**Compiler Design (CS335), Spring 2024**
**Indian Institute of Technology Kanpur**
**Homework Assignment Number 2**

*Student Name:* Siddhant Jakhotiya
*Roll Number:* 211030
*Date:* February 14, 2024

**QUESTION**

**1**

## 1 Proving that the grammar is not LL(1)

Given the context-free grammar with the following productions:

$$Function \rightarrow Type\mathbf{id}(Arguments)$$
$$Type \rightarrow \mathbf{id}$$
$$Type \rightarrow Type\mathbf{*}$$
$$Arguments \rightarrow ArgList$$
$$Arguments \rightarrow \epsilon$$
$$ArgList \rightarrow Type\mathbf{id,}ArgList$$
$$ArgList \rightarrow Type\mathbf{id}$$

It can be easily observed that the grammar has left recursion ($Type \rightarrow Type\mathbf{*}$) and cannot be parsed by a top-down parser and is thus not LL(1).

We can also show that the above grammar is not LL(1) by constructing the predictive parsing table M and observing that there is some entry in the table which corresponds to more than one possible productions. We first start by computing the FIRST and FOLLOW sets.

**Computing FIRST and FOLLOW Sets**

**FIRST Sets**

- FIRST($\mathbf{T}$) = {$\mathbf{T}$} where T are all the terminal symbols i.e. id()*,$ and $\epsilon$.

- FIRST($Type$) = {$\mathbf{id}$} as we have the production $Type \rightarrow \mathbf{id}$

- FIRST($ArgList$) = FIRST($Type$) = {$\mathbf{id}$}

- FIRST($Arguments$) = FIRST($ArgList$) $\cup$ {$\epsilon$} = {$\mathbf{id}, \epsilon$}

- FIRST($Function$) = FIRST($Type$) = {$\mathbf{id}$}

**FOLLOW Sets**

- FOLLOW($Function$) = {$\$$}

- FOLLOW($Type$) = {$\mathbf{id}$, $\mathbf{*}$} since Type is followed by id and *.

- FOLLOW($Arguments$) = {)}

- FOLLOW($ArgList$) = {)} since we have the production $Arguments \rightarrow ArgList$.

**Constructing the Predictive Parsing Table M**

We fill the table entries M[A, a] for a nonterminal A and terminal a with the production $A \to \alpha$ if:

- $a \in \text{FIRST}(\alpha)$, or

- $\epsilon \in \text{FIRST}(\alpha)$ and $a \in \text{FOLLOW}(A)$.

| NON-TERMINAL | INPUT SYMBOL | | | | | |
|---|---|---|---|---|---|---|
| | id | * | ( | ) | , | $ |
| *Function* | *Function* $\to$ *Type***id** (*Arguments*) | | | | | |
| *Type* | *Type* $\to$ **id** <br> *Type* $\to$ *Type** | | | | | |
| *Arguments* | *Arguments* $\to$ *ArgList* | | | *Arguments* $\to \epsilon$ | | |
| *ArgList* | *ArgList* $\to$ *Type***id**, *ArgList* <br> *Arglist* $\to$ *Type***id** | | | | | |

Table 1: Predictive Parsing Table for the Given CFG

As it can be evidently seen, the table doesn't have a unique production associated for M[*Type*,**id**] and M[*ArgList*, **id**]. Thus it is not LL(1).

# 2 Constructing the LL(1) grammar

We first eliminate left-recursion in the grammar. The left-recursive derivation is *Type*$\to$*Type**. We introduce a new non-terminal $Type_1$ and the productions are as follows:

$$Function \to Type\mathbf{id}(Arguments)$$
$$Type \to \mathbf{id}\,Type_1$$
$$Type_1 \to {}^*Type_1 \mid \epsilon$$
$$Arguments \to ArgList \mid \epsilon$$
$$ArgList \to Type\mathbf{id} \mid Type\mathbf{id,}ArgList$$

Now, we perform left-factoring to remove the ambiguity about which production to use from *ArgList*. For this purpose, we introduce another non-terminal $ArgList_1$ and the productions are as follows:

$$Function \rightarrow Type\mathbf{id}(Arguments)$$
$$Type \rightarrow \mathbf{id}\,Type_1$$
$$Type_1 \rightarrow \mathbf{*}\,Type_1 \mid \epsilon$$
$$Arguments \rightarrow ArgList \mid \epsilon$$
$$ArgList \rightarrow Type\mathbf{id}\,ArgList_1$$
$$ArgList_1 \rightarrow \epsilon \mid \mathbf{,}\,ArgList$$

The resulting grammar should be LL(1).

# 3    FIRST and FOLLOW sets for the transformed grammar

**FIRST Sets**

- FIRST($\mathbf{T}$) = {$\mathbf{T}$} for all terminals i.e. id()*,$ and $\epsilon$.

- FIRST($Type$) = {$\mathbf{id}$} = FIRST($Function$) = FIRST($ArgList$)

- FIRST($Type_1$) = {$\mathbf{*}, \epsilon$}

- FIRST($Arguments$) = {$\mathbf{id}, \epsilon$}

- FIRST($ArgList_1$) = {$\mathbf{,}, \epsilon$}

**FOLLOW Sets**

- FOLLOW($Function$) = {$\$$}

- FOLLOW($Type$) = {$\mathbf{id}$}

- FOLLOW($Type_1$) = {$\mathbf{id}$}

- FOLLOW($Arguments$) = {$)$}

- FOLLOW($ArgList$) = {$)$}

- FOLLOW($ArgList_1$) = {$)$}

# 4 Predictive Parsing Table for transformed grammar

| NON-TERMINAL | INPUT SYMBOL | | | | | |
|---|---|---|---|---|---|---|
| | **id** | **\*** | **(** | **)** | **,** | **$** |
| $Function$ | $Function \rightarrow Type\,\mathbf{id}(Arguments)$ | | | | | |
| $Type$ | $Type \rightarrow \mathbf{id}\,Type_1$ | | | | | |
| $Type_1$ | $Type_1 \rightarrow \epsilon$ | $Type_1 \rightarrow \mathbf{*}\,Type_1$ | | | | |
| $Arguments$ | $Arguments \rightarrow ArgList$ | | | $Arguments \rightarrow \epsilon$ | | |
| $ArgList$ | $ArgList \rightarrow Type\,\mathbf{id}\,ArgList_1$ | | | | | |
| $ArgList_1$ | | | | $ArgList_1 \rightarrow \epsilon$ | $ArgList_1 \rightarrow \mathbf{,}\,ArgList$ | |

Table 2: Predictive Parsing Table for the Transformed CFG

As it can be seen, each cell in the table corresponds to a unique transition thus we have a LL(1) grammar.

*Student Name:* Siddhant Jakhotiya
*Roll Number:* 211030
*Date:* February 14, 2024

We have been given the CFG

$$S \rightarrow LM \mid Lp \mid qLr \mid sr \mid qsp$$
$$L \rightarrow aMb \mid s \mid t$$
$$M \rightarrow t$$

We first augment this grammar by adding the non-terminal S' s.t.

$$S' \rightarrow S$$
$$S \rightarrow LM \mid Lp \mid qLr \mid sr \mid qsp$$
$$L \rightarrow aMb \mid s \mid t$$
$$M \rightarrow t$$

# 1   Is the CFG SLR(1)?

**FIRST and FOLLOW Sets**

**FIRST Sets**

- FIRST($\mathbf{T}$) = {$\mathbf{T}$} for all terminals T $\in$ {a, b, p, q, r, s, t, \$}

- FIRST(M) = {t}

- FIRST(L) = {a, s, t}

- FIRST(S) = FIRST(L) $\cup$ {q, s} = {a, q, s, t} (which is also FIRST(S'), note that S' was simply introduced for augmentation)

**FOLLOW Sets**

- FOLLOW(S) = {\$} [which is also FOLLOW(S')]

- FOLLOW(L) = FIRST(M) $\cup$ {p, r} = {p, r, t} (since we have S $\rightarrow$LM | Lp | qLr)

- FOLLOW(M) = FOLLOW(S) $\cup$ {b} = {\$, b} (since we have S $\rightarrow$ LM and L $\rightarrow$ aMb)

### LR(0) canonical collection

We begin with the computation of CLOSURE($\{$[S' $\to$ .S]$\}$) using the definition of CLOSURE. It results in the item set:

$$S' \to .S$$
$$S \to .LM$$
$$S \to .Lp$$
$$S \to .qLr$$
$$S \to .sr$$
$$S \to .qsp$$
$$L \to .aMb$$
$$L \to .s$$
$$L \to .t$$

Similarly, we iteratively build the item sets using the GOTO function and iterating over all grammar symbols to obtain the LR(0) canonical collection:

**I₀**
$S' \to .S$
$S \to .LM$
$S \to .Lp$
$S \to .qLr$
$S \to .sr$
$S \to .qsp$
$L \to .aMb$
$L \to .s$
$L \to .t$

**I₁**
$S' \to S.$

**I₂**
$S \to L.M$
$S \to L.p$
$M \to .t$

**I₃**
$S \to q.Lr$
$S \to q.sp$
$L \to .aMb$
$L \to .s$
$L \to .t$

**I₄**
$S \to s.r$
$L \to s.$

**I₅**
$L \to a.Mb$
$M \to .t$

**I₆**
$L \to t.$

**I₇**
$S \to LM.$

**I₈**
$S \to Lp.$

**I₉**
$M \to t.$

**I₁₀**
$S \to qL.r$

**I₁₁**
$S \to qs.p$
$L \to s.$

**I₁₂**
$S \to sr.$

**I₁₃**
$L \to aM.b$

**I₁₄**
$S \to qLr.$

**I₁₅**
$S \to qsp.$

**I₁₆**
$L \to aMb.$

This is the LR(0) canonical collection. To construct the LR(0) automaton:

6

The states of the automaton are coded as the sets of items from the LR(0) canonical collection and the transitions are given by the GOTO function, with the start state being CLOSURE({[S' → .S]}). We introduce the transitions by considering the definition of the GOTO function:

GOTO(I,X) is defined as the closure of the set of all items [A → αX.β] such that [A → α.Xβ] is in I, where I is a set of items and X is a grammar symbol. Thus we have this automaton:



Here, the start state corresponds to $\mathbf{I}_0$.

## SLR Parsing Table

Using the LR(0) canonical items and the GOTO function and FOLLOW(A) for all non-terminals A, we determine the parsing actions for each state (where state i corresponds to $\mathrm{I}_i$).

In the table, s$i$ corresponds to "shift j" and r$i$ corresponds to using the $i^{th}$ reduction (the numbers associated with each production are defined below the table). **acc** means the accepting state and all empty entries correspond to error cases.

| STATE | ACTION | | | | | | | | GOTO | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | a | b | p | q | r | s | t | $ | S | L | M |
| 0 | s5 | | | s3 | | s4 | s6 | | 1 | 2 | |
| 1 | | | | | | | | **acc** | | | |
| 2 | | | s8 | | | | s9 | | | | 7 |
| 3 | s5 | | | | | s11 | s6 | | | 10 | |
| 4 | | | r8 | | s12 or r8 | | r8 | | | | |
| 5 | | | | | | | s9 | | | | 13 |
| 6 | | | r9 | | r9 | | r9 | | | | |
| 7 | | | | | | | | r2 | | | |
| 8 | | | | | | | | r3 | | | |
| 9 | | r10 | | | | | | r10 | | | |
| 10 | | | | | s14 | | | | | | |
| 11 | | | s15 or r8 | | r8 | | r8 | | | | |
| 12 | | | | | | | | r5 | | | |
| 13 | | s16 | | | | | | | | | |
| 14 | | | | | | | | r4 | | | |
| 15 | | | | | | | | r6 | | | |
| 16 | | | r7 | | r7 | | r7 | | | | |

Table 3: SLR Parsing Table for the CFG

$$S' \rightarrow S \ (1)$$
$$S \rightarrow LM \ (2) \mid Lp \ (3) \mid qLr \ (4) \mid sr \ (5) \mid qsp \ (6)$$
$$L \rightarrow aMb \ (7) \mid s \ (8) \mid t \ (9)$$
$$M \rightarrow t \ (10)$$

Here, the numbers in the bracket are used to number each production to denote it compactly in the parsing table.

As it can be clearly seen from the table, the grammar is not SLR(1) since it cannot take a unique parsing decision on the input r when in state 4 or on the input p when in state 11. This is because SLR doesn't have enough mechanisms to remember about what it has previously read to reach the current state.

## 2 Is the CFG LALR(1)?

The augmented grammar remains the same.

$$S' \rightarrow S$$
$$S \rightarrow LM \mid Lp \mid qLr \mid sr \mid qsp$$
$$L \rightarrow aMb \mid s \mid t$$
$$M \rightarrow t$$

### FIRST and FOLLOW Sets

#### FIRST Sets

- FIRST($\mathbf{T}$) = {$\mathbf{T}$} for all terminals T $\in$ {a, b, p, q, r, s, t, \$}

- FIRST(M) = {t}

- FIRST(L) = {a, s, t}

- FIRST(S) = FIRST(L) $\cup$ {q, s} = {a, q, s, t} (which is also FIRST(S'), note that S' was simply introduced for augmentation)

#### FOLLOW Sets

- FOLLOW(S) = {\$} [which is also FOLLOW(S')]

- FOLLOW(L) = FIRST(M) $\cup$ {p, r} = {p, r, t} (since we have S $\rightarrow$LM | Lp | qLr)

- FOLLOW(M) = FOLLOW(S) $\cup$ {b} = {\$, b} (since we have S $\rightarrow$ LM and L $\rightarrow$ aMb)

#### LALR collection

We first construct the LR(1) collection of items based on the definitions of CLOSURE and GOTO functions:

**I₇**
$S \rightarrow LM.$, $

**I₂**
$S \rightarrow L.M$, $
$S \rightarrow L.p$, $
$M \rightarrow .t$, $

**I₈**
$S \rightarrow Lp.$, $

**I₁**
$S' \rightarrow S.$, $

**I₉**
$M \rightarrow t.$, $

**I₁₇**
$L \rightarrow aMb.$, p/t

**I₀**
$S' \rightarrow .S$, $
$S \rightarrow .LM$, $
$S \rightarrow .Lp$, $
$S \rightarrow .qLr$, $
$S \rightarrow .sr$, $
$S \rightarrow .qsp$, $
$L \rightarrow .aMb$, p/t
$L \rightarrow .s$, p/t
$L \rightarrow .t$, p/t

**I₃**
$L \rightarrow a.Mb$, p/t
$M \rightarrow .t$, b

**I₁₀**
$L \rightarrow aM.b$, p/t

**I₁₁**
$M \rightarrow t.$, b

**I₁₈**
$S \rightarrow qLr.$, $

**I₄**
$S \rightarrow q.Lr$, $
$S \rightarrow q.sp$, $
$L \rightarrow .aMb$, r
$L \rightarrow .s$, r
$L \rightarrow .t$, r

**I₁₂**
$S \rightarrow qL.r$, $

**I₁₃**
$L \rightarrow a.Mb$, r
$M \rightarrow .t$, b

**I₁₉**
$L \rightarrow aM.b$, r

**I₆**
$L \rightarrow t.$, p/t

**I₁₄**
$S \rightarrow qs.p$, $
$L \rightarrow s.$, r

**I₂₁**
$L \rightarrow aMb.$, r

**I₅**
$S \rightarrow s.r$, $
$L \rightarrow s.$, p/t

**I₁₅**
$L \rightarrow t.$, r

**I₂₀**
$S \rightarrow qsp.$, $

**I₁₆**
$S \rightarrow sr.$, $

Given the algorithmic constructon of the LALR collection of states by merging the common cores, we have the following collection:

**I₁**
$S' \to S., \$$

**I₇**
$S \to LM., \$$

**I₀**
$S' \to .S, \$$
$S \to .LM, \$$
$S \to .Lp, \$$
$S \to .qLr, \$$
$S \to .sr, \$$
$S \to .qsp, \$$
$L \to .aMb, \text{p/t}$
$L \to .s, \text{p/t}$
$L \to .t, \text{p/t}$

**I₂**
$S \to L.M, \$$
$S \to L.p, \$$
$M \to .t, \$$

**I₈**
$S \to Lp., \$$

**I₁₇,₂₁**
$L \to aMb., \text{p/r/t}$

**I₉,₁₁**
$M \to t., \text{b/\$}$

**I₃,₁₃**
$L \to a.Mb, \text{p/r/t}$
$M \to .t, \text{b}$

**I₁₀,₁₉**
$L \to aM.b, \text{p/r/t}$

**I₆,₁₅**
$L \to t., \text{p/r/t}$

**I₄**
$S \to q.Lr, \$$
$S \to q.sp, \$$
$L \to .aMb, \text{r}$
$L \to .s, \text{r}$
$L \to .t, \text{r}$

**I₁₂**
$S \to qL.r, \$$

**I₁₈**
$S \to qLr., \$$

**I₅**
$S \to s.r, \$$
$L \to s., \text{p/t}$

**I₁₄**
$S \to qs.p, \$$
$L \to s., \text{r}$

**I₂₀**
$S \to qsp., \$$

**I₁₆**
$S \to sr., \$$

## LR(1) automaton

We add the GOTO transitions to the above collection of atoms to get the automaton. The start state corresponds to CLOSURE({[S' → .S]}), which has been marked as **I₀** in the diagram.

**accept**

$

**I₁**
$S' \to S., \$$

**I₀**
$S' \to .S, \$$
$S \to .LM, \$$
$S \to .Lp, \$$
$S \to .qLr, \$$
$S \to .sr, \$$
$S \to .qsp, \$$
$L \to .aMb, \text{p/t}$
$L \to .s, \text{p/t}$
$L \to .t, \text{p/t}$

**I₂**
$S \to L.M, \$$
$S \to L.p, \$$
$M \to .t, \$$

**I₇**
$S \to LM., \$$

**I₈**
$S \to Lp., \$$

**I₉**
$M \to t., \$$

**I₁₇**
$L \to aMb., \text{p/t}$

**I₃**
$L \to a.Mb, \text{p/t}$
$M \to .t, \text{b}$

**I₁₀**
$L \to aM.b, \text{p/t}$

**I₁₁**
$M \to t., \text{b}$

**I₁₈**
$S \to qLr., \$$

**I₁₂**
$S \to qL.r, \$$

**I₄**
$S \to q.Lr, \$$
$S \to q.sp, \$$
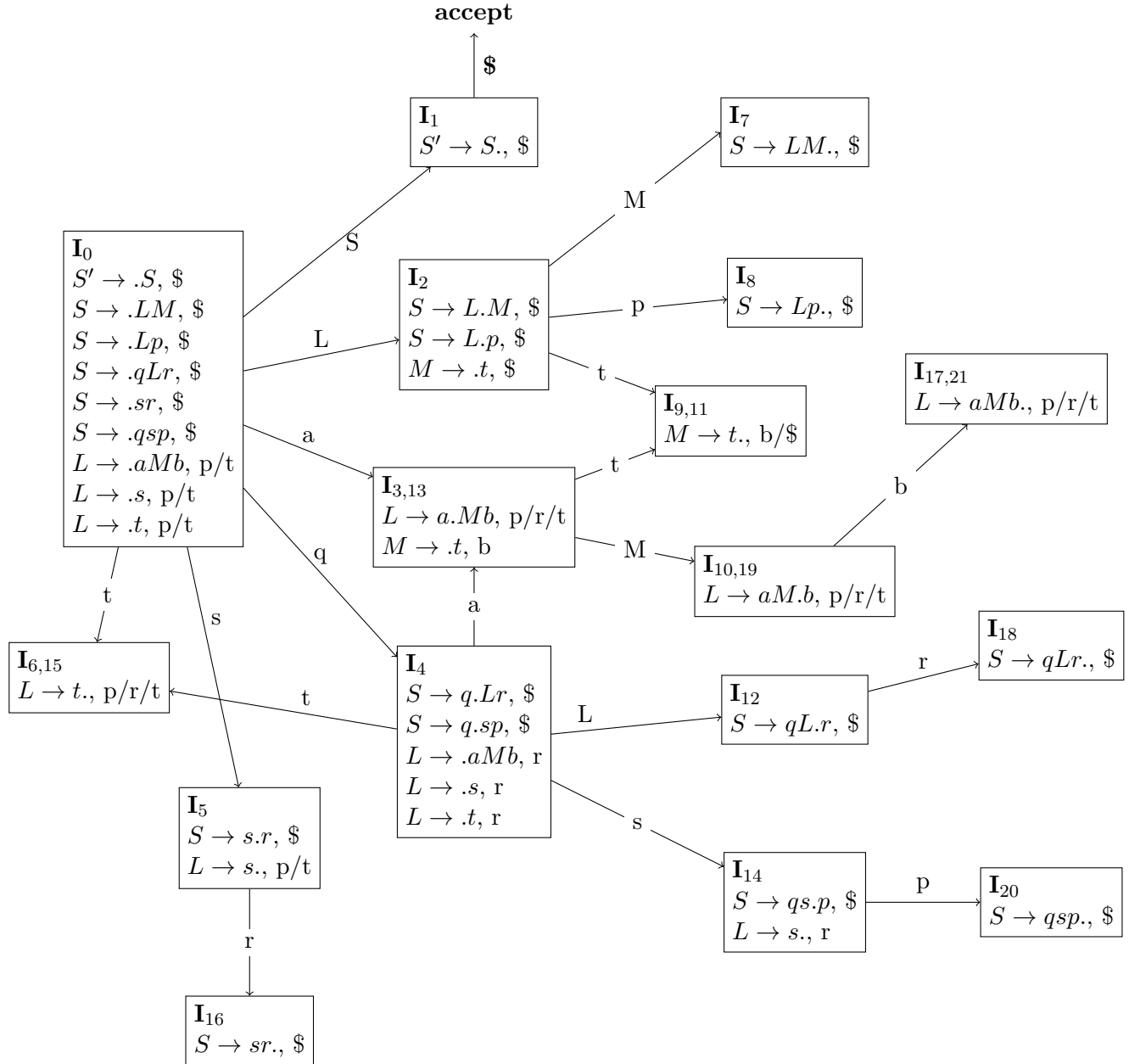$L \to .aMb, \text{r}$
$L \to .s, \text{r}$
$L \to .t, \text{r}$

**I₁₃**
$L \to a.Mb, \text{r}$
$M \to .t, \text{b}$

**I₁₉**
$L \to aM.b, \text{r}$

**I₆**
$L \to t., \text{p/t}$

**I₁₄**
$S \to qs.p, \$$
$L \to s., \text{r}$

**I₂₁**
$L \to aMb., \text{r}$

**I₅**
$S \to s.r, \$$
$L \to s., \text{p/t}$

**I₁₅**
$L \to t., \text{r}$

**I₂₀**
$S \to qsp., \$$

**I₁₆**
$S \to sr., \$$

Transitions (labels): I₁ → accept on $; I₀ → I₁ on S; I₀ → I₂ on L; I₀ → I₃ on a; I₀ → I₄ on q; I₀ → I₅ on s; I₀ → I₆ on t; I₂ → I₇ on M; I₂ → I₈ on p; I₂ → I₉ on t; I₃ → I₁₀ on M; I₃ → I₁₁ on t; I₁₀ → I₁₇ on b; I₄ → I₁₂ on L; I₄ → I₁₃ on a; I₄ → I₁₄ on s; I₄ → I₁₅ on t; I₁₂ → I₁₈ on r; I₁₃ → I₁₉ on M; I₁₃ → I₁₁ on t; I₁₉ → I₂₁ on b; I₁₄ → I₂₀ on p; I₅ → I₁₆ on r.

## LALR(1) Parsing Table

Now that we have the collection of items, we need the ACTION and GOTO values for the states. For the sake of brevity, the automaton with the reduced sets is also drawn:

**accept**

$\mathbf{I}_0$
$S' \rightarrow .S, \$$
$S \rightarrow .LM, \$$
$S \rightarrow .Lp, \$$
$S \rightarrow .qLr, \$$
$S \rightarrow .sr, \$$
$S \rightarrow .qsp, \$$
$L \rightarrow .aMb, \text{p/t}$
$L \rightarrow .s, \text{p/t}$
$L \rightarrow .t, \text{p/t}$

$\mathbf{I}_1$
$S' \rightarrow S., \$$

$\mathbf{I}_7$
$S \rightarrow LM., \$$

$\mathbf{I}_2$
$S \rightarrow L.M, \$$
$S \rightarrow L.p, \$$
$M \rightarrow .t, \$$

$\mathbf{I}_8$
$S \rightarrow Lp., \$$

$\mathbf{I}_{9,11}$
$M \rightarrow t., \text{b/}\$$

$\mathbf{I}_{17,21}$
$L \rightarrow aMb., \text{p/r/t}$

$\mathbf{I}_{3,13}$
$L \rightarrow a.Mb, \text{p/r/t}$
$M \rightarrow .t, \text{b}$

$\mathbf{I}_{10,19}$
$L \rightarrow aM.b, \text{p/r/t}$

$\mathbf{I}_{6,15}$
$L \rightarrow t., \text{p/r/t}$

$\mathbf{I}_4$
$S \rightarrow q.Lr, \$$
$S \rightarrow q.sp, \$$
$L \rightarrow .aMb, \text{r}$
$L \rightarrow .s, \text{r}$
$L \rightarrow .t, \text{r}$

$\mathbf{I}_{12}$
$S \rightarrow qL.r, \$$

$\mathbf{I}_{18}$
$S \rightarrow qLr., \$$

$\mathbf{I}_5$
$S \rightarrow s.r, \$$
$L \rightarrow s., \text{p/t}$

$\mathbf{I}_{14}$
$S \rightarrow qs.p, \$$
$L \rightarrow s., \text{r}$

$\mathbf{I}_{20}$
$S \rightarrow qsp., \$$

$\mathbf{I}_{16}$
$S \rightarrow sr., \$$

The transitions have been numbered as follows:

$$S' \rightarrow S \ (1)$$
$$S \rightarrow LM \ (2) \mid Lp \ (3) \mid qLr \ (4) \mid sr \ (5) \mid qsp \ (6)$$
$$L \rightarrow aMb \ (7) \mid s \ (8) \mid t \ (9)$$
$$M \rightarrow t \ (10)$$

The ACTION and GOTO values are computed as per the algorithm and displayed below. 's$i$' refers to "shift i", 'r$i$' refers to using the i$^{th}$ transition as a reduction. **acc** refers to the accept action. All the blank entries become error cases.

| | ACTION | | | | | | | | GOTO | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **STATE** | a | b | p | q | r | s | t | $ | S | L | M |
| 0 | s3,13 | | | s4 | | s5 | s6,15 | | 1 | 2 | |
| 1 | | | | | | | | **acc** | | | |
| 2 | | | s8 | | | | s9,11 | | | | 7 |
| 3,13 | | | | | | | s9,11 | | | | 10,19 |
| 4 | s3,13 | | | | | s14 | s6,15 | | | 12 | |
| 5 | | | r8 | | s16 | | r8 | | | | |
| 6,15 | | | r9 | | r9 | | r9 | | | | |
| 7 | | | | | | | | r2 | | | |
| 8 | | | | | | | | r3 | | | |
| 9,11 | | r10 | | | | | | r10 | | | |
| 10,19 | | s17,21 | | | | | | | | | |
| 12 | | | | | s18 | | | | | | |
| 14 | | | s20 | | r8 | | | | | | |
| 16 | | | | | | | | r5 | | | |
| 17,21 | | | r7 | | r7 | | r7 | | | | |
| 18 | | | | | | | | r4 | | | |
| 20 | | | | | | | | r6 | | | |

Table 4: LALR Parsing Table for the CFG

As it can be seen, for each state on each input, we can take a unique parsing decision. Thus the grammar is LALR(1).

*Student Name:* Siddhant Jakhotiya
*Roll Number:* 211030
*Date:* February 14, 2024

## 1   Problem Description

The goal was the generate a parser and display the desired output for the given markup schema.

## 2   Instructions to Run

The subdirectory 'problem3' contains 3 files: prob3.l, prob3.y and prob3.sh. To run, ensure that the present working directory is problem3. Now,

- Open prob3.sh. It contains a variable called *file_name*. Provide the relative file path for the testcase as the value for this variable. Ensure that the entire filename (with the file extension) is written. Eg. if the filename is *test.html*, this entire name should be written.

- Now, run prob3.sh. It can be executed in 2 ways:

    - ./prob3.sh: It displays the output on the terminal (stdout)
    - ./prob3.sh -f: It redirects the output to file_name.output. Note that the .output file will be in the same subdirectory as the original testcase.

- The script automatically deletes all the additional files that it has created.

## 3   Corner/Error Cases and Format of Output

- The parser continues parsing the input till it encounters the first error. It reports the error and terminates the execution. All misspelt/random tags as well as words outside the placeholder for a given text are ignored.

- The output format is as instructed in the PDF.

- The parser detects the following errors:

    - Marks are out of range for the given question type.
    - Marks are not enclosed within ".
    - There are too few/many choices.
    - If there are 0 correct choices for a question.
    - If there are >1 correct choices for a singleselect type question.
    - Number of correct choices > number of choices for a multiselect type question.
    - Missing opening and closing tags for quiz/singleselect/multiselect/choice/correct. Note that stray closing tags are also highlighted by reporting that the opening tag is missing!

- The following forms of input required modifications to the grammar and in my opinion are worthy of being included in the testcases:

  - **An empty program**: An empty file should also be syntactically valid and thus accepted by the parser without reporting any error.

  - **An integer present outside the placeholder**: The following input required non-trivial modifications to the grammar to parse correctly. Note that it is syntactically correct and should result in no errors:

    ```
    <quiz>
        <singleselect marks="2"> This is a question
            <choice>20</choice>10
            <choice>15</choice>
            <choice>10</choice>
            <correcct>15</correct>
        </singleselect>
    </quiz>
    ```

    i.e. the additional 10 (a **number**) present outside the tokens shouldn't result in an error. Along similar lines, the input:

    ```
    <quiz>
        <singleselect marks=" 1"> This is a question
            <choice>20</choice>10
            <choice>15</choice>
            <choice>10</choice>12
            15</correct>
        </singleselect>
    </quiz>
    ```

    should result in a 'missing opening correct tag' error. This also required some modification to the grammar for correct error reporting.

  - **Attribute inside any of the valid tag**: As instructed in the PDF, we have been told that any attribute except marks inside any placeholder has to be ignored. Thus an input of the following type:

    ```
    <quiz ignore this text>
        <singleselect marks=" 2"> This is a question
            <choice>A</choice>
            <choice asdasd>B</choice>
            <choice>C</choICE>
        </singleselect>
    </quiz>
    ```

    Should not result in any error.