



*Parsing V*  
*Introduction to LR(1) Parsers*



## LR(1) Parsers

---

- LR(1) parsers are table-driven, shift-reduce parsers that use a limited right context (1 word) for handle recognition
- LR(1) parsers recognize languages that have an LR(1) grammar

### *Informal definition:*

A grammar is LR(1) if, given a rightmost derivation

$$S \Rightarrow \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n \Rightarrow \textit{sentence}$$

We can

1. *isolate the handle of each right-sentential form  $\gamma_i$ , and*
2. *determine the production by which to reduce,*

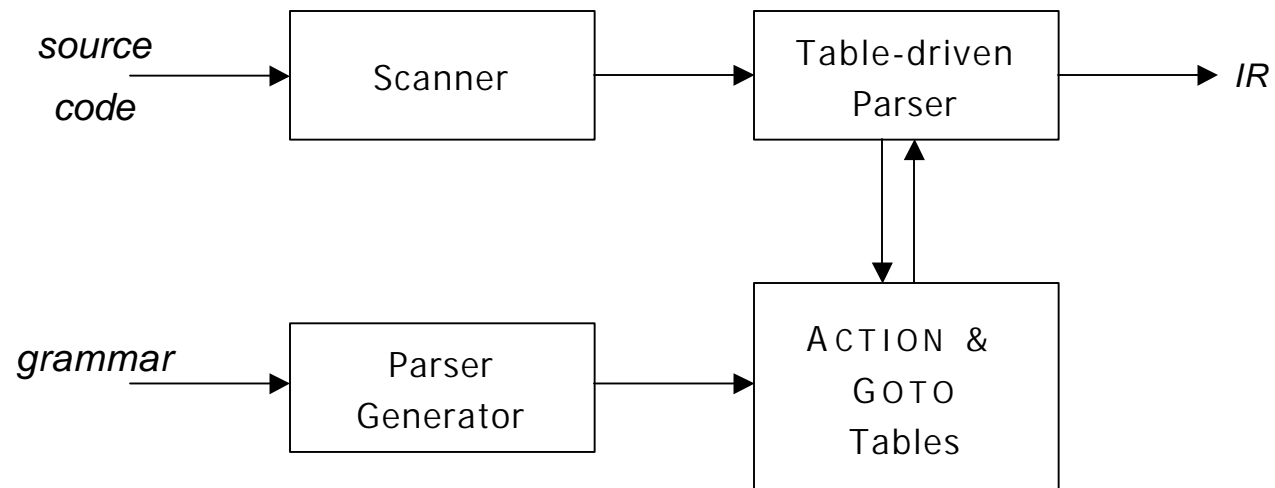
by scanning  $\gamma_i$  from left-to-right, going at most 1 symbol beyond the right end of the handle of  $\gamma_i$



## LR(1) Parsers

---

A table-driven LR(1) parser looks like



Tables can be built by hand

*(homework # 2)*

It is a perfect task to automate

## LR(1) Parsers

(the skeleton parser)



```
push INVALID
push  $s_0$ 
token  $\leftarrow$  next_token()
repeat forever
  s  $\leftarrow$  top of stack
  if ACTION[s,token] = "reduce  $N \rightarrow \beta$ "
    then
      pop  $2 * |\beta|$  symbols
      s  $\leftarrow$  top of stack
      push N
      push GOTO[s,N]
  else if ACTION[s,token] = "shift  $s_i$ "
    then
      push token ; push  $s_i$ 
      token  $\leftarrow$  next_token()
  else if ACTION[s,token] = "accept"
    and token = EOF
    then break;
  else report a syntax error
report success
```

### The skeleton parser

- uses ACTION & GOTO
- does |*words*| shifts
- does |derivation| reductions
- does 1 accept
- detects errors by failure of 3 other cases



## LR(1) Parsers (parse tables)

---

To make a parser for  $L(G)$ , need a set of tables

The grammar

|   |                   |   |                       |
|---|-------------------|---|-----------------------|
| 1 | <i>Goal</i>       | → | SheepNoise            |
| 2 | <i>SheepNoise</i> | → | SheepNoise <u>baa</u> |
| 3 |                   |   | <u>baa</u>            |

The tables

| ACTION |          |            |
|--------|----------|------------|
| State  | EOF      | <u>baa</u> |
| 0      | —        | shift 2    |
| 1      | accept   | shift 3    |
| 2      | reduce 3 | reduce 3   |
| 3      | reduce 2 | reduce 2   |

| GOTO  |                   |
|-------|-------------------|
| State | <i>SheepNoise</i> |
| 0     | 1                 |
| 1     | 0                 |
| 2     | 0                 |
| 3     | 0                 |



## Example Parses

The string "baa"

| Inp u t                 | S t a c k                                     | A c t i o n   |
|-------------------------|---|---------------|
| <u>b a a</u> <u>EOF</u> | \$ s <sub>0</sub>                             | s h i f t 2   |
| <u>EOF</u>              | \$ s <sub>0</sub> <u>b a a</u> s <sub>2</sub> | r e d u c e 3 |
| <u>EOF</u>              | \$ s <sub>0</sub> <i>S N</i> s <sub>1</sub>   | a c c e p t   |

We cannot have a syntax error with *SN*, because it only has 1 terminal symbol!

"baa woof" is a lexical problem, not a syntax error!

The string "baa"

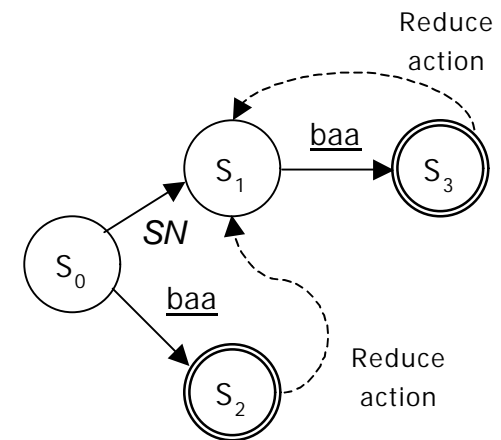
| Inp u t                              | S t a c k   | A c t i o n   |
|--------------------------------------|---|---------------|
| <u>b a a</u> <u>b a a</u> <u>EOF</u> | \$ s <sub>0</sub>   | s h i f t 2   |
| <u>b a a</u> <u>EOF</u>              | \$ s <sub>0</sub> <u>b a a</u> s <sub>2</sub>                           | r e d u c e 3 |
| <u>b a a</u> <u>EOF</u>              | \$ s <sub>0</sub> <i>S N</i> s <sub>1</sub>                             | s h i f t 3   |
| <u>EOF</u>                           | \$ s <sub>0</sub> <i>S N</i> s <sub>1</sub> <u>b a a</u> s <sub>3</sub> | r e d u c e 2 |
| <u>EOF</u>                           | \$ s <sub>0</sub> <i>S N</i> s <sub>1</sub>                             | a c c e p t   |



## LR(1) Parsers

How does this LR(1) stuff work?

- Unambiguous grammar  $\Rightarrow$  unique rightmost derivation
- Keep upper fringe on a stack
  - > All active handles include TOS
  - > Shift inputs until TOS is right end of a handle
- Language of handles is regular (finite)
  - > Build a handle-recognizing DFA
  - > ACTION & GOTO tables encode the DFA
- To match subterms, recurse & leave DFA's state on stack
- Final state in DFA  $\Rightarrow$  a *reduce* action
  - > New state is GOTO[*lhs*,state at TOS]
  - > For *SN*, this takes the DFA to  $S_1$



Control DFA for SN



## Building LR(1) Parsers

---

How do we generate the ACTION and GOTO tables?

- Use the grammar to build a model of the DFA
- Use the model to build ACTION & GOTO tables
- If construction succeeds, the grammar is LR(1)

The Big Picture

- Model the state of the parser
- Use two functions *goto(s,N)* and *closure(s)*
  - > *goto()* is analogous to *move()* in the subset construction
  - > *closure()* adds information to round out a state
- Build up the states and transition functions of the DFA
- Use this information to fill in the ACTION and GOTO tables



## $LR(k)$ items

---

An  $LR(k)$  item is a pair  $[A, B]$ , where

$A$  is a production  $\alpha \rightarrow \beta\gamma\delta$  with a  $\cdot$  at some position in the *rhs*

$B$  is a lookahead string of length  $k$  (words or EOF)

The  $\cdot$  in an item indicates the position of the top of the stack

$[\alpha \rightarrow \cdot \beta\gamma\delta, \underline{a}]$  means that the input seen so far is consistent with the use of  $\alpha \rightarrow \beta\gamma\delta$  immediately after the symbol on top of the stack

$[\alpha \rightarrow \beta\gamma \cdot \delta, \underline{a}]$  means that the input seen so far is consistent with the use of  $\alpha \rightarrow \beta\gamma\delta$  at this point in the parse, and that the parser has already recognized  $\beta\gamma$ .

$[\alpha \rightarrow \beta\gamma\delta \cdot, \underline{a}]$  means that the parser has seen  $\beta\gamma\delta$ , and that a lookahead symbol of  $\underline{a}$  is consistent with reducing to  $\alpha$ .

The table construction algorithm uses items to represent valid configurations of an  $LR(1)$  parser



## LR(1) Items

---

The production  $\alpha \rightarrow \cdot \beta \gamma \delta$ , with lookahead  $\underline{a}$ , generates 4 items

$$[\alpha \rightarrow \cdot \beta \gamma \delta, \underline{a}], [\alpha \rightarrow \beta \cdot \gamma \delta, \underline{a}], [\alpha \rightarrow \beta \gamma \cdot \delta, \underline{a}], \text{ \& } [\alpha \rightarrow \beta \gamma \delta \cdot, \underline{a}]$$

The set of LR(1) items for a grammar is finite

What's the point of all these lookahead symbols?

- Carry them along to choose correct reduction *(if a choice occurs)*
  - Lookaheads are bookkeeping, unless item has  $\cdot$  at right end
    - > Has no direct use in  $[\alpha \rightarrow \beta \gamma \cdot \delta, \underline{a}]$
    - > In  $[\alpha \rightarrow \beta \gamma \delta \cdot, \underline{a}]$ , a lookahead of  $\underline{a}$  implies a reduction by  $\alpha \rightarrow \beta \gamma \delta$
    - > For  $\{ [\alpha \rightarrow \gamma \cdot, \underline{a}], [\beta \rightarrow \gamma \cdot \delta, \underline{b}] \}$ ,  $\underline{a} \Rightarrow$  **reduce** to  $\alpha$ ;  $\text{FIRST}(\delta) \Rightarrow$  **shift**
- $\Rightarrow$  Limited right context is enough to pick the actions



## LR(1) Table Construction

---

High-level overview

- 1 Build the canonical collection of sets of LR(1) Items,  $I$ 
  - a Begin in an appropriate state,  $i_0$ 
    - $[S' \rightarrow \bullet S, \underline{EOF}]$ , along with any equivalent items
    - Derive equivalent items as  $closure(i_0)$
  - b Repeatedly compute, for each  $i_k$ , and each  $\mathbf{a}$ ,  $goto(i_k, \mathbf{a})$ 
    - If the set is not already in the collection, add it
    - Record all the transitions created by  $goto( )$

This eventually reaches a fixed point
- 2 Fill in the table from the collection of sets of LR(1) items

*The canonical collection completely encodes the transition diagram for the handle-finding DFA*



## Back to Finding Handles

---

Revisiting an issue from last class

Parser in a state where the stack (the fringe) was

$Expr \_ Term$

With lookahead of  $\_*$

How did it choose to expand  $Term$  rather than reduce to  $Expr$ ?

- *Lookahead* symbol is the key
- With lookahead of  $\_+$  or  $\_-$ , parser should reduce to  $Expr$
- With lookahead of  $\_*$  or  $\_ /$ , parser should shift
- Parser uses lookahead to decide
- All this context from the grammar is encoded in the handle recognizing mechanism

Remember this slide from last lecture?



Back to  $x - 2 * y$

| Stack  | Input                               | Handle | Action   |
|--|-------------------------------------|--------|----------|
| \$   | <u>id</u> - <u>nu m</u> * <u>id</u> | non e  | shif t   |
| \$ <u>id</u>                                     | - <u>nu m</u> * <u>id</u>           | 9,1    | red . 9  |
| \$ <i>Fac to r</i>                               | - <u>nu m</u> * <u>id</u>           | 7,1    | red . 7  |
| \$ <i>Te rm</i>                                  | - <u>nu m</u> * <u>id</u>           | 4,1    | red . 4  |
| \$ <i>Ex pr</i>                                  | - <u>nu m</u> * <u>id</u>           | non e  | shif t   |
| \$ <i>Ex pr</i> -                                | <u>nu m</u> * <u>id</u>             | non e  | shif t   |
| \$ <i>Ex pr</i> - <u>nu m</u>                    | * <u>id</u>                         | 8,3    | red . 8  |
| \$ <i>Ex pr</i> - <i>Fac to r</i>                | * <u>id</u>                         | 7,3    | red . 7  |
| \$ <i>Ex pr</i> - <i>Te rm</i>                   | * <u>id</u>                         | non e  | shif t   |
| \$ <i>Ex pr</i> - <i>Te rm</i> *                 | <u>id</u>                           | non e  | shif t   |
| \$ <i>Ex pr</i> - <i>Te rm</i> * <u>id</u>       |                                     | 9,5    | red . 9  |
| \$ <i>Ex pr</i> - <i>Te rm</i> * <i>Fac to r</i> |                                     | 5,5    | red . 5  |
| \$ <i>Ex pr</i> - <i>Te rm</i>                   |                                     | 3,3    | red . 3  |
| \$ <i>Ex pr</i>                                  |                                     | 1,1    | red . 1  |
| \$ <i>Goal</i>                                   |                                     | non e  | ac c ept |

shift here

reduce here

1. Shift until TOS is the right end of a handle
2. Find the left end of the handle & reduce



## Next Class

---

- Algorithms for FIRST, *goto*, & *closure*
- Work an example — a simplified expression grammar

To prepare

- Look at the book or web pages
- Work through the *SheepNoise* example



## Computing FIRST Sets

Define FIRST as

- If  $\alpha \Rightarrow^* \underline{a}\beta$ ,  $\underline{a} \in T$ ,  $\beta \in (T \cup NT)^*$ , then  $\underline{a} \in \text{FIRST}(\alpha)$
- If  $\alpha \Rightarrow^* \epsilon$ , then  $\epsilon \in \text{FIRST}(\alpha)$

To compute FIRST

- Use a fixed-point method
- $\text{FIRST}(\alpha) \in 2^{(T \cup \epsilon)}$
- Loop is monotonic

$\Rightarrow$  Algorithm halts

For *SheepNoise*:

$\text{FIRST}(\textit{Goal}) = \{ \underline{\textit{baa}} \}$

$\text{FIRST}(\textit{SM}) = \{ \underline{\textit{baa}} \}$

$\text{FIRST}(\underline{\textit{baa}}) = \{ \underline{\textit{baa}} \}$

```

For each  $\alpha \in T$ 
     $\text{FIRST}(\alpha) \leftarrow \alpha$ 
For each  $\alpha \in NT$ ,
     $\text{FIRST}(\alpha) \leftarrow \emptyset$ 
While (FIRST sets are still changing)
    for each  $p \in P$ , of the form  $\mathbf{a} @ \mathbf{b}$ ,
        if  $\mathbf{b}$  is  $\epsilon$  then
             $\text{FIRST}(\mathbf{a}) \leftarrow \text{FIRST}(\mathbf{a}) \cup \{ \epsilon \}$ 
        else if  $\mathbf{b}$  is  $\mathbf{b}_1 \mathbf{b}_2 \dots \mathbf{b}_k$  then
             $\text{FIRST}(\mathbf{a}) \leftarrow \text{FIRST}(\mathbf{a}) \cup \text{FIRST}(\mathbf{b}_1)$ 
             $i \leftarrow 1$ 
            while ( $\epsilon \in \text{FIRST}(\mathbf{b}_i)$  &  $i < k$ )
                 $\text{FIRST}(\mathbf{a}) \leftarrow \text{FIRST}(\mathbf{a}) \cup \text{FIRST}(\mathbf{b}_{i+1})$ 
                 $i \leftarrow i + 1$ 

```