

CALCULO EFICIENTE DE CONJUNTOS DE LOOKAHEADS LALR(1) EN MICROCOMPUTADORES

Arturo Montes Sinning
Rodrigo López Beltrán
Universidad de los Andes
Bogotá - Colombia

1. Introducción

Durante el primer semestre de 1983 se comenzó un proyecto de software, en el departamento de sistemas de la Universidad de los Andes, cuyo objetivo era producir un generador de analizadores sintácticos LALR para ser utilizado como herramienta didáctica en el curso de construcción de compiladores de nuestro currículum. Dadas las restricciones de uso del equipo del Centro de Cálculo de la universidad (en donde ya se disponía de un generador de ese estilo), era necesario implementar el sistema sobre los microcomputadores que soportan buena parte de la atención a estudiantes. Una primera versión [MON85], que se ha venido utilizando desde comienzos de 1985, funciona sobre microcomputadores con sistema operacional MS-DOS.

Las ideas básicas para el desarrollo del proyecto fueron tomadas del trabajo de DeRemer y Pennello[DeR82], pero el resultado final contiene algunas mejoras con respecto al cálculo de los LookAheads (Algoritmo 3.6), aparte de ciertas ideas acerca de información útil para la resolución de los eventuales conflictos que se presentan en una gramática (Sección 4). El documento original de este trabajo [MON85], no resalta claramente las virtudes antes mencionadas pues se pierde en lo intrincado del formalismo utilizado, que es el propuesto en [DeR82]. Este artículo aclara las ideas de [MON85] gracias al excelente formalismo expuesto por Park, Chang y Choe [PAR85]. En las secciones 2 y 3 se exponen los algoritmos de cálculo del autómeta LR(0) y los LookAheads, respectivamente, junto con la

introducción del formalismo "a la Park". En la sección 4 se esbozan las ideas con respecto a la resolución de conflictos y en la sección 5 se presentan algunas conclusiones.

2. Generación del Autómata LR(0)

Para la mayoría de las definiciones que aparecen a continuación utilizaremos la terminología habitual (ver [AHO77]). Algunas, que son la base de los resultados presentados en el artículo, han sido tomadas de [PAR85] y [DeR79].

Definición 2.1

Las nociones de símbolo y de cadena se suponen conocidas. Un vocabulario V es un conjunto de símbolos. V^* denota el conjunto de todas las cadenas de símbolos de V . V^+ denota $V^* - \{ \lambda \}$, donde λ es la cadena nula. La longitud de una cadena alfa la denotaremos por $|alfa|$. El primer símbolo de una cadena alfa lo denotaremos por $\text{Prim}(alfa)$; la cadena que sigue al $\text{Prim}(alfa)$ será $\text{Rest}(alfa)$ (el resto de alfa); el último símbolo lo designaremos por $\text{Ult}(alfa)$.

Definición 2.2

Si R es una relación, R^* denota la clausura reflexiva y transitiva de R , y R^+ denota la clausura transitiva de R .

Definición 2.3

Una gramática G , es una cuádrupla $G = \langle N, T, S, P \rangle$ donde:

- . N es un conjunto de símbolos llamados no-terminales.
- . T es un conjunto de símbolos llamados símbolos terminales.
- . T es disjunto de N , $S \in N$, y $V = N \cup T$.
- . P es un subconjunto de $N \times V^*$ llamado Producciones de G . Si $(A, w) \in P$ escribiremos $A \rightarrow w$ donde A se denomina la parte izquierda y w la parte derecha de una producción.

A continuación establecemos las siguientes convenciones que serán usuales en este artículo a menos que se especifique lo contrario:

S, A, B, C, \dots	\in	N
X	\in	V
t, a, b, c, \dots	\in	T
$\dots x, y, z$	\in	T^*
<u>alfa</u> , <u>beta</u> , <u>gamma</u> ..	\in	V^* .

Definición 2.4

\Rightarrow_d es una relación de V^* en V^* tal que si alfa $\in V^*$ y $z \in T^*$, entonces alfa $Az \Rightarrow_d$ alfawz si y sólo si $A \rightarrow w \in P$. Al subíndice d lo eliminaremos de la notación, y cuando aparezca el símbolo \Rightarrow entenderemos que se trata de \Rightarrow_d . \Rightarrow^* y \Rightarrow^+ denotan la clausura reflexiva y transitiva de la relación \Rightarrow y ambas se leerán "deriva".

Definición 2.5

A es un no-terminal anulable si y sólo si, $A \Rightarrow^* \lambda$.

Definición 2.6

alfa es una forma sentencial si y sólo si, $\text{alfa} \in V^*$ y $S \Rightarrow^* \text{alfa}$. Se llamará una sentencia si $\text{alfa} \in T^*$.

Definición 2.7

$L(G)$, el lenguaje generado por G , es el conjunto de sentencias. En otras palabras, $L(G) = \{ x \in T^* \mid S \Rightarrow^+ x \}$.

Definición 2.8

G es una gramática reducida, si y sólo si, para toda forma sentencial $S \Rightarrow^+ \text{alfa}\beta$, existe $y \in T^*$ tal que $A \Rightarrow^* y$.

Definición 2.9

Para una gramática arbitraria G , podemos construir una gramática aumentada G' que contiene la producción $S' \rightarrow S\#$ en donde $S' \notin N$ y $\#$ es un nuevo símbolo terminal. Ahora $L(G') = L(G)\#$ y además se garantiza que el símbolo S' no aparece en la parte derecha de ninguna producción.

En lo sucesivo, supondremos que todas las gramáticas son reducidas y aumentadas.

Definición 2.10

Una tripleta (A, alfa, β) es un ítem, si y sólo si, existe $A \rightarrow \text{alfa} \beta \in P$, donde $(A, \text{alfa}, \beta) \in N \times V^* \times V^*$. El ítem (A, alfa, β) será denotado por $A \rightarrow \text{alfa} \beta$. Si $\text{alfa} = \lambda$ lo llamaremos un ítem inicial y si $\beta = \lambda$ lo llamaremos un ítem final y si alfa y β son λ lo denominaremos un ítem nulo.

Definición 2.11

Una Tabla es un conjunto de ítems.

Definición 2.12

Si S es una Tabla,

$$\text{CLAUSURA}(S) = S \cup \{ A \rightarrow \cdot w \mid B \rightarrow \text{alfa} \cdot A \beta \in \text{CLAUSURA}(S) \}.$$

Definición 2.13

Si S es una Tabla y X es elemento de V ,

$$\text{GOTO}(X, S) = \text{CLAUSURA}(\{ A \rightarrow \text{alfa} X \cdot \beta \mid A \rightarrow \text{alfa} \cdot X \beta \in S \}).$$

Definición 2.14

Si S es una Tabla, $KERNEL(S)$ es el mínimo K , subconjunto de S , tal que $CLAUSURA(K) = CLAUSURA(S)$.

Definición 2.15

Sea S una Tabla, $SUCESORES(S) = \{ GOTO(X, CLAUSURA(S)) \mid X \in V \}$.

Definición 2.16

$TP(G)$, el conjunto de tablas de parsing LR(0) de una gramática G , es

$$TP(G) = \{ CLAUSURA(\{S \rightarrow .S'\#\}) \} \cup \{ CLAUSURA(K) \mid K \in SUCESORES(K') \text{ y } K' \in TP(G) \}.$$

Definición 2.17

Sea $A \rightarrow \underline{\text{alfa.Xbeta}}$ un ítem. La función GEN se define como,

$$GEN(A \rightarrow \underline{\text{alfa.Xbeta}}) = \emptyset \text{ si } X \in T \\ \text{ó} \\ A \rightarrow \underline{\text{alfa.Xbeta}} \text{ es un ítem final.}$$

$GEN(A \rightarrow \underline{\text{alfa.Xbeta}}) = \{ X \rightarrow .w \mid X \rightarrow w \in P \}$ en cualquier otro caso. Por extensión, si S es una Tabla, entonces

$$GEN(S) = \cup \{ GEN^*(A \rightarrow \underline{\text{alfa.Xbeta}}) \mid A \rightarrow \underline{\text{alfa.Xbeta}} \in S \}.$$

Lema 2.1

Sea $A \rightarrow \underline{\text{alfa.Xbeta}}$ un ítem,

$$CLAUSURA(\{A \rightarrow \underline{\text{alfa.X beta}}\}) = GEN^*(A \rightarrow \underline{\text{alfa.Xbeta}}).$$

Del lema y de la definición anteriores se observa que el cálculo de la $CLAUSURA$ de una tabla, se puede realizar calculando GEN^* . Por otro lado, GEN nos sugiere una relación (que llamaremos I) entre los no-terminales:

$$B I C \text{ si y solamente si } B \rightarrow C_{\text{gamma}} \in P.$$

Al grafo dirigido asociado con la relación I lo llamaremos el I -grafo.

Lema 2.2

Sea $A \rightarrow \underline{\text{alfa.Bbeta}}$ un ítem, entonces

$$GEN^*(A \rightarrow \underline{\text{alfa.Bbeta}}) = \{ C \rightarrow .w \mid B I^* C \}.$$

Definición 2.18

El Autómata LR(0) para una gramática G (que en adelante

llamaremos $ALR(0)$ es una cuádrupla $ALR(0) = (Q, i^{\circ}, SUCC, TP(G))$ donde

- . $TP(G)$ es el conjunto de tablas de parsing de G .
- . Q es un conjunto de índices de una enumeración de $TP(G)$.
- . i° es el índice del primer elemento de la enumeración de $TP(G)$ y corresponde a $CLAUSURA(\{S' \rightarrow .S\#})$.
- . $SUCC$ es una función de $Q \times V$ en $TP(G)$ tal que

$$SUCC(p, X) = GOTO(X, Cp)$$

donde Cp denota el p -ésimo elemento en la enumeración de $TP(G)$. De ahora en adelante cuando referenciamos un elemento Cp de $TP(G)$ nos estaremos refiriendo al p -ésimo elemento de la enumeración de $TP(G)$ y con Kp al $KERNEL(Cp)$. Cuando no se introduzca ambigüedad, utilizaremos un índice p como sinónimo de Cp .

El siguiente algoritmo nos calculará $TP(G)$ para una gramática G cualquiera:

Algoritmo 2.1

```

Begin
  TP(G) := {CLAUSURA({S' -> .S#})} ;
  For Each K ∈ TP(G) tal que sus SUCESORES no hayan sido
    calculados dO_Loop
    For Each K' ∈ SUCESORES(K) tal que K' ∉ TP(G) dO_Loop
      TP(G) := TP(G) U { K' } ;
    End_Loop ;
  Marcar a K indicando que sus SUCESORES ya han sido
    calculados ;
  End_Loop ;
End ;

```

Si guardamos en una cola los elementos de $TP(G)$ cuyos sucesores no han sido calculados, obtenemos una primera mejora y el algoritmo quedaría de la siguiente forma:

Algoritmo 2.2

QK : Es una cola de los $KERNELs$ de los elementos de $TP(G)$ cuyos $SUCESORES$ no han sido calculados.

```

Begin
  QK := <{S' -> .S#}> ;
  TP(G) := {CLAUSURA({S' -> .S#})} ;
  WHILE Not(Empty(QK)) dO_Loop
    K := DeQueue(QK) ;
    For Each K' ∈ SUCESORES(K) tal que K' ∉ TP(G) dO_Loop
      TP(G) := TP(G) U { K' } ;
      EnQueue(K', QK) ;
    End_Loop ;
  End_Loop ;
End ;

```

Se observa, que la operación de averiguar si una Tabla K se encuentra en $TP(G)$, es determinante en la complejidad del

algoritmo. Por lo tanto, respondiendo a ella de una manera eficiente podemos minimizar el costo del cálculo de la tabla de parsing.

Definición 2.19

Sea $A \rightarrow \text{alfa.Xbeta}$ un ítem, se define la función ESTADOS como,

$$\text{ESTADOS}(A \rightarrow \text{alfa.Xbeta}) = \{ p \in Q \mid A \rightarrow \text{alfa.Xbeta} \in K_p \}.$$

donde K_p es el KERNEL del p -ésimo elemento en la enumeración de $TP(G)$.

Lema 2.3

Sea $q \in Q$, entonces

$$Kq = \{ A \rightarrow \text{alfa.Xbeta} \mid q \in \text{ESTADOS}(A \rightarrow \text{alfa.Xbeta}) \}.$$

Definición 2.20

Sea C_p una Tabla de $TP(G)$ y $q \in Q$. Diremos que q es un estado candidato de C_p si y solamente si Kq es un superconjunto de S .

Lema 2.4

$q \in Q$ es un estado candidato de S , si y sólo si,

$$\text{Para todo } A \rightarrow \text{alfa.Xbeta} \in S, q \in \text{ESTADOS}(A \rightarrow \text{alfa.Xbeta}).$$

Por el lema anterior tenemos que el costo de responder a la pregunta $P \in TP(G)$ es lineal con respecto a $\#\{q \mid q \text{ es un estado candidato de } P\}$, que es lo mismo que $\#\{q \mid Kp \subset Kq\}$.

3. Cálculo de los LookAheads

Definición 3.1

Sea $q \in Q$ y $\text{alfa} \in V^*$, se extiende SUCC a cadenas de V^*

$$\text{SUCC}(q, \text{alfa}) = q \text{ si } \text{alfa} \text{ es } \lambda,$$

$$\text{SUCC}(q, \text{alfa}) = \text{SUCC}(\text{SUCC}(q, \text{Prim}(\text{alfa})), \text{Rest}(\text{alfa})).$$

Definición 3.2

Sea $q \in Q$, $A \in N$. La pareja (q, A) es una "transición no terminal" si y sólo si $\text{SUCC}(q, A)$ está definido.

Definición 3.3.

Sea $q \in Q$ y $A \rightarrow \text{alfa.beta}$ un ítem,

$$LA(q, A \rightarrow \underline{\text{alfa}}.\underline{\text{beta}}) = \{ t \in T^+ \mid S^* \Rightarrow^* \underline{\text{gamma}}A\underline{\text{sigma}}^* \Rightarrow^* \underline{\text{gamma}} \underline{\text{alfa}} \underline{\text{beta}} \underline{\text{sigma}}^* \Rightarrow^* \underline{\text{gamma}} \underline{\text{alfa}} \underline{\text{beta}} t \Rightarrow^* x, \\ x \in L(G) \text{ y } \text{SUCC}(i^*, \underline{\text{gamma}}, \underline{\text{alfa}}) = Cq \}.$$

Definición 3.4

Sea $q \in Q$ y $A \rightarrow w \in P$,

$$LA(q, A \rightarrow w) = LA(q, A \rightarrow w.).$$

Definición 3.5

Sea $q \in Q$ y $t \in T$,

$$\text{Reduce}(q, t) = \{ A \rightarrow w \in P \mid t \in LA(q, A \rightarrow w) \}.$$

Definición 3.6

Sea $\underline{\text{alfa}} \in V^*$,

$$\text{FIRST}(\underline{\text{alfa}}) = \{ t \in T \mid \underline{\text{alfa}} \Rightarrow^* tw \text{ para algún } w \in T^* \}.$$

Si $\underline{\text{alfa}}$ es anulable entonces $\underline{\text{lambda}}$ también pertenece a $\text{FIRST}(\underline{\text{alfa}})$.

Definición 3.7

Sean L, M conjuntos de cadenas,

$$L \otimes M = L \text{ si } \underline{\text{lambda}} \text{ no pertenece a } L \\ \delta \\ L \otimes M = L \cup M \text{ en caso contrario.}$$

Definición 3.8

Sean P y $C - N$,

$$\text{PATH}(P, C) = \cup \{ \text{FIRST}(\underline{\text{beta}}_n \dots \underline{\text{beta}}_2 \underline{\text{beta}}_1) \mid B_0 = B, B_n = C, \\ B_0 \rightarrow B_1 \underline{\text{beta}}_1 \in P, B_1 \rightarrow B_2 \underline{\text{beta}}_2 \in P, \dots, \\ B_{n-1} \rightarrow B_n \underline{\text{beta}}_n \in P \}.$$

Definición 3.9

Sea $p \in Q$ y $\underline{\text{alfa}} \in V^*$,

$$\text{PRED}(p, \underline{\text{alfa}}) = \{ q \in Q \mid p = \text{SUCC}(q, \underline{\text{alfa}}) \}.$$

Lema 3.1 [PAR85]

Sean $p, q \in Q$ tales que $q \in \text{PRED}(p, \underline{\text{alfa}}_1)$,
 $A \rightarrow \underline{\text{alfa}}_1.\underline{\text{alfa}}_2 \in C_p$, $t \in LA(p, A \rightarrow \underline{\text{alfa}}_1.\underline{\text{alfa}}_2)$, $A \neq S$.
 Entonces debe haber por lo menos un ítem $B \rightarrow \underline{\text{beta}}_1.A'\underline{\text{beta}}_2 \in K_q$
 que satisface la siguiente condición:

$$t \in \{ a \mid a \in \text{PATH}(A', A) @ \text{FIRST}(\text{beta}2t), \\ q \in \text{PRED}(p, \text{alfa}1), A' I^* A, \\ B \rightarrow \text{beta}1.A'\text{beta}2 \in Kq \}.$$

Teorema 3.1 [PAR85]

Sea $p \in Q$ y $A \rightarrow \text{alfa}1.\text{alfa}2 \in Cp$, con $A \neq S$, entonces:

$$\text{LA}(p, A \rightarrow \text{alfa}1.\text{alfa}2) = \{ t \mid t \in \text{PATH}(A', A) @ \text{FIRST}(\text{beta}2) @ \\ \text{LA}(q, B \rightarrow \text{beta}1.A'\text{beta}2), \\ q \in \text{PRED}(p, \text{alfa}1), A' I^* A, B \rightarrow \text{beta}1.A'\text{beta}2 \in Kq \}.$$

Colocando la ecuación anterior en un lenguaje más ameno,

$$\text{LA}(p, A \rightarrow \text{alfa}1.\text{alfa}2) = \\ \cup \{ \cup \{ \text{PATH}(A', A) @ \text{FIRST}(\text{beta}2) @ \text{LA}(q, B \rightarrow \text{beta}1.A'\text{beta}2) \mid \\ A' I^* A, B \rightarrow \text{beta}1.A'\text{beta}2 \in Kq \\ \} \mid \\ q \in \text{PRED}(p, \text{alfa}1) \}.$$

Esta ecuación muestra explícitamente la idea principal del algoritmo para el cálculo de los LookAheads.

Algoritmo 3.1

Function $\text{LA}(p, A \rightarrow \text{alfa}1.\text{alfa}2)$ Return Set Of T Is

```

Begin
  LA := Ø ;
  For Each q ∈ PRED(p, alfa1) do Loop
    For Each B → beta1.A'beta2 ∈ Kq, A' I* A do Loop
      LA := LA U PATH(A', A) ;
      do_When (lambda ∈ PATH(A', A))
        LA := LA U FIRST(beta2) ;
      do_When (lambda ∈ FIRST(beta2))
        LA := LA U LA(q, B → beta1.A'beta2);
      End_d0 ;
    End_d0 ;
  End_Loop ;
End_Loop ;
End ;

```

Se observa, que el algoritmo anterior requiere una modificación para evitar llamadas recursivas al mismo estado con el mismo ítem. Para ello, construimos una relación que nos evite casos como el anterior, de tal manera que dicha relación nos permita, por así decirlo, construir un orden de evaluación para el cálculo de los LAs.

Tomando las ideas propuestas por DeRemer y Pennello [DeR79], redefinimos el conjunto FOLLOW de la siguiente manera:

Definición 3.10

Sea $q \in Q$ y $A \in N$,

$FOLLOW(q,A) = \{ a \mid a \in FIRST(\beta_2) @ LA(q,B \rightarrow \beta_1.A'\beta_2),$
 $B \rightarrow \beta_1.A'\beta_2 \in Kq, A I^* A' \}.$

Desprendiéndose el siguiente lema:

Lema 3.2 [PAR85]

Sean $p,q \in Q$ y $A \rightarrow \alpha_1.\alpha_2 \in Cp,$

$LA(p,A \rightarrow \alpha_1.\alpha_2) = \{ a \mid a \in FOLLOW(q,A),$
 $q \in PRED(p,\alpha_1) \}.$

Reescribiendo los FOLLOWS con base en los FOLLOWS,

Lema 3.3 [PAR85]

Sea $p,q \in Q,$

$FOLLOW(q,A) = \{ a \mid a \in FIRST(\beta_2) @ FOLLOW(r,B),$
 $B \rightarrow \beta_1.A\beta_2 \in Kq,$
 $r \in PRED(q,\beta_1),$
 $C \rightarrow \gamma_1.B\gamma_2 \in Cr \}.$

Si $\lambda \in FIRST(\beta_2),$ tenemos $FOLLOW(q,A)$ incluye al $FOLLOW(r,B).$

Definición 3.11

(q,A) INCLUDES (r,B) si y solamente si se cumple la condición anterior.

El cálculo de los LAs según este método envuelve evaluar la clausura reflexiva de la relación INCLUDES. El cálculo de los LAs para un ítem final $A \rightarrow w.$ de un estado $p,$ es la unión de los $FOLLOW(q,A)$ donde $q \in PRED(p,w).$ Obsérvese que

$FOLLOW(p,A) = \cup FOLLOW(q,B)$ tales que (p,A) INCLUDES* $(q,B).$

Algoritmo 3.2 [DeR82]

Construcción del grafo debido a la relación INCLUDES.

```

For Each  $q \in Q$  dO_Loop
  For Each  $B \rightarrow \beta_1.A\beta_2 \in Cq$  dO_Loop
    FOLLOW( $q,A$ ) := FIRST( $\beta_2$ ) ;
    INCLUDES( $q,A$ ) :=  $\emptyset$  ;
    dO_When( $\lambda \in FIRST(\beta_2)$ )
      For Each  $r \in PRED(q,\beta_1)$  dO_Loop
        For Each  $C \rightarrow \gamma_1.B\gamma_2 \in Cr$  dO_Loop
          INCLUDES( $q,A$ ) := INCLUDES( $q,A$ ) U { ( $r,B$ ) } ;
        End_Loop ;
      End_Loop ;
    End_dO ;
  End_Loop ;
End_Loop ;

```

Algoritmo 3.3 [DeR82]

Cálculo de los LookAheads usando FOLLOWS,

```

For Each p ∈ Q dO_Loop
  For Each A → w. ∈ Kp dO_Loop
    LA(q,A → w.) := ∅ ;
    For Each q ∈ PRED(p,w) dO_Loop
      LA(q,A → w.) := LA(q,A → w.) U FOLLOW(q,A) ;
    End_Loop ;
  End_Loop ;
End_Loop ;

```

El problema del método de DeRemer y Pennello es que la relación INCLUDES llega a ser demasiado grande, debido a la cantidad de ítems iniciales que pueden hacer aparecer arcos de la misma (en especial cuando hay producciones de la forma $A \rightarrow B$). Podemos obtener una reducción sustancial de dicha relación si sólo adicionamos arcos a través de los ítems presentes en el kernel exclusivamente. Y según el método desarrollado por J.C.H. Park, K. M. Choe y C. H. Chang [PAR85] tenemos,

Lema 3.4 [PAR85]

Sean $p, q \in Q$ y $A \rightarrow \underline{\text{alfa1.alfa2}} \in C_p$,

$$LA(p, A \rightarrow \underline{\text{alfa1.alfa2}}) = \{ a | a \in \text{PATH}(A', A) @ \text{FOLLOW}(q, A'), \\ q \in \text{PRED}(p, \underline{\text{alfa1}}), A' I^* A, \\ B \rightarrow \underline{\text{beta1.A'beta2}} \in K_q \}.$$

Lema 3.5 [PAR85]

Sean $p, q \in Q$,

$$\text{FOLLOW}(q, A') = \{ a | a \in \text{FIRST}(\underline{\text{beta2}}) @ \text{PATH}(P', B) @ \text{FOLLOW}(r, B'), \\ B \rightarrow \underline{\text{beta1.A'beta2}} \in K_q, \\ r \in \text{PRED}(q, \underline{\text{beta1}}), B' I^* B, \\ C \rightarrow \underline{\text{gamma1.B'gamma2}} \in K_r \}.$$

Si $\underline{\text{lambda}} \in \text{FIRST}(\underline{\text{beta2}})$ y $\underline{\text{lambda}} \in \text{PATH}(B', B)$, tenemos $\text{FOLLOW}(q, A')$ incluye al $\text{FOLLOW}(r, B')$.

El algoritmo para generar la relación INCLUDES quedaría:

Algoritmo 3.4 [PAR85]

```

For Each q ∈ Q dO_Loop
  For Each R -> beta1.Abeta2 ∈ Kq dO_Loop
    FOLLOW(q,A) := FIRST(beta2) ;
    INCLUDES(q,A) := ∅ ;
    dO_When(lambda ∈ FIRST(beta2))
      For Each r ∈ PRED(q,beta1) dO_Loop
        For Each C -> gamma1.Bgamma2 ∈ Kr, B' I* B
          dO_Loop
            FOLLOW(q,A') := FOLLOW(q,A') U PATH(B',B);
            dO_When(lambda ∈ PATH(B',B))
              INCLUDES(q,A) := INCLUDES(q,A) U { (r,B) }
            End_dO ;
          End_Loop ;
        End_Loop ;
      End_Loop ;
    End_dO ;
  End_Loop ;
End_Loop ;

```

Y el cálculo de los LookAheads,

Algoritmo 3.5 [PAR85]

```

For Each p ∈ Q dO_Loop
  For Each A -> w. ∈ Kp dO_Loop
    LA(q,A -> w.) := ∅ ;
    For Each q ∈ PRED(p,w), B -> beta1A beta2 ∈ Kq dO_Loop
      LA(q,A -> w.) := LA(q,A -> w.) U PATH(A',A) ;
      dO_When(lambda ∈ PATH(B',B))
        LA(q,A -> w.) := LA(q,A -> w.) U FOLLOW(q,A) ;
      End_dO ;
    End_Loop ;
  End_Loop ;
End_Loop ;

```

Todos los algoritmos anteriores calculan los FOLLOW de las transiciones no-terminales y construyen completamente la relación INCLUDES. Si se desea implementar dichos algoritmos en un microcomputador sería deseable construir un orden de evaluación sobre dicha relación para ir desechando FOLLOWS que no se necesitarán más en el cálculo de los LookAheads.

Extendiendo la relación INCLUDES entre no-terminales,

Definición 3.12

Si $A, B \in N$,

A INCLUDES B si y solamente si existe un ítem de la forma
 $R \rightarrow \underline{\text{alfa.Abeta}}$ tal que $\underline{\text{lambda}} \in \text{FIRST}(\underline{\text{beta}})$.

Esta relación no depende del autómata sino de la gramática. Si

esta idea va a ser usada en el algoritmo 3.5, $|\alpha| > 0$.

Definición 3.13

Si $A \in N$,

$$\text{ITEMS}(A) = U \{ U \{ B \rightarrow \alpha.A\beta \mid \lambda \in \text{FIRST}(\beta) \} \mid A \text{ INCLUDES } B \}.$$

Habiendo construido la relación INCLUDES para no-terminales sólo resta construir un orden de evaluación utilizando dicha relación. Recordando cómo sería un orden de evaluación usando un grafo de dependencia, tendríamos las siguientes restricciones sobre dicho orden, donde Orden(A) significa el orden de evaluación de los LAS para los ítems finales que tienen a A en su parte izquierda.

- Orden(B) < Orden(A) si y solamente si A INCLUDES* B y no B INCLUDES* A.
- Orden(B) = Orden(A) si y solamente si A INCLUDES* B y B INCLUDES* A. A y B tienen el mismo orden de evaluación si y sólo si A y B se encuentran en un ciclo de la relación INCLUDES.

De aquí encontramos que la relación INCLUDES entre las transiciones no-terminales de dos no-terminales A y B, no tiene necesidad de construirse, siempre y cuando ni A ni B tengan el mismo orden de evaluación. Si por ejemplo, existen (p,A) y (q,B) tales que (p,A) INCLUDES (q,B) no tenemos necesidad de construirla pues sabemos que Orden(B) < Orden(A) y por lo tanto los FOLLOWS de las transiciones no-terminales de B se han calculado. Simplemente, a FOLLOW(p,A) se le adiciona el FOLLOW(q,B).

Definición 3.14

Sea i la i -ésima posición en el orden de evaluación de los LAS,

$$G_i = \{ A \mid \text{Orden}(A) = i \}.$$

El algoritmo para el cálculo de los LAS usando un orden de evaluación quedaría,

Algoritmo 3.6

```

For Each i in 1..Max Orden de evaluación dO_Loop
  (* Inicialización de los FOLLOWS y construcción de
    INCLUDES para los no terminales en este orden *)
  For Each A ∈ Gi
    dO_Loop
      For Each (p,A) ! SUCC(p,A) está definida
        dO_Loop
          FOLLOW(p,A) := ∅ ;
          INCLUDES(p,A) := ∅ ;
        End_Loop
      For Each B → alfa.A beta ∈ ITEMS(A)
        dO_Loop
          For Each (q,B) ! SUCC(q,B) está definida
            dO_Loop
              p := SUCC(q,alfa) ;
              FOLLOW(p,A) := FOLLOW(p,A) U FIRST(beta) ;
              dO_When ( beta =>* lambda )
                For Each r ∈ PRED(p,alfa)
                  dO_Loop
                    For Each C → beta1.B'beta2 ∈ Kr
                      donde B' I* B
                        dO_Loop
                          FOLLOW(p,A) := FOLLOW(p,A) U PATH(B',B) ;
                          dO_When ( lambda ∈ PATH(B',B) )
                            dO_When ( B' ∈ Gi )
                              INCLUDES(p,A) := INCLUDES(p,A) U
                                { (r,B') } ;
                            Else_dO
                              FOLLOW(p,A) := FOLLOW(p,A) U
                                FOLLOW(r,B') ;
                            End_dO ;
                          End_dO ;
                        End_Loop ;
                      End_Loop ;
                    End_dO ;
                  End_Loop ;
                End_dO ;
              End_Loop ;
            End_Loop ;
          End_Loop ;
        End_Loop ;
      End_Loop ;
    End_Loop ;
  (* Evaluación de la relación INCLUDES para los no-
    terminales presentes en este orden *)
  (* Cálculo de los LAs para los items finales de los no-
    terminales envueltos en este orden *)
  For Each A ∈ Gi
    dO_Loop
      For Each A → w. ∈ Cp dO_Loop
        LA(q,A → w.) := ∅ ;
        For Each q ∈ PRED(p,w), B → beta1A beta2 - Kq
          dO_Loop
            LA(q,A → w.) := LA(q,A → w.) U PATH(A',A) ;
            dO_When ( lambda ∈ PATH(B',B) )
              LA(q,A → w.) := LA(q,A → w.) U FOLLOW(q,A) ;
            End_dO ;
          End_Loop ;
        End_Loop ;
      End_Loop ;
    End_Loop ;
  End_Loop ;

```

```

        End_Loop ;
    End_Loop ;
End_Loop ;

```

Este algoritmo es válido para las dos mejoras (DeRemer y Pennello y Park), pues, para la mejora de DeRemer y Pennello es suficiente eliminar desde la (*Línea 1*) hasta la (Línea 7*) y en el caso de Park calcular el PATH de acuerdo con la definición. La mejora básica de este algoritmo sobre el de Park está en la generación del orden de evaluación. Recordando las estadísticas de DeRemer y Pennello, para gramáticas prácticas (PASCAL, MODULA, etc) el número de no-terminales envueltos en ciclos en la relación INCLUDES es muy pequeño (PASCAL 7 no-terminales), y en nuestro algoritmo sólo se crea la relación INCLUDES para los no-terminales presentes en el orden de evaluación. Normalmente, el número de transiciones no-terminales en un ciclo de la relación INCLUDES (para transiciones no terminales) tiene una cota superior en el número de arcos que aparecen en la relación INCLUDES (entre los no-terminales) presentes en el mismo orden de evaluación. Si sumamos todas las transiciones no-terminales de los no-terminales presentes en el mismo orden de evaluación y multiplicamos por la cota superior anterior obtendremos una cota superior (un poco exagerada) del número de arcos presentes en la relación INCLUDES para este algoritmo.

4. Resolviendo Conflictos

En esta sección esbozaremos algunas ideas que pueden emplearse en la construcción de algoritmos que ilustren el por qué se han producido ciertos conflictos en un autómata. La información puede utilizarse entonces para intentar resolver los conflictos, pero el éxito en esta labor dependerá de la astucia del usuario para desembarazarse de ellos.

Definición 4.1

Sea $q \in Q$ y $t \in T$.

q tiene un conflicto con t si y sólo si, $SUCC(q,t)$ está definido y $|\text{Reduce}(q,t)| > 0$ y/o $|\text{Reduce}(q,t)| > 1$. Se llamará conflicto SHIFT-REDUCE si $SUCC(q,t)$ está definido y Conflicto REDUCE-REDUCE si $|\text{Reduce}(q,t)| > 1$. Un conflicto a la vez puede ser SHIFT-REDUCE o REDUCE-REDUCE.

En el caso de un conflicto SHIFT-REDUCE debemos mostrar dos tipos de información, una debido a la componente del conflicto $SUCC(q,t)$, que la denominaremos SHIFT-TRACE y otra debido a cada una de las producciones presentes en $\text{Reduce}(q,t)$, que la denominaremos REDUCE TRACE.

Definición 4.2

Sea $q \in Q$, $A \rightarrow \underline{\text{alfa.beta}} \in Cq$, decimos

$$CI(q, A \rightarrow \underline{\text{alfa.tbeta}}) = \{ \underline{\text{gamma}} \in V^* \mid \text{SUCC}(i^\circ, \underline{\text{gamma alfa}}) = q \}.$$

Sea q un estado que tenga un conflicto SHIFT-REDUCE con t . En qué consiste el SHIFT-TRACE ?. Si $\text{SUCC}(q, t)$ está definido es porque existe por lo menos un ítem de la forma $A \rightarrow \underline{\text{alfa.tbeta}} \in Cq$. Entonces, hay que justificar por qué $A \rightarrow \underline{\text{alfa.tbeta}} \in Cq$, por lo tanto el SHIFT-TRACE es igual a $CI(q, A \rightarrow \underline{\text{alfa.tbeta}})$ para cada ítem de la forma $A \rightarrow \underline{\text{alfa.tbeta}} \in Cq$.

Sin embargo, CI es un conjunto demasiado grande (puede ser no finito), y a veces no es necesario presentar toda la información para entender el por qué del conflicto. Por ejemplo, presentando un elemento de CI que tenga la mínima longitud es suficiente (el contexto izquierdo más corto).

A continuación presentaremos un algoritmo que calcula el contexto izquierdo más corto, usando exclusivamente la función SUCC de un autómata $\text{ALR}(0)$ cualquiera. Suponiendo que la numeración de los nodos del grafo asociado a dicho autómata se hubiera recorrido por extensión. La manera como se generó $\text{TP}(G)$ en el algoritmo 3.2 nos permite obtener dicha numeración.

Definición 4.2

Sea $\underline{\text{alfa}} \in V^*$ y $p \in Q$,

$$\text{MINPRED}(p, \underline{\text{alfa}}) = p \text{ si } \underline{\text{alfa}} \text{ es } \underline{\text{lambda}}$$

$$\text{MINPRED}(p, \underline{\text{alfa}}) = \text{MINPRED}(\text{MIN}(\{ q \mid q \in \text{PRED}(p, \text{Prim}(\underline{\text{alfa}})) \}), \text{Rest}(\underline{\text{alfa}})).$$

Definición 4.3

Sea $p, q \in Q$ y $X \in V$,

$$\text{ROTULO}(q) = X \text{ si y sólo si } \text{SUCC}(p, X) = q.$$

Teniendo las definiciones anteriores,

Algoritmo 4.1

Function $\text{CIC}(p, A \rightarrow \underline{\text{alfa.tbeta}})$ return V^* is

Begin

$\text{CIC} := \underline{\text{lambda}} \quad ;$

$q := \text{MINPRED}(q, \underline{\text{alfa}}) \quad ;$

WHILE $q \neq i^\circ$ dO_Loop

$\text{CIC} := \text{Append}(\text{CIC}, \text{ROTULO}(q)) \quad ;$

$q := \text{MINPRED}(q, \text{ROTULO}(q)) \quad ;$

End_Loop

return (CIC) ;

End CIC ;

Definición 4.4

Sean $B, C \in N$, $t \in T$, $t \in \text{PATH}(B, C)$ y, por la definición de $\text{PATH}(B, C)$, existen $B_0=B \dots, B_n=C$ y $B_0 \rightarrow B_1 \beta_1 \in P$, $B_1 \rightarrow B_2 \beta_2 \in P, \dots, B_{n-1} \rightarrow B_n \beta_n \in P$ y $t \in \text{FIRST}(\beta_n \dots \beta_1)$. Por último, existiría un $i \in [0..n]$ tal que $t \notin \text{FIRST}(\beta_n \dots \beta_{i+1})$ y $t \in \text{FIRST}(\beta_i)$.

A,

- $\beta_n \beta_{n-1} \dots \beta_{i+1} \Rightarrow^* \lambda$, los llamaremos los λ $\text{PATH}(B, C)$ de t para todo i que satisfaga la anterior condición ;
- $\beta_i \Rightarrow^* t \beta_i'$, los llamaremos las $\text{FIRST-PATH}(B, C)$ derivaciones de t para λ $\text{PATH}(B, C)$ de t ;
- $B_0 \Rightarrow^* B_1 \beta_n \dots \beta_1$ la $\text{PATH}(B, C)$ derivación,
- $B_0 \rightarrow B_1 \beta_1$ la producción contribuyente.

los llamaremos PATH-TRACES de t en $\text{PATH}(B, C)$.

Definición 4.5

Sea $q \in Q$, $A \in N$, $t \in T$, decimos que,

$$\beta \Rightarrow^* t \beta'$$

Es una FIRST derivación de t en $\text{FOLLOW}(q, A)$ si y sólo si existe $R \rightarrow \alpha \beta \in Kq$, tal que $t \in \text{FIRST}(\beta)$.

El siguiente algoritmo nos devolverá todos los $B \rightarrow \alpha \beta \in Kq$ que tengan una FIRST derivación de t en el $\text{FOLLOW}(q, A)$.

Algoritmo 4.2

Function $\text{FIRST_deriv}(q, A, t)$ return $N \times V^* \times V^* \text{ Is}$

```

Begin
  FIRST_deriv :=  $\emptyset$  ;
  For Each  $B \rightarrow \alpha \beta \in Kq$  |  $t \in \text{FIRST}(\beta)$ 
  do Loop
    FIRST_deriv :=  $\text{FIRST\_deriv} \cup \{ B \rightarrow \alpha \beta \}$  ;
  End Loop ;
  return (FIRST_deriv) ;
End ;

```

Definición 4.6

Sea $q \in Q$, $A \in N$, $t \in T$, y $B \rightarrow \alpha \beta \in Kq$, tal que $t \in \text{FIRST}(\beta)$, $\lambda \in \text{FIRST}(\beta)$ y $B' \stackrel{*}{\Rightarrow} B$ tal que $t \in \text{PATH}(B', B)$,

A los PATH-TRACES de t en $\text{PATH}(B', B)$ los denominaremos las PATHs derivaciones de t en $\text{FOLLOW}(q, A)$ y a $B \rightarrow \alpha \beta$ el ítem contribuyente.

Algoritmo 4.3

```

Function  PATH_deriv(q,A,t) return ( NxV#xV# )xN Is
Begin
  PATH_deriv := ∅ ;
  For Each B -> alfa.Abeta ∈ Kq ! beta =>#
  do Loop
    PATH_deriv := PATH_deriv U
    { (B -> alfa.Abeta, {B' | t ∈ PATH(B',B)}) }
  End Loop ;
  return (PATH_deriv) ;
End ;

```

La anterior función halla los ítems contribuyentes con sus respectivos no terminales que a través del PATH adicionan un símbolo terminal cualquiera al FOLLOW de la debida transición no terminal. Esta información es necesaria para producir todas las PATH-derivaciones de un terminal en el FOLLOW de una transición no terminal cualquiera.

Definición 4.7

Sea $q \in Q$, $A \in N$, $t \in T$, decimos que, $FOLLOW(q,A)$ incluye directamente a t , si y sólo si, hay una FIRST derivación de t en $FOLLOW(q,A)$ o hay una PATH derivación de t en $FOLLOW(q,A)$.

Definición 4.8

$(p,A \rightarrow w)$ LOOKBACK (q,B) si y solamente si existe B tal que $\lambda \in PATH(B,A)$, $q \in PRED(p,w)$.

Si $t \in LA(p,A \rightarrow w)$, entonces o existe $B \in N$ tal que $t \in PATH(B,A)$ o existen r y C tales que

$(p,A \rightarrow w)$ LOOKBACK (q,B) INCLUDES* (r,C)

y $FOLLOW(r,C)$ incluye directamente a t .

Ambos casos pueden suceder, pero nos interesa justificar la ocurrencia de al menos uno de ellos. Al primer caso, los vamos a llamar una inclusión directa de t en $LA(p,A \rightarrow w)$; y al segundo, una inclusión indirecta de t en $LA(p,A \rightarrow w)$. Las justificaciones de $t \in LA(p,A \rightarrow w)$ son todas las derivaciones de las inclusiones directas e indirectas y según las definiciones anteriores, se están abarcando todos los casos. Como siempre, para propósitos prácticos solo es deseable mostrar una de tales derivaciones: La primera que se encuentre según el método de búsqueda que sea empleado. No daremos algoritmos completos para la obtención directa de esta información, sino que nuestra intención es mostrar ideas generales para su desarrollo así como los algoritmos básicos para obtenerla.

5. Conclusiones

Las ideas desarrolladas en este trabajo, junto con las tomadas de los artículos descritos en la introducción y en la bibliografía, han sido hechas realidad en un sistema generador de analizadores sintácticos que ha sido usado como material de laboratorio en el curso de compiladores de la Universidad de los Andes (Bogotá, Colombia) con bastante éxito. La principal virtud del sistema es el tiempo de respuesta, debido en parte a la reducción de cálculos obtenida gracias al algoritmo 3.6.

El sistema, como se encuentra actualmente, recibe como entrada una gramática en forma de BNF o BNF-extendida y produce como salida un programa fuente manejador de la tabla de parsing, escrito en TURBOPASCAL, una descripción completa del autómata y las tablas, así como información para ayudar a resolver los conflictos.

Debido al excelente tiempo de respuesta y a la cantidad de información que proporciona, se ha convertido en una herramienta didáctica muy apreciada por los estudiantes, sobre todo por el hecho de ser utilizable sobre microcomputadores lo que le da una flexibilidad de acceso difícil de obtener sobre un computador grande (máquinas en las que tradicionalmente se ha implementado este tipo de software). Por otro lado, además de haber sido usado como material de laboratorio, el compiler-compiler ha sido probado para diferentes gramáticas "reales", entre las cuales tenemos: MODULA, TURBOPASCAL, un subconjunto de ADA (no hay problemas en la generación de ADA), una serie de gramáticas para generación de reportes y la gramática de PASCAL ANSI, con resultados satisfactorios, demostrando que el sistema es útil en un ambiente de producción real y es algo más que una facilidad didáctica.

En el estado actual, los fuentes están siendo pasados a C para obtener un mejor desempeño, aunque el tiempo de respuesta es excelente (sobre todo teniendo en cuenta que se utiliza en microcomputadores standard de 16 bits). Adicionalmente, se ha pensado agregarle otros módulos que comprenden: generador automático de tablas de símbolos, un generador automático de analizadores léxicos, un compactador de tablas esparcidas, un sencillo evaluador de atributos, un algoritmo general para generación de código para expresiones y un lenguaje para la especificación de lenguajes de máquina.

Referencias

- [AH077] Aho A, Ullman J.
"Principles of Compiler Design"
Addison-Wesley, Reading Mass., 1977.

-
- [DeR82] DeRemer F, Pennello T.J.
"Efficient Computation of LALR(1) LookAhead Sets"
ACM TOPLS Vol. 4 No. 4, octubre 1982
- [MON85] Montes A.
"Diseño e Implementación de un Generador de Compiladores"
Proyecto de Grado, Dpto. de Sistemas, UNIANDES, Bogotá,
1985.
- [PAR85] Park J, Choe K, Chang C.
"A New Analysis of LALR Formalisms"
ACM TOPLS Vol. 7 No. 1, enero 1985.