);
%eofval}
L = [A-Za-z]
D = [0-9]
Ident = \{L\}(\{L\}|\{D\})*
Iconst = [-+]?{D}+
Blank = [ \t \n] +
%%
          { return sf.newSymbol("PLUS", sym.PLUS); }
\+
          { return sf.newSymbol("MINUS", sym.MINUS); }
\-
          { return sf.newSymbol("ASTER", sym.ASTER); }
\*
\/
          { return sf.newSymbol("SLASH", sym.SLASH); }
          { return sf.newSymbol("LPAR", sym.LT); }
\<
          { return sf.newSymbol("RPAR", sym.GT); }
\>
\(
          { return sf.newSymbol("LPAR", sym.LPAR); }
\)
          { return sf.newSymbol("RPAR", sym.RPAR); }
          { return sf.newSymbol("LPAR", sym.LBRA); }
\{
\}
          { return sf.newSymbol("RPAR", sym.RBRA); }
          { return sf.newSymbol("SEMI", sym.SEMI); }
\;
```

JFlexとCupの連携 (3)

その後は基本的にこれまでのJFlex と同じですが、ただし返す値は上述のように Symbol オブジェクトを返すようにします。そのときのパラメタとして、1番目は表示用の文字列、2番目がトークンの値 (クラス sym で定義されているものを使います)、3番目は Ident と Const についてのみ、値としてトークンの文字列を渡しています。これは後で使います。

main は次のようになります。ここでは何も動作を指定していないので、構文エラーがあればエラーが出るというだけです(つまり認識器です)。

JFlexとCupの連携 (4)

```
import java.util.*;
import java.io.*;
import java_cup.runtime.*;
public class Sam81 {
 public static void main(String[] args) throws Exception {
   parser p1 = new parser(new Lexer(new FileReader(args[0]))
   p1.parse();
 }
}
 なおComplexSymbolFactoryはCUPランタイムが提供する記
号用クラスで、これを標準で渡すことになっています。
 実行のようすを一応のせましょう。
% cup sam81.cup
(Cupのメッセージ)
% jflex sam81.jflex
(JFlexのメッセージ)
% javac Sam81.java
% java Sam81 test.min ←後述
            ←出力なし: エラーがないことは分かる
%
```

属性とアクション

前節まででは、構文解析はしても何も動作がついていないので、 出力が何もありませんでした。Cupで(YaccやBisonもそうで すが)実際にコンパイラを作るにはどうするのでしょう。それに は、次の2つの道具だてがあります。

- (a) それぞれの構文記号(端記号、非端記号)に属性(attribute) をつけられる。
- (b) 構文規則に動作 (action) を付随させられる。
- (a)については、記号ごとにその記号が値を伴うことができる、ということです。これは実装としては、構文解析に使うスタックと並行してもう1つ意味スタック (semantic stack)というスタックを用意し、構文記号 (実際には対応する状態)を入れるのと同じ場所の意味スタック側に属性の値を格納します。Cupでは属性値はJavaのオブジェクト型である必要があります。

たとえば今回は、IDENTとICONSTは字句解析からの文字列(既にyytext()の値を渡すように作ってありました)を属性値とし、端記号はすべてTree.Nodeオブジェクトを属性値とすれば、木構造を組み立てることができます。

(b) については、それぞれの構文規則の右端に「 $\{: \cdots:\}$ 」という形で囲んだJavaコードを記述しておくと、そのコードはその構文規則が還元されるときに実行されます。上向き解析なので、下から順に木のノードを組み立てて行くのが自然な使い方です。

属性とアクション (2)

しかし、子の構文規則で組み立てたノードを親側から取り出すのにはどうすればいいでしょうか? そこで属性値を使います。 Cupでは構文規則の右辺に現れる端記号、非端記号には「:名前」という指定がつけられ、アクション中ではその名前の変数が、対応する記号の属性値を保持しています。そして、RESULTという特別な名前に代入したものが、その構文規則の左辺の記号の属性値となります。

では具体例を見てみましょう。文法は先の例と同じですが、属性とアクションが追加されています。まず属性の型として、IDENTとICONSTはString、非端記号はTree.Nodeを指定しています。そして、規則の右辺で属性を受け渡す記号は:xとか:yなど、受け渡す変数を指定しています。

progの属性値はstatlistの属性値そのままです。statlistは、空に対応するときは空のTree.Seqを作り、そのあと1つ文が現れるごとにその文のノードを追加し、値としてはTree.Seqを受け渡して行きます。

属性とアクション (3)

```
import java_cup.runtime.*;
terminal READ, PRINT, IF, WHILE, ASSIGN, SEMI, LPAR, RPAR, LE
terminal GT, LT, PLUS, MINUS, ASTER, SLASH, UMINUS;
terminal String ICONST, IDENT;
non terminal Tree. Node prog, statlist, stat, expr;
precedence nonassoc LT, GT;
precedence left PLUS, MINUS;
precedence left ASTER, SLASH;
precedence left UMINUS;
                            \{: RESULT = x; :\}
prog ::= statlist:x
statlist ::= statlist:x stat:y {: RESULT = x; x.add(y); :}
                             {: RESULT = new Tree.Seq(); :}
stat ::= IDENT:x ASSIGN expr:y SEMI
        {: RESULT = new Tree.Assign(new Tree.Var(x), y); :}
       | READ IDENT:x SEMI {: RESULT = new Tree.Read(new Tree)
       | PRINT expr:x SEMI {: RESULT = new Tree.Print(x); :}
       | WHILE LPAR expr:x RPAR stat:y {: RESULT = new Tree.V
       | LBRA statlist:x RBRA {: RESULT = x; :}
expr ::= expr:x PLUS expr:y {: RESULT = new Tree.Add(x, y);
```

```
| expr:x MINUS expr:y {: RESULT = new Tree.Sub(x, y);
| expr:x ASTER expr:y {: RESULT = new Tree.Mul(x, y);
| expr:x SLASH expr:y {: RESULT = new Tree.Div(x, y);
| expr:x GT expr:y {: RESULT = new Tree.Gt(x, y); :
| expr:x LT expr:y {: RESULT = new Tree.Lt(x, y); :
| MINUS expr:x {: RESULT = new Tree.Neg(x); :}
| IDENT:x {: RESULT = new Tree.Var(x); :}
| ICONST:x {: RESULT = new Tree.Lit(Integer.parseInt(x); :}
| LPAR expr:x RPAR {: RESULT = x; :}
| ;
```

属性とアクション (4)

それぞれの文や式はこれまでやってきたように、各種のノードを作って自分の値として返すだけです。そのとき子ノードの値は、構文規則の属性として受け渡されてくる値を使えばよいわけです。このようにして、パーサジェネレータを使うことでコンパクトに動作のついたパーサを構成し、抽象構文木を組み立てることができるわけです。

最後にmain()ですが、組み立てた木構造はparse()が返すSymbol オブジェクトの変数 value に入っているので、それを Tree. Node 型にキャストして eval() することで実行できます。

```
import java.util.*;
import java.io.*;
import java_cup.runtime.*;

public class Sam82 {
   public static void main(String[] args) throws Exception {
     parser p1 = new parser(new Lexer(new FileReader(args[0]))
        Symbol s1 = p1.parse();
        ((Tree.Node)s1.value).eval();
   }
}
```

属性とアクション (5)

では、最大公約数のプログラムを実行してみましょう。

```
% cat test.min
read x; read y;
while(x - y) {
   if(x > y) { x = x - y; }
   if(x < y) { y = y - x; }
}
print x;
% java Sam82 test.min
x? 60
y? 18
6
%</pre>
```

- 演習8-4 例題の文法のif文にelse部をつけてみなさい。曖昧な文法と曖昧でない文法の両方でできるとなおよい (JFlexでELSE を追加する必要。曖昧な文法は Cup のマニュアルをよく読む必要があるかも。低い順位で「right ELSE」を設定するとよいかも)。
- 演習 8-5 例題の文法の if 文や while 文を end のある形に直して みなさい。if 文については、if-elsif-elsif-else-end の構文も きちんとサポートすること。
- 演習8-6 例題に自分の好きな言語機能(制御構造でもそれ以外でも)を追加してみなさい。
- 演習8-7 実行のされかたは例題と同じだが、書き方の見た目がまったく違う構文になっている言語を実装してみなさい。

課題 8A

今回の演習問題から(小問を)1つ以上選び、プログラムを作成しなさい。作成したプログラムについてレポートを作成し、久野(y-kuno@uec.ac.jp)までPDFを送付してください。LaTeXの使用を強く希望します。レポートは次の内容を含むこと。期限は次回授業前日一杯。レポートおよびその評点はクラス内で公開します。

- タイトル 「システムソフトウェア特論 課題 # 7」、学籍番号、氏名、提出日付。
- ●課題の再掲 レポートを読む人がどの課題をやったのか分かる程度にやった課題を要約して説明してください。
- 方針 ― その課題をどのような方針でやろうと考えたか。
- 成果物 ― プログラムとその説明および実行例。
- 考察 ― 課題をやってみて分かったこと、気付いたことなど。
- 以下のアンケートの解答。
 - Q1. 上向き解析の原理は分かりましたか。LR オートマトン の構成は無理としても、オートマトンができたらそれを 使って動作するやり方が分かればよいです。
 - Q2. cup を使ってパーサを構築するのはどうでしたか。
 - Q3. リフレクション (課題をやってみて気付いたこと)、感想、 要望など。