

CS3300 - Compiler Design

Parsing

V. Krishna Nandivada

IIT Madras

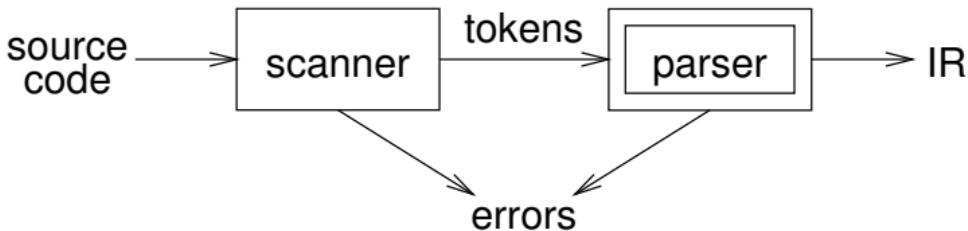
Acknowledgement

These slides borrow liberal portions of text verbatim from Antony L. Hosking @ Purdue, Jens Palsberg @ UCLA and the Dragon book.

Copyright ©2025 by Antony L. Hosking. *Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or fee. Request permission to publish from hosking@cs.purdue.edu.*



The role of the parser



A parser

- performs context-free syntax analysis
- guides context-sensitive analysis
- constructs an intermediate representation
- produces meaningful error messages
- attempts error correction

For the next several classes, we will look at parser construction



Syntax analysis by using a CFG

Context-free syntax is specified with a *context-free grammar*.

Formally, a CFG G is a 4-tuple (V_t, V_n, S, P) , where:

V_t is the set of *terminal symbols* in the grammar.

For our purposes, V_t is the set of tokens returned by the scanner.

V_n , the *nonterminals*, is a set of syntactic variables that denote sets of (sub)strings occurring in the language.

These are used to impose a structure on the grammar.

S is a distinguished nonterminal ($S \in V_n$) denoting the entire set of strings in $L(G)$.

This is sometimes called a *goal symbol*.

P is a finite set of *productions* specifying how terminals and non-terminals can be combined to form strings in the language.

Each production must have a single non-terminal on its left hand side.

The set $V = V_t \cup V_n$ is called the *vocabulary* of G .



Notation and terminology

- $a, b, c, \dots \in V_t$
- $A, B, C, \dots \in V_n$
- $U, V, W, \dots \in V$
- $\alpha, \beta, \gamma, \dots \in V^*$
- $u, v, w, \dots \in V_t^*$

If $A \rightarrow \gamma$ then $\alpha A \beta \Rightarrow \alpha \gamma \beta$ is a *single-step derivation* using $A \rightarrow \gamma$

Similarly, \rightarrow^* and \Rightarrow^+ denote derivations of ≥ 0 and ≥ 1 steps

If $S \rightarrow^* \beta$ then β is said to be a *sentential form* of G

$L(G) = \{w \in V_t^* \mid S \Rightarrow^+ w\}$, $w \in L(G)$ is called a *sentence* of G

Note, $L(G) = \{\beta \in V^* \mid S \rightarrow^* \beta\} \cap V_t^*$



Syntax analysis

Grammars are often written in Backus-Naur form (BNF).

Example:

| | | | |
|---|-------------------------------|-------|---|
| 1 | $\langle \text{goal} \rangle$ | $::=$ | $\langle \text{expr} \rangle$ |
| 2 | $\langle \text{expr} \rangle$ | $::=$ | $\langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$ |
| 3 | | | num |
| 4 | | | id |
| 5 | $\langle \text{op} \rangle$ | $::=$ | + |
| 6 | | | - |
| 7 | | | * |
| 8 | | | / |

This describes simple expressions over numbers and identifiers.

In a BNF for a grammar, we represent

- 1 non-terminals with angle brackets or capital letters
- 2 terminals with typewriter font or underline
- 3 productions as in the example



Derivations

We can view the productions of a CFG as rewriting rules.
Using our example CFG (for $x + 2 * y$):

$$\begin{aligned}\langle \text{goal} \rangle &\Rightarrow \langle \text{expr} \rangle \\&\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{num}, 2 \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{num}, 2 \rangle * \langle \text{expr} \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle\end{aligned}$$

We have derived the sentence $x + 2 * y$.

We denote this $\langle \text{goal} \rangle \xrightarrow{*} \text{id} + \text{num} * \text{id}$.

Such a sequence of rewrites is a *derivation* or a *parse*.

The process of discovering a derivation is called *parsing*.



Derivations

*At each step, we chose a non-terminal to replace.
This choice can lead to different derivations.*

Two are of particular interest:

leftmost derivation

the leftmost non-terminal is replaced at each step

rightmost derivation

the rightmost non-terminal is replaced at each step

The previous example was a leftmost derivation.



Rightmost derivation

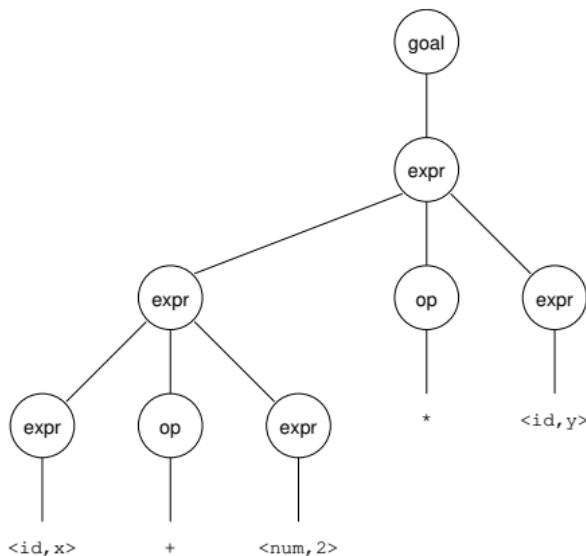
For the string $x + 2 * y$:

$$\begin{aligned}\langle \text{goal} \rangle &\Rightarrow \langle \text{expr} \rangle \\&\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\&\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle\end{aligned}$$

Again, $\langle \text{goal} \rangle \Rightarrow^* \text{id} + \text{num} * \text{id}$.



Precedence



*Treewalk evaluation computes $(x + 2) * y$*
— the “wrong” answer!
Should be $x + (2 * y)$



Precedence

*These two derivations point out a problem with the grammar.
It has no notion of precedence, or implied order of evaluation.
To add precedence takes additional machinery:*

| | | | |
|---|---------------------------------|-------|---|
| 1 | $\langle \text{goal} \rangle$ | $::=$ | $\langle \text{expr} \rangle$ |
| 2 | $\langle \text{expr} \rangle$ | $::=$ | $\langle \text{expr} \rangle + \langle \text{term} \rangle$ |
| 3 | | | $\langle \text{expr} \rangle - \langle \text{term} \rangle$ |
| 4 | | | $\langle \text{term} \rangle$ |
| 5 | $\langle \text{term} \rangle$ | $::=$ | $\langle \text{term} \rangle * \langle \text{factor} \rangle$ |
| 6 | | | $\langle \text{term} \rangle / \langle \text{factor} \rangle$ |
| 7 | | | $\langle \text{factor} \rangle$ |
| 8 | $\langle \text{factor} \rangle$ | $::=$ | num |
| 9 | | | id |

This grammar enforces a precedence on the derivation:

- terms *must* be derived from expressions
- forces the “correct” tree



Precedence

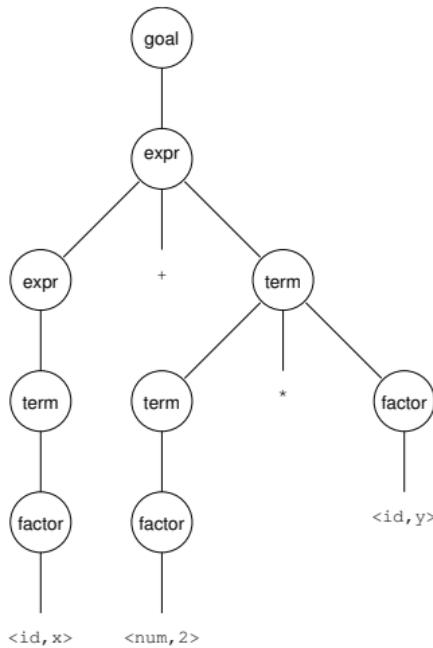
Now, for the string $x + 2 * y$:

$$\begin{aligned}\langle \text{goal} \rangle &\Rightarrow \langle \text{expr} \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{factor} \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{factor} \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{expr} \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{term} \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{factor} \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle \\&\Rightarrow \langle \text{id}, x \rangle + \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle\end{aligned}$$

Again, $\langle \text{goal} \rangle \Rightarrow^* \text{id} + \text{num} * \text{id}$, but this time, we build the desired tree.



Precedence



*Treewalk evaluation computes $x + (2 * y)$*



Ambiguity

If a grammar has more than one derivation for a single sentential form, then it is *ambiguous*

Example:

```
 $\langle \text{stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{then } \langle \text{stmt} \rangle$ 
|  $\text{if } \langle \text{expr} \rangle \text{then } \langle \text{stmt} \rangle \text{else } \langle \text{stmt} \rangle$ 
| \ other stmts
```

Consider deriving the sentential form:

if E_1 then if E_2 then S_1 else S_2

It has two derivations.

This ambiguity is purely grammatical.

It is a *context-free* ambiguity.



Ambiguity

May be able to eliminate ambiguities by rearranging the grammar:

```
<stmt>      ::=  <matched>
                  |
                  |  <unmatched>
<matched>    ::=  if <expr> then <matched> else <matched>
                  |
                  |  other stmts
<unmatched> ::=  if <expr> then <stmt>
                  |
                  |  if <expr> then <matched> else <unmatched>
```

This generates the same language as the ambiguous grammar, but applies the common sense rule:

match each else with the closest unmatched then

This is most likely the language designer's intent.



Ambiguity

Ambiguity is often due to confusion in the context-free specification.
Context-sensitive confusions can arise from *overloading*.

Example:

$$a = f(17)$$

In many Algol/Scala-like languages, f could be a function or subscripted variable. Disambiguating this statement requires context:

- need *values* of declarations
- not *context-free*
- really an issue of *type*

Rather than complicate parsing, we will handle this separately.



Scanning vs. parsing

Where do we draw the line?

```
term   ::=  [a-zA-Z][([a-zA-Z] | [0-9])]*  
        |  0 | [1-9][0-9]*  
op     ::=  + | - | * | /  
expr   ::=  (term op)*term
```

Regular expressions are used to classify:

- identifiers, numbers, keywords
- REs are more concise and simpler for tokens than a grammar
- more efficient scanners can be built from REs (DFAs) than grammars

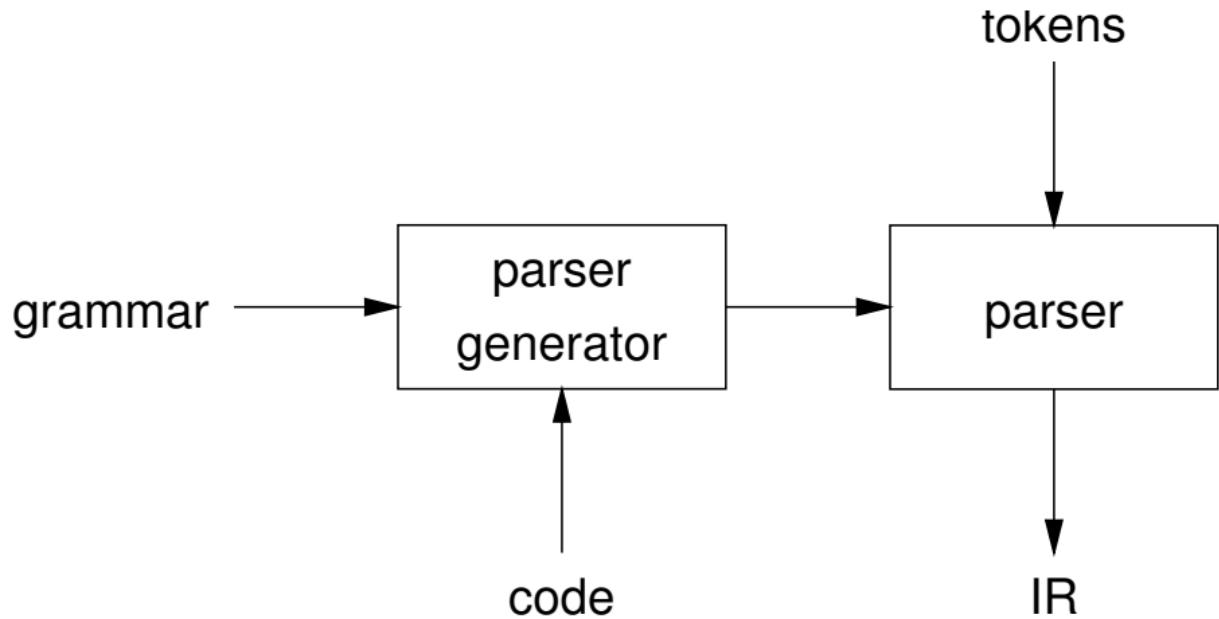
Context-free grammars are used to count:

- brackets: (), begin...end, if...then...else
- imparting structure: expressions

Syntactic analysis is complicated enough: grammar for C has around 200 productions. Factoring out lexical analysis as a separate phase makes compiler more manageable.



Parsing: the big picture



Our goal is a flexible parser generator system



Different ways of parsing: Top-down Vs Bottom-up

Top-down parsers

- start at the root of derivation tree and fill in
- picks a production and tries to match the input
- may require backtracking
- some grammars are backtrack-free (*predictive*)

Bottom-up parsers

- start at the leaves and fill in
- start in a state valid for legal first tokens
- as input is consumed, change state to encode possibilities (*recognize valid prefixes*)
- use a stack to store both state and sentential forms



Top-down parsing

A top-down parser starts with the root of the parse tree, labelled with the start or goal symbol of the grammar.

To build a parse, it repeats the following steps until the fringe of the parse tree matches the input string

- ① At a node labelled A , select a production $A \rightarrow \alpha$ and construct the appropriate child for each symbol of α
- ② When a terminal is added to the fringe that doesn't match the input string, backtrack
- ③ Find next node to be expanded (must have a label in V_n)

The key is selecting the right production in step 1.

If the parser makes a wrong step, the “derivation” process does not terminate.

Why is it bad?



Left-recursion

Top-down parsers cannot handle left-recursion in a grammar

Formally, a grammar is *left-recursive* if

$\exists A \in V_n \text{ such that } A \Rightarrow^+ A\alpha \text{ for some string } \alpha$

Our simple expression grammar is left-recursive



Eliminating left-recursion

To remove left-recursion, we can transform the grammar

Consider the grammar fragment:

$$\begin{aligned}\langle \text{foo} \rangle &::= \langle \text{foo} \rangle \alpha \\ &\quad | \\ &\quad \beta\end{aligned}$$

where α and β do not start with $\langle \text{foo} \rangle$

We can rewrite this as:

$$\begin{aligned}\langle \text{foo} \rangle &::= \beta \langle \text{bar} \rangle \\ \langle \text{bar} \rangle &::= \alpha \langle \text{bar} \rangle \\ &\quad | \\ &\quad \varepsilon\end{aligned}$$

where $\langle \text{bar} \rangle$ is a new non-terminal

This fragment contains no left-recursion



How much lookahead is needed?

We saw that top-down parsers may need to backtrack when they select the wrong production

Do we need arbitrary lookahead to parse CFGs?

- in general, yes
- use the Earley or Cocke-Younger, Kasami algorithms

Fortunately

- large subclasses of CFGs can be parsed with limited lookahead
- most programming language constructs can be expressed in a grammar that falls in these subclasses

Among the interesting subclasses are:

LL(1): left to right scan, left-most derivation, **1**-token lookahead; and

LR(1): left to right scan, reversed right-most derivation, **1**-token lookahead



Predictive parsing

Basic idea:

- For any two productions $A \rightarrow \alpha \mid \beta$, we would like a distinct way of choosing the correct production to expand.
- For some RHS $\alpha \in G$, define $\text{FIRST}(\alpha)$ as the set of tokens that appear first in some string derived from α .
- That is, for some $w \in V_t^*$, $w \in \text{FIRST}(\alpha)$ iff. $\alpha \Rightarrow^* w\gamma$.
- *Key property:*
Whenever two productions $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like
 - $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$
- This would allow the parser to make a correct choice with a lookahead of only one symbol!



Left factoring

What if a grammar does not have this property?

Sometimes, we can transform a grammar to have this property.

For each non-terminal A find the longest prefix α common to two or more of its alternatives.

if $\alpha \neq \epsilon$ then replace all of the A productions

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n$
with

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

where A' is a new non-terminal.

Repeat until no two alternatives for a single non-terminal have a common prefix.



Example

There are two non-terminals
to left factor:

$$\begin{aligned}\langle \text{expr} \rangle &::= \langle \text{term} \rangle + \langle \text{expr} \rangle \\ &\quad | \\ &\quad \langle \text{term} \rangle - \langle \text{expr} \rangle \\ &\quad | \\ &\quad \langle \text{term} \rangle\end{aligned}$$

$$\begin{aligned}\langle \text{term} \rangle &::= \langle \text{factor} \rangle * \langle \text{term} \rangle \\ &\quad | \\ &\quad \langle \text{factor} \rangle / \langle \text{term} \rangle \\ &\quad | \\ &\quad \langle \text{factor} \rangle\end{aligned}$$

Applying the transformation:

$$\begin{aligned}\langle \text{expr} \rangle &::= \langle \text{term} \rangle \langle \text{expr}' \rangle \\ \langle \text{expr}' \rangle &::= + \langle \text{expr} \rangle \\ &\quad | \\ &\quad - \langle \text{expr} \rangle \\ &\quad | \\ &\quad \varepsilon\end{aligned}$$

$$\begin{aligned}\langle \text{term} \rangle &::= \langle \text{factor} \rangle \langle \text{term}' \rangle \\ \langle \text{term}' \rangle &::= * \langle \text{term} \rangle \\ &\quad | \\ &\quad / \langle \text{term} \rangle \\ &\quad | \\ &\quad \varepsilon\end{aligned}$$



Indirect Left-recursion elimination

Given a left-factored CFG, to eliminate left-recursion:

- 1 **Input:** Grammar G with no *cycles* and no ϵ productions.
- 2 **Output:** Equivalent grammar with no left-recursion. **begin**
- 3 Arrange the non terminals in some order A_1, A_2, \dots, A_n ;
- 4 **foreach** $i = 1 \dots n$ **do**
- 5 **foreach** $j = 1 \dots i - 1$ **do**
- 6 Say the i^{th} production is: $A_i \rightarrow A_j \gamma$;
- 7 and $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$;
- 8 Replace, the i^{th} production by:
- 9 $A_i \rightarrow \delta_1 \gamma | \delta_2 \gamma | \dots | \delta_n \gamma$;
- 10 Eliminate immediate left recursion in A_i ;



Generality

Question:

By left factoring and eliminating left-recursion, can we transform an arbitrary context-free grammar to a form where it can be predictively parsed with a single token lookahead?

Answer:

Given a context-free grammar that doesn't meet our conditions, it is undecidable whether an equivalent grammar exists that does meet our conditions.

Many *context-free languages* do not have such a grammar:

$$\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}$$

Must look past an arbitrary number of *a*'s to discover the 0 or the 1 and so determine the derivation.



Recursive descent parsing

```
1 int A()
2 begin
3     foreach production of the form  $A \rightarrow X_1X_2X_3 \cdots X_k$  do
4         for  $i = 1$  to  $k$  do
5             if  $X_i$  is a non-terminal then
6                 if  $(X_i()) \neq 0$  then
7                     backtrack; break; // Try the next
                           production
8             else if  $X_i$  matches the current input symbol  $a$  then
9                 advance the input to the next symbol;
10            else
11                backtrack; break; // Try the next production
12            if  $i \neq k+1$  then
13                return 0; // Success
14        return 1; // Failure
```



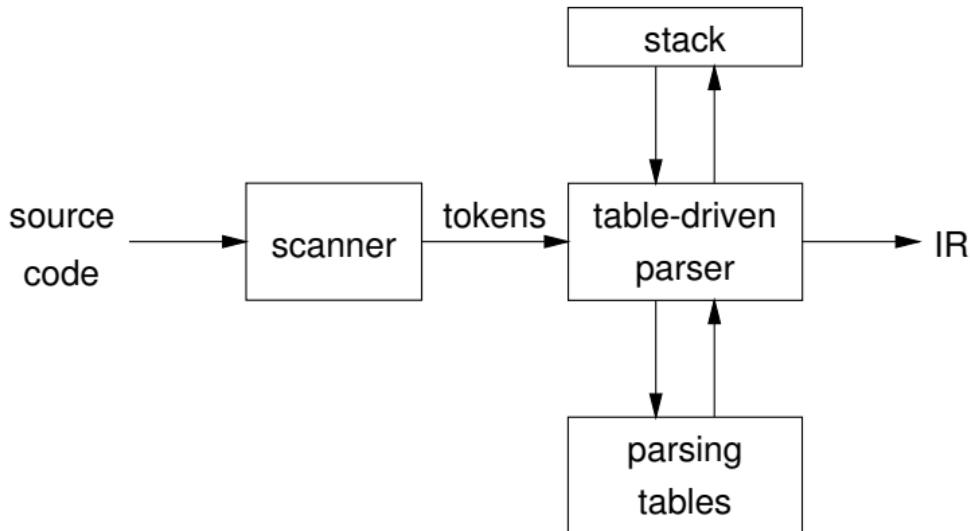
Recursive descent parsing

- Backtracks in general – in practise may not do much.
- How to backtrack?
- A left recursive grammar will lead to infinite loop.



Non-recursive predictive parsing

Now, a predictive parser looks like:

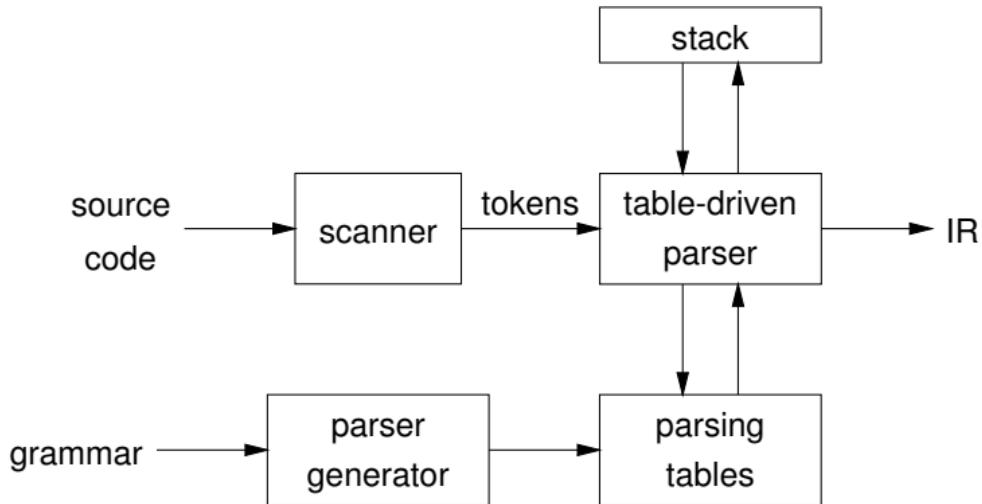


Rather than writing recursive code, we build tables.
Why? *Building tables can be automated, easily.*



Table-driven parsers

A parser generator system often looks like:



- This is true for both top-down (LL) and bottom-up (LR) parsers
- This also uses a stack – but mainly to remember part of the input string; no recursion.



FIRST

For a string of grammar symbols α , define $\text{FIRST}(\alpha)$ as:

- the set of terminals that begin strings derived from α :
 $\{a \in V_t \mid \alpha \Rightarrow^* a\beta\}$
- If $\alpha \Rightarrow^* \epsilon$ then $\epsilon \in \text{FIRST}(\alpha)$

$\text{FIRST}(\alpha)$ contains the tokens valid in the initial position in α

To build $\text{FIRST}(X)$:

- If $X \in V_t$ then $\text{FIRST}(X)$ is $\{X\}$
- If $X \rightarrow \epsilon$ then add ϵ to $\text{FIRST}(X)$
- If $X \rightarrow Y_1 Y_2 \cdots Y_k$:
 - Put $\text{FIRST}(Y_1) - \{\epsilon\}$ in $\text{FIRST}(X)$
 - $\forall i : 1 < i \leq k$, if $\epsilon \in \text{FIRST}(Y_1) \cap \cdots \cap \text{FIRST}(Y_{i-1})$
(i.e., $Y_1 \cdots Y_{i-1} \Rightarrow^* \epsilon$)
then put $\text{FIRST}(Y_i) - \{\epsilon\}$ in $\text{FIRST}(X)$
 - If $\epsilon \in \text{FIRST}(Y_1) \cap \cdots \cap \text{FIRST}(Y_k)$ then put ϵ in $\text{FIRST}(X)$

Repeat until no more additions can be made.



FOLLOW

For a non-terminal A , define $\text{FOLLOW}(A)$ as

the set of terminals that can appear immediately to the right of A in some sentential form

Thus, a non-terminal's FOLLOW set specifies the tokens that can legally appear after it.

A terminal symbol has no FOLLOW set.

To build $\text{FOLLOW}(A)$:

- ① Put $\$$ in $\text{FOLLOW}(\langle \text{goal} \rangle)$
- ② If $A \rightarrow \alpha B \beta$:
 - ① Put $\text{FIRST}(\beta) - \{\epsilon\}$ in $\text{FOLLOW}(B)$
 - ② If $\beta = \epsilon$ (i.e., $A \rightarrow \alpha B$) or $\epsilon \in \text{FIRST}(\beta)$ (i.e., $\beta \Rightarrow^* \epsilon$) then put $\text{FOLLOW}(A)$ in $\text{FOLLOW}(B)$

Repeat until no more additions can be made



Previous definition

A grammar G is LL(1) iff. for all non-terminals A , each distinct pair of productions $A \rightarrow \beta$ and $A \rightarrow \gamma$ satisfy the condition $\text{FIRST}(\beta) \cap \text{FIRST}(\gamma) = \phi$.

What if $A \Rightarrow^* \epsilon$?

Revised definition

A grammar G is LL(1) iff. for each set of productions $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$:

- ① $\text{FIRST}(\alpha_1), \text{FIRST}(\alpha_2), \dots, \text{FIRST}(\alpha_n)$ are all pairwise disjoint
- ② If $\alpha_i \Rightarrow^* \epsilon$ then $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \phi, \forall 1 \leq j \leq n, i \neq j$.

If G is ϵ -free, condition 1 is sufficient.



Provable facts about LL(1) grammars:

- 1 No left-recursive grammar is LL(1)
- 2 No ambiguous grammar is LL(1)
- 3 Some languages have no LL(1) grammar
- 4 A ϵ -free grammar where each alternative expansion for A begins with a distinct terminal is a *simple* LL(1) grammar.

Example

- $S \rightarrow aS \mid a$ is not LL(1) because $\text{FIRST}(aS) = \text{FIRST}(a) = \{a\}$
- $S \rightarrow aS'$
 $S' \rightarrow aS' \mid \epsilon$
accepts the same language and is LL(1)



LL(1) parse table construction

Input: Grammar G

Output: Parsing table M

Method:

① \forall productions $A \rightarrow \alpha$:

 ① $\forall a \in \text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$

 ② If $\epsilon \in \text{FIRST}(\alpha)$:

 ① $\forall b \in \text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$

 ② If $\$ \in \text{FOLLOW}(A)$ then add $A \rightarrow \alpha$ to $M[A, \$]$

② Set each undefined entry of M to `error`

If $\exists M[A, a]$ with multiple entries then grammar is not LL(1).

Note: recall $a, b \in V_t$, so $a, b \neq \epsilon$



Example

Our long-suffering expression grammar:

$$\begin{array}{lll|lll} 1. & S & \rightarrow E & 6. & T & \rightarrow FT' \\ 2. & E & \rightarrow TE' & 7. & T' & \rightarrow *T \\ 3. & E' & \rightarrow +E & 8. & & | /T \\ 4. & & | -E & 9. & & | \epsilon \\ 5. & & | \epsilon & 10. & F & \rightarrow \text{num} \\ & & & & & | \text{id} \end{array}$$

| | FIRST | FOLLOW | id | num | + | - | * | / | \$ |
|------|------------------|--------------|----|-----|---|---|---|---|----|
| S | num, id | \$ | 1 | 1 | — | — | — | — | — |
| E | num, id | \$ | 2 | 2 | — | — | — | — | — |
| E' | $\epsilon, +, -$ | \$ | — | — | 3 | 4 | — | — | 5 |
| T | num, id | $+,-,\$$ | 6 | 6 | — | — | — | — | — |
| T' | $\epsilon, *, /$ | $+,-,\$$ | — | — | 9 | 9 | 7 | 8 | 9 |
| F | num, id | $+,-,*,/,\$$ | 11 | 10 | — | — | — | — | — |
| id | id | — | | | | | | | |
| num | num | — | | | | | | | |
| * | * | — | | | | | | | |
| / | / | — | | | | | | | |
| + | + | — | | | | | | | |
| — | — | — | | | | | | | |



Table driven Predictive parsing

Input: A string w and a parsing table M for a grammar G

Output: If w is in $L(G)$, a leftmost derivation of w ; otherwise, indicate an error

- 1 push $\$$ onto the stack; push S onto the stack;
- 2 a points to the input tape;
- 3 $X = \text{stack.top}();$
- 4 **while** $X \neq \$$ **do**
 - 5 **if** X is a **then**
 - 6 $\text{stack.pop}(); \text{inp}++;$
 - 7 **else if** X is a terminal **then**
 - 8 $\text{error}();$
 - 9 **else if** $M[X, a]$ is an error entry **then**
 - 10 $\text{error}();$
 - 11 **else if** $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$ **then**
 - 12 output the production $X \rightarrow Y_1 Y_2 \dots Y_k$;
 - 13 $\text{stack.pop}();$
 - 14 push Y_k, Y_{k-1}, \dots, Y_1 in that order;
 - 15 $X = \text{stack.top}();$



A grammar that is not LL(1)

```
 $\langle \text{stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle$ 
|  $\text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$ 
| ...
```

Left-factored: $\langle \text{stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \langle \text{stmt}' \rangle | \dots$ Now,
 $\langle \text{stmt}' \rangle ::= \text{else } \langle \text{stmt} \rangle | \epsilon$

$\text{FIRST}(\langle \text{stmt}' \rangle) = \{\epsilon, \text{else}\}$

Also, $\text{FOLLOW}(\langle \text{stmt}' \rangle) = \{\text{else}, \$\}$

But, $\text{FIRST}(\langle \text{stmt}' \rangle) \cap \text{FOLLOW}(\langle \text{stmt}' \rangle) = \{\text{else}\} \neq \emptyset$

On seeing `else`, there is a conflict between choosing

$\langle \text{stmt}' \rangle ::= \text{else } \langle \text{stmt} \rangle \text{ and } \langle \text{stmt}' \rangle ::= \epsilon$

\Rightarrow grammar is not LL(1)!

The fix:

Put priority on $\langle \text{stmt}' \rangle ::= \text{else } \langle \text{stmt} \rangle$ to associate `else` with closest previous `then`.



Another example of painful left-factoring

- Here is a typical example where a programming language fails to be LL(1):

$$\begin{aligned} \text{stmt} &\rightarrow \text{asginment} \mid \text{call} \mid \text{other} \\ \text{assignment} &\rightarrow \text{id} := \text{exp} \\ \text{call} &\rightarrow \text{id} (\text{exp-list}) \end{aligned}$$

- This grammar is not in a form that can be left factored. We must first replace assignment and call by the right-hand sides of their defining productions:

$$\text{statement} \rightarrow \text{id} := \text{exp} \mid \text{id} (\text{exp-list}) \mid \text{other}$$

- We left factor:

$$\begin{aligned} \text{statement} &\rightarrow \text{id} \text{ stmt}' \mid \text{other} \\ \text{stmt}' &\rightarrow := \text{exp} \mid (\text{exp-list}) \end{aligned}$$

- See how the grammar obscures the language semantics.



Error recovery in Predictive Parsing

- An error is detected when the terminal on top of the stack does not match the next input symbol or $M[A, a] = \text{error}$.

Panic mode error recovery

- Skip input symbols till a “synchronizing” token appears.

Q: How to identify a synchronizing token?

Some heuristics:

- All symbols in $\text{FOLLOW}(A)$ in the synchronizing set for the non-terminal A .
- Semicolon after a Stmt production: assignmentStmt;
assignmentStmt;
- If a terminal on top of the stack cannot be matched? –
 - pop the terminal.
 - issue a message that the terminal was inserted.

Q: How about error messages?



Some definitions

Recall

- For a grammar G , with start symbol S , any string α such that $S \Rightarrow^* \alpha$ is called a *sentential form*
- If $\alpha \in V_t^*$, then α is called a *sentence* in $L(G)$
- Otherwise it is just a sentential form (not a sentence in $L(G)$)

A *left-sentential form* is a sentential form that occurs in the leftmost derivation of some sentence.

A *right-sentential form* is a sentential form that occurs in the rightmost derivation of some sentence.

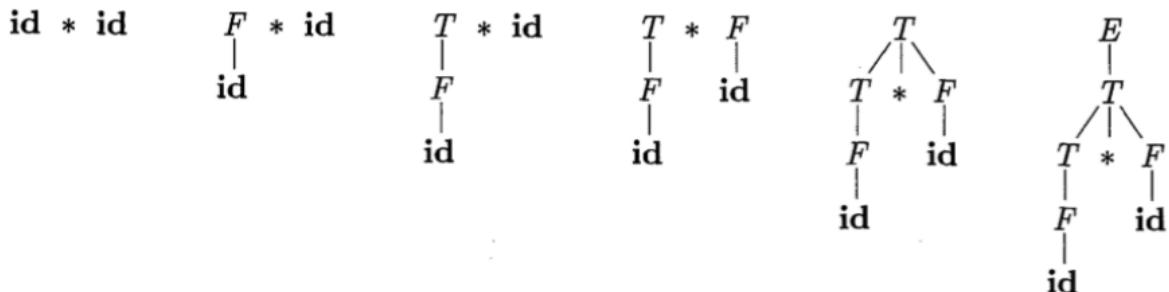
An unambiguous grammar will have a unique leftmost/rightmost derivation.



Bottom-up parsing

Goal:

Given an input string w and a grammar G , construct a parse tree by starting at the leaves and working to the root.



Reductions Vs Derivations

Reduction:

- At each reduction step, a specific substring matching the body of a production is replaced by the non-terminal at the head of the production.

Key decisions

- When to reduce?
- What production rule to apply?

Reduction Vs Derivations

- Recall: In derivation: a non-terminal in a sentential form is replaced by the body of one of its productions.
- A reduction is reverse of a step in derivation.
- Bottom-up parsing is the process of “reducing” a string w to the start symbol.
- Goal of bottom-up parsing: build derivation tree in reverse.



Example

Consider the grammar

$$\begin{array}{rcl} 1 & S & \rightarrow aABe \\ 2 & A & \rightarrow Abc \\ 3 & & | \\ 4 & B & \rightarrow d \end{array}$$

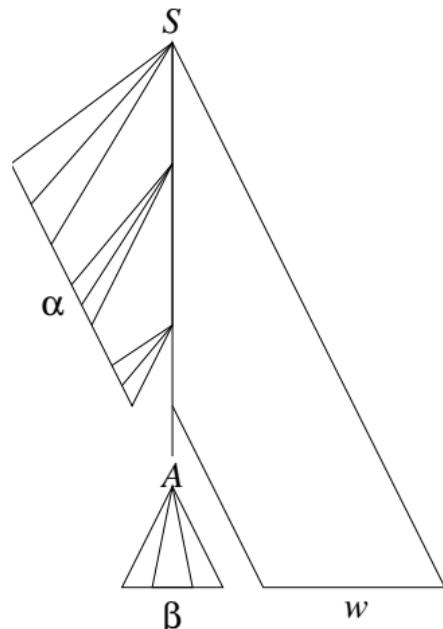
and the input string abbcde

| Prod'n. | Sentential Form |
|---------|-----------------|
| 3 | a b bcde |
| 2 | a Abc de |
| 4 | aA d e |
| 1 | aABe |
| - | S |

The trick appears to be scanning the input and finding valid sentential forms.



Handles



The handle $A \rightarrow \beta$ in the parse tree
for $\alpha\beta w$

Informally, a “handle” is

- a substring that matches the body of a production (not necessarily the first one),
- and reducing this handle, represents one step of reduction (or reverse rightmost derivation).



Handles

Theorem:

If G is unambiguous then every right-sentential form has a unique handle.

Proof: (by definition)

- ① G is unambiguous \Rightarrow rightmost derivation is unique
- ② \Rightarrow a unique production $A \rightarrow \beta$ applied to take γ_{i-1} to γ_i
- ③ \Rightarrow a unique position k at which $A \rightarrow \beta$ is applied
- ④ \Rightarrow a unique handle $A \rightarrow \beta$



Example

The left-recursive expression grammar (*original form*)

| | |
|---|---|
| 1 | $\langle \text{goal} \rangle ::= \langle \text{expr} \rangle$ |
| 2 | $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$ |
| 3 | $\langle \text{expr} \rangle - \langle \text{term} \rangle$ |
| 4 | $\langle \text{term} \rangle$ |
| 5 | $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$ |
| 6 | $\langle \text{term} \rangle / \langle \text{factor} \rangle$ |
| 7 | $\langle \text{factor} \rangle$ |
| 8 | $\langle \text{factor} \rangle ::= \text{num}$ |
| 9 | id |

| Prod'n. | Sentential Form |
|---------|---|
| - | $\langle \text{goal} \rangle$ |
| 1 | <u>$\langle \text{expr} \rangle$</u> |
| 3 | <u>$\langle \text{expr} \rangle - \langle \text{term} \rangle$</u> |
| 5 | <u>$\langle \text{expr} \rangle - \langle \text{term} \rangle * \langle \text{factor} \rangle$</u> |
| 9 | $\langle \text{expr} \rangle - \langle \text{term} \rangle * \underline{\text{id}}$ |
| 7 | $\langle \text{expr} \rangle - \underline{\langle \text{factor} \rangle} * \text{id}$ |
| 8 | $\langle \text{expr} \rangle - \underline{\text{num}} * \text{id}$ |
| 4 | $\langle \text{term} \rangle - \text{num} * \text{id}$ |
| 7 | $\langle \text{factor} \rangle - \text{num} * \text{id}$ |
| 9 | <u>$\text{id} - \text{num} * \text{id}$</u> |



Handle-pruning

The process to construct a bottom-up parse is called *handle-pruning*.
To construct a rightmost derivation

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n = w$$

we set i to n and apply the following simple algorithm

for $i = n$ downto 1

- ① find the handle $A_i \rightarrow \beta_i$ in γ_i
- ② replace β_i with A_i to generate γ_{i-1}

This takes $2n$ steps, where n is the length of the derivation



Stack implementation

One scheme to implement a handle-pruning, bottom-up parser is called a *shift-reduce* parser.

Shift-reduce parsers use a *stack* and an *input buffer*

- 1 initialize stack with \$
- 2 Repeat until the top of the stack is the goal symbol and the input token is \$
 - a) *find the handle*
if we don't have a handle on top of the stack, *shift* an input symbol onto the stack
 - b) *prune the handle*
if we have a handle $A \rightarrow \beta$ on the stack, *reduce*
 - i) pop $|\beta|$ symbols off the stack
 - ii) push A onto the stack



Example: back to $x - 2 * y$

| | Stack | Input | Action |
|---|----------------------------|-------------------------------|--------|
| 1 | $S \rightarrow E$ | | |
| 2 | $E \rightarrow E + T$ | \$ | S |
| 3 | $ E - T$ | \$ id | R9 |
| 4 | $ T$ | \$ <factor> | R7 |
| 5 | $T \rightarrow T * F$ | \$ <term> | R4 |
| 6 | $ T / F$ | \$ <expr> | S |
| 7 | $ F$ | \$ <expr> - | S |
| 8 | $F \rightarrow \text{num}$ | \$ <expr> - num | R8 |
| 9 | $ \text{id}$ | \$ <expr> - <factor> | R7 |
| | | \$ <expr> - <term> | * id |
| | | \$ <expr> - <term> * | S |
| | | \$ <expr> - <term> * id | S |
| | | \$ <expr> - <term> * <factor> | R9 |
| | | \$ <expr> - <term> | R5 |
| | | \$ <expr> | R3 |
| | | \$ <goal> | R1 |
| | | | A |



Shift-reduce parsing

Shift-reduce parsers are simple to understand

A shift-reduce parser has just four canonical actions:

- ① *shift* — next input symbol is shifted onto the top of the stack
- ② *reduce* — right end of handle is on top of stack;
locate left end of handle within the stack;
pop handle off stack and push appropriate non-terminal LHS
- ③ *accept* — terminate parsing and signal success
- ④ *error* — call an error recovery routine

Key insight: recognize handles with a DFA:

- DFA transitions shift states instead of symbols
- accepting states trigger reductions

May have Shift-Reduce Conflicts.



LR parsing

The skeleton parser:

```
push  $s_0$ 
token  $\leftarrow$  next_token()
repeat forever
    s  $\leftarrow$  top of stack
    if action[s, token] = "shift  $s_i$ " then
        push  $s_i$ 
        token  $\leftarrow$  next_token()
    else if action[s, token] = "reduce  $A \rightarrow \beta$ " then
        pop | $\beta$ | states
        s'  $\leftarrow$  top of stack
        push goto[s', A]
    else if action[s, token] = "accept" then
        return
    else error()
```

“How many ops?”: k shifts, l reduces, and 1 accept, where k is length of input string and l is length of reverse rightmost derivation



Example tables

| state | ACTION | GOTO | | |
|-------|------------|------|-----|-----|
| | id + * \$ | E | T | F |
| 0 | s4 - - - | 1 | 2 | 3 |
| 1 | - - - acc | --- | --- | --- |
| 2 | - s5 - r3 | --- | --- | --- |
| 3 | - r5 s6 r5 | --- | --- | --- |
| 4 | - r6 r6 r6 | --- | --- | --- |
| 5 | s4 - - - | 7 | 2 | 3 |
| 6 | s4 - - - | - | 8 | 3 |
| 7 | - - - r2 | --- | --- | --- |
| 8 | - r4 - r4 | --- | --- | --- |

The Grammar

| | |
|---|---------------------------|
| 1 | $S \rightarrow E$ |
| 2 | $E \rightarrow T + E$ |
| 3 | T |
| 4 | $T \rightarrow F * T$ |
| 5 | F |
| 6 | $F \rightarrow \text{id}$ |

Note: This is a simple little right-recursive grammar. It is *not* the same grammar as in previous lectures.



Example using the tables

| Stack | Input | Action |
|------------|--------------|--------|
| \$ 0 | id* id+ id\$ | s4 |
| \$ 0 4 | * id+ id\$ | r6 |
| \$ 0 3 | * id+ id\$ | s6 |
| \$ 0 3 6 | id+ id\$ | s4 |
| \$ 0 3 6 4 | + id\$ | r6 |
| \$ 0 3 6 3 | + id\$ | r5 |
| \$ 0 3 6 8 | + id\$ | r4 |
| \$ 0 2 | + id\$ | s5 |
| \$ 0 2 5 | id\$ | s4 |
| \$ 0 2 5 4 | \$ | r6 |
| \$ 0 2 5 3 | \$ | r5 |
| \$ 0 2 5 2 | \$ | r3 |
| \$ 0 2 5 7 | \$ | r2 |
| \$ 0 1 | \$ | acc |



Informally, we say that a grammar G is LR(k) if, given a rightmost derivation

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \cdots \Rightarrow \gamma_n = w,$$

we can, for each right-sentential form in the derivation:

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

by scanning γ_i from left to right, going at most k symbols beyond the right end of the handle of γ_i .



LR(k) grammars

Formally, a grammar G is LR(k) iff.:

- 1 $S \Rightarrow_{\text{rm}}^* \alpha Aw \Rightarrow_{\text{rm}} \alpha \beta w$, and
- 2 $S \Rightarrow_{\text{rm}}^* \gamma Bx \Rightarrow_{\text{rm}} \alpha \beta y$, and
- 3 $\text{FIRST}_k(w) = \text{FIRST}_k(y)$

$$\Rightarrow \alpha Ay = \gamma Bx$$

i.e., Assume sentential forms $\alpha \beta w$ and $\alpha \beta y$, with common prefix $\alpha \beta$ and common k -symbol lookahead $\text{FIRST}_k(y) = \text{FIRST}_k(w)$, such that $\alpha \beta w$ reduces to αAw and $\alpha \beta y$ reduces to γBx .

But, the common prefix means $\alpha \beta y$ also reduces to αAy , for the same result.

Thus $\alpha Ay = \gamma Bx$.



Why study LR grammars?

LR(1) grammars are often used to construct parsers.

We call these parsers LR(1) parsers.

- virtually all context-free programming language constructs can be expressed in an LR(1) form
- LR grammars are the most general grammars parsable by a deterministic, bottom-up parser
- efficient parsers can be implemented for LR(1) grammars
- LR parsers detect an error as soon as possible in a left-to-right scan of the input
- LR grammars describe a proper superset of the languages recognized by predictive (i.e., LL) parsers

LL(k): recognize use of a production $A \rightarrow \beta$ seeing first k symbols derived from β

LR(k): recognize the handle β after seeing everything derived from β plus k lookahead symbols



LR parsing

Three common algorithms to build tables for an “LR” parser:

① SLR(1)

- smallest class of grammars
- smallest tables (number of states)
- simple, fast construction

② LR(1)

- full set of LR(1) grammars
- largest tables (number of states)
- slow, large construction

③ LALR(1)

- intermediate sized set of grammars
- same number of states as SLR(1)
- canonical construction is slow and large
- better construction techniques exist



SLR vs. LR/LALR

An LR(1) parser for either Algol or Pascal has several thousand states, while an SLR(1) or LALR(1) parser for the same language may have several hundred states.



LR(k) items

The table construction algorithms use sets of LR(k) *items* or *configurations* to represent the possible states in a parse.

An LR(k) item is a pair $[\alpha, \beta]$, where

- α is a production from G with a \bullet at some position in the RHS, marking how much of the RHS of a production has already been seen
- β is a lookahead string containing k symbols (terminals or $\$$)

Two cases of interest are $k = 0$ and $k = 1$:

- LR(0) items play a key role in the SLR(1) table construction algorithm.
- LR(1) items play a key role in the LR(1) and LALR(1) table construction algorithms.



Example

The • indicates how much of an item we have seen at a given state in the parse:

$[A \rightarrow \bullet XYZ]$ indicates that the parser is looking for a string that can be derived from XYZ

$[A \rightarrow XY \bullet Z]$ indicates that the parser has seen a string derived from XY and is looking for one derivable from Z

LR(0) items: (*no lookahead*)

$A \rightarrow XYZ$ generates 4 LR(0) items:

- ① $[A \rightarrow \bullet XYZ]$
- ② $[A \rightarrow X \bullet YZ]$
- ③ $[A \rightarrow XY \bullet Z]$
- ④ $[A \rightarrow XYZ \bullet]$



The characteristic finite state machine (CFSM)

The CFSM for a grammar is a DFA which recognizes *viable prefixes* of right-sentential forms:

A viable prefix is any prefix that does not extend beyond the handle.

It accepts when a handle has been discovered and needs to be reduced.

To construct the CFSM we need two functions:

- CLOSURE(I) to build its states
- GOTO(I, X) to determine its transitions



CLOSURE

Given an item $[A \rightarrow \alpha \bullet B\beta]$, its closure contains the item and any other items that can generate legal substrings to follow α .
Thus, if the parser has viable prefix α on its stack, the input should reduce to $B\beta$ (or γ for some other item $[B \rightarrow \bullet\gamma]$ in the closure).

```
function CLOSURE ( $I$ )
repeat
    if  $[A \rightarrow \alpha \bullet B\beta] \in I$ 
        add  $[B \rightarrow \bullet\gamma]$  to  $I$ 
until no more items can be added to  $I$ 
return  $I$ 
```



GOTO

Let I be a set of LR(0) items and X be a grammar symbol.

Then, $\text{GOTO}(I, X)$ is the closure of the set of all items

$[A \rightarrow \alpha X \bullet \beta]$ such that $[A \rightarrow \alpha \bullet X \beta] \in I$

If I is the set of valid items for some viable prefix γ , then $\text{GOTO}(I, X)$ is the set of valid items for the viable prefix γX .

$\text{GOTO}(I, X)$ represents state after recognizing X in state I .

function $\text{GOTO}(I, X)$

 let J be the set of items $[A \rightarrow \alpha X \bullet \beta]$

 such that $[A \rightarrow \alpha \bullet X \beta] \in I$

 return $\text{CLOSURE}(J)$



Building the LR(0) item sets

We start the construction with the item $[S' \rightarrow \bullet S \$]$, where

- S' is the start symbol of the augmented grammar G'
- S is the start symbol of G
- $\$$ represents EOF

To compute the collection of sets of LR(0) items

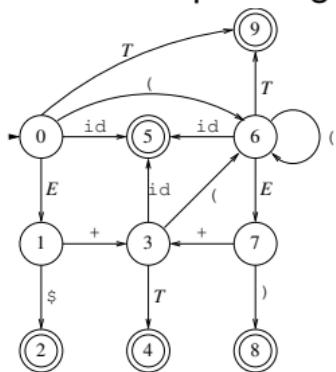
```
function items ( $G'$ )
     $s_0 \leftarrow \text{CLOSURE}(\{[S' \rightarrow \bullet S \$]\})$ 
     $C \leftarrow \{s_0\}$ 
    repeat
        for each set of items  $s \in C$ 
            for each grammar symbol  $X$ 
                if  $\text{GOTO}(s, X) \neq \emptyset$  and  $\text{GOTO}(s, X) \notin C$ 
                    add  $\text{GOTO}(s, X)$  to  $C$ 
    until no more item sets can be added to  $C$ 
    return  $C$ 
```



LR(0) example

| | |
|---|---------------------------|
| 1 | $S \rightarrow E\$$ |
| 2 | $E \rightarrow E + T$ |
| 3 | T |
| 4 | $T \rightarrow \text{id}$ |
| 5 | (E) |

The corresponding CFSM:



| | |
|-------------------------------------|---|
| $I_0 : S \rightarrow \bullet E \$$ | $I_4 : E \rightarrow E + T \bullet$ |
| $E \rightarrow \bullet E + T$ | $I_5 : T \rightarrow \text{id} \bullet$ |
| $E \rightarrow \bullet T$ | $I_6 : T \rightarrow (\bullet E)$ |
| $T \rightarrow \bullet \text{id}$ | $E \rightarrow \bullet E + T$ |
| $T \rightarrow \bullet (E)$ | $E \rightarrow \bullet T$ |
| $I_1 : S \rightarrow E \bullet \$$ | $T \rightarrow \bullet \text{id}$ |
| $E \rightarrow E \bullet + T$ | $T \rightarrow \bullet (E)$ |
| $I_2 : S \rightarrow E \$ \bullet$ | $I_7 : T \rightarrow (E \bullet)$ |
| $I_3 : E \rightarrow E + \bullet T$ | $E \rightarrow E \bullet + T$ |
| $T \rightarrow \bullet \text{id}$ | $I_8 : T \rightarrow (E) \bullet$ |
| $T \rightarrow \bullet (E)$ | $I_9 : E \rightarrow T \bullet$ |

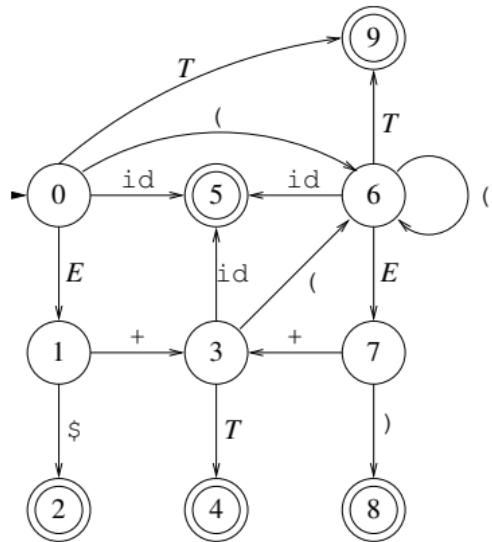


Constructing the LR(0) parsing table

- 1 construct the collection of sets of LR(0) items for G'
- 2 state i of the CFSM is constructed from I_i
 - 1 $[A \rightarrow \alpha \bullet a\beta] \in I_i$ and $\text{GOTO}(I_i, a) = I_j$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"shift } j\text{"}$
 - 2 $[A \rightarrow \alpha \bullet] \in I_i, A \neq S'$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"reduce } A \rightarrow \alpha\text", \forall a$
 - 3 $[S' \rightarrow S\$ \bullet] \in I_i$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"accept"}, \forall a$
- 3 $\text{GOTO}(I_i, A) = I_j$
 $\Rightarrow \text{GOTO}[i, A] \leftarrow j$
- 4 set undefined entries in ACTION and GOTO to "error"
- 5 initial state of parser s_0 is $\text{CLOSURE}([S' \rightarrow \bullet S\$])$



LR(0) example



| state | ACTION | | | | | GOTO | |
|-------|---------------------|----------------|----------------|-----------------|----|----------------|----------------|
| | <code>id</code> | <code>(</code> | <code>+</code> | <code>\$</code> | | <code>S</code> | <code>E</code> |
| 0 | s5 s6 | - | - | - | - | -1 | 9 |
| 1 | - | - | - | s3 s2 | - | - | - |
| 2 | acc acc acc acc acc | - | - | - | - | - | - |
| 3 | s5 s6 | - | - | - | - | - | 4 |
| 4 | r2 r2 | r2 | r2 | r2 | r2 | - | - |
| 5 | r4 r4 | r4 | r4 | r4 | r4 | - | - |
| 6 | s5 s6 | - | - | - | - | 7 | 9 |
| 7 | - | - | s8 s3 | - | - | - | - |
| 8 | r5 r5 | r5 | r5 | r5 | r5 | - | - |
| 9 | r3 r3 | r3 | r3 | r3 | r3 | - | - |



Conflicts in the ACTION table

If the LR(0) parsing table contains any multiply-defined ACTION entries then G is not LR(0)

Two conflicts arise:

shift-reduce: both shift and reduce possible in same item set

reduce-reduce: more than one distinct reduce action possible in same item set

Conflicts can be resolved through *lookahead* in ACTION. Consider:

- $A \rightarrow \epsilon \mid a\alpha$
 \Rightarrow shift-reduce conflict
- $a := b + c * d$
 requires lookahead to avoid shift-reduce conflict after shifting c
 (need to see $*$ to give precedence over $+$)



SLR(1): simple lookahead LR

Add lookahead after building LR(0) item sets

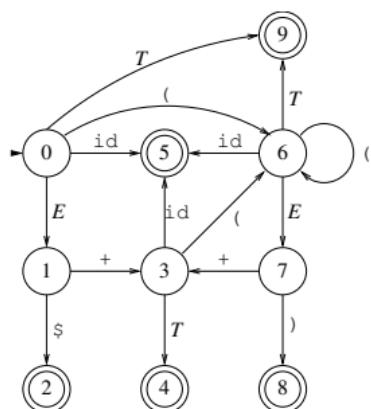
Constructing the SLR(1) parsing table:

- ① construct the collection of sets of LR(0) items for G'
- ② state i of the CFSM is constructed from I_i
 - ① $[A \rightarrow \alpha \bullet a\beta] \in I_i$ and $\text{GOTO}(I_i, a) = I_j$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"shift } j\text{"}, \underline{\forall a \neq \$}$
 - ② $[A \rightarrow \alpha \bullet] \in I_i, A \neq S'$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"reduce } A \rightarrow \alpha\text", \underline{\forall a \in \text{FOLLOW}(A)}$
 - ③ $[S' \rightarrow S \bullet \$] \in I_i$
 $\Rightarrow \text{ACTION}[i, \$] \leftarrow \text{"accept"}$
- ③ $\text{GOTO}(I_i, A) = I_j$
 $\Rightarrow \text{GOTO}[i, A] \leftarrow j$
- ④ set undefined entries in ACTION and GOTO to "error"
- ⑤ initial state of parser s_0 is $\text{CLOSURE}([S' \rightarrow \bullet S\$])$



From previous example

| | |
|---|-----------------------|
| 1 | $S \rightarrow E\$$ |
| 2 | $E \rightarrow E + T$ |
| 3 | T |
| 4 | $T \rightarrow id$ |
| 5 | (E) |



$$\text{FOLLOW}(E) = \text{FOLLOW}(T) = \{\$, +,)\}$$

| state | ACTION | GOTO |
|-------|--------------|--------------|
| | id () + \\$ | <i>S E T</i> |
| 0 | s5 s6 - - - | - 1 9 |
| 1 | - - - s3 acc | - - - |
| 2 | - - - - - | - - - |
| 3 | s5 s6 - - - | - - 4 |
| 4 | - - r2 r2 r2 | - - - |
| 5 | - - r4 r4 r4 | - - - |
| 6 | s5 s6 - - - | - 7 9 |
| 7 | - - s8 s3 - | - - - |
| 8 | - - r5 r5 r5 | - - - |
| 9 | - - r3 r3 r3 | - - - |



Example: A grammar that is not LR(0)

| | |
|---|---------------------------|
| 1 | $S \rightarrow E\$$ |
| 2 | $E \rightarrow E + T$ |
| 3 | T |
| 4 | $T \rightarrow T * F$ |
| 5 | F |
| 6 | $F \rightarrow \text{id}$ |
| 7 | (E) |

| | FOLLOW |
|-----|-------------------|
| E | $\{+,), \$\}$ |
| T | $\{+, *,), \$\}$ |
| F | $\{+, *,), \$\}$ |

| | |
|---|--|
| $I_0 : S \rightarrow \bullet E \$$ | $I_6 : F \rightarrow (\bullet E)$ |
| $E \rightarrow \bullet E + T$ | $E \rightarrow \bullet E + T$ |
| $E \rightarrow \bullet T$ | $E \rightarrow \bullet T$ |
| $T \rightarrow \bullet T * F$ | $T \rightarrow \bullet T * F$ |
| $T \rightarrow \bullet F$ | $T \rightarrow \bullet F$ |
| $F \rightarrow \bullet \text{id}$ | $F \rightarrow \bullet \text{id}$ |
| $F \rightarrow \bullet (E)$ | $F \rightarrow \bullet (E)$ |
| $I_1 : S \rightarrow E \bullet \$$ | $I_7 : E \rightarrow T \bullet$ |
| $E \rightarrow E \bullet + T$ | $T \rightarrow T \bullet * F$ |
| $I_2 : S \rightarrow E \$ \bullet$ | $I_8 : T \rightarrow T * \bullet F$ |
| $I_3 : E \rightarrow E + \bullet T$ | $F \rightarrow \bullet \text{id}$ |
| $T \rightarrow \bullet T * F$ | $F \rightarrow \bullet (E)$ |
| $T \rightarrow \bullet F$ | $I_9 : T \rightarrow T * F \bullet$ |
| $F \rightarrow \bullet \text{id}$ | $I_{10} : F \rightarrow (E) \bullet$ |
| $F \rightarrow \bullet (E)$ | $I_{11} : E \rightarrow E + T \bullet$ |
| $I_4 : T \rightarrow F \bullet$ | $T \rightarrow T \bullet * F$ |
| $I_5 : F \rightarrow \text{id} \bullet$ | $I_{12} : F \rightarrow (E \bullet)$ |
| | $E \rightarrow E \bullet + T$ |



Example: But it is SLR(1)

| state | ACTION | | | | | | GOTO | | | |
|-------|--------|----|----|----|-----|-----|------|----|----|---|
| | + | * | id | (|) | \$ | S | E | T | F |
| 0 | - | - | s5 | s6 | - | - | - | 1 | 7 | 4 |
| 1 | s3 | - | - | - | - | acc | - | - | - | - |
| 2 | - | - | - | - | - | - | - | - | - | - |
| 3 | - | - | s5 | s6 | - | - | - | - | 11 | 4 |
| 4 | r5 | r5 | - | - | r5 | r5 | - | - | - | - |
| 5 | r6 | r6 | - | - | r6 | r6 | - | - | - | - |
| 6 | - | - | s5 | s6 | - | - | - | 12 | 7 | 4 |
| 7 | r3 | s8 | - | - | r3 | r3 | - | - | - | - |
| 8 | - | - | s5 | s6 | - | - | - | - | - | 9 |
| 9 | r4 | r4 | - | - | r4 | r4 | - | - | - | - |
| 10 | r7 | r7 | - | - | r7 | r7 | - | - | - | - |
| 11 | r2 | s8 | - | - | r2 | r2 | - | - | - | - |
| 12 | s3 | - | - | - | s10 | - | - | - | - | - |



Example: A grammar that is not SLR(1)

Consider:

$$\begin{array}{l} S \rightarrow L = R \\ | \\ L \end{array} \quad \begin{array}{l} \rightarrow *R \\ | \\ id \end{array} \quad \begin{array}{l} R \rightarrow L \end{array}$$

Its LR(0) item sets:

$$\begin{array}{lll} I_0 : S' \rightarrow \bullet S \$ & I_5 : L \rightarrow * \bullet R \\ S \rightarrow \bullet L = R & R \rightarrow \bullet L \\ S \rightarrow \bullet R & L \rightarrow \bullet * R \\ L \rightarrow \bullet * R & L \rightarrow \bullet id \\ L \rightarrow \bullet id & I_6 : S \rightarrow L = \bullet R \\ R \rightarrow \bullet L & R \rightarrow \bullet L \\ I_1 : S' \rightarrow S \bullet \$ & L \rightarrow \bullet * R \\ I_2 : S \rightarrow L \bullet = R & L \rightarrow \bullet id \\ R \rightarrow L \bullet & I_7 : L \rightarrow * R \bullet \\ I_3 : S \rightarrow R \bullet & I_8 : R \rightarrow L \bullet \\ I_4 : L \rightarrow id \bullet & I_9 : S \rightarrow L = R \bullet \end{array}$$

Now consider $I_2 : S \rightarrow L \bullet = R \in \text{FOLLOW}(R)$ ($S \Rightarrow L = R \Rightarrow *R = R$)



LR(1) items

Recall: An LR(k) item is a pair $[\alpha, \beta]$, where

- α is a production from G with a \bullet at some position in the RHS, marking how much of the RHS of a production has been seen
- β is a lookahead string containing k symbols (terminals or $\$$)

What about LR(1) items?

- All the lookahead strings are constrained to have length 1
- Look something like $[A \rightarrow X \bullet YZ, a]$



What's the point of the lookahead symbols?

- carry along to choose correct reduction when there is a choice
- lookaheads are bookkeeping, unless item has • at right end:
 - in $[A \rightarrow X \bullet YZ, a]$, a has no direct use
 - in $[A \rightarrow XYZ\bullet, a]$, a is useful
- allows use of grammars that are not *uniquely invertible*[†]

The point: For $[A \rightarrow \alpha\bullet, a]$ and $[B \rightarrow \alpha\bullet, b]$, we can decide between reducing to A or B by looking at limited right context

[†]No two productions have the same RHS



closure1(I)

Given an item $[A \rightarrow \alpha \bullet B\beta, a]$, its closure contains the item and any other items that can generate legal substrings to follow α . Thus, if the parser has viable prefix α on its stack, the input should reduce to $B\beta$ (or γ for some other item $[B \rightarrow \bullet\gamma, b]$ in the closure).

```
function closure1( $I$ )
repeat
    if  $[A \rightarrow \alpha \bullet B\beta, a] \in I$ 
        add  $[B \rightarrow \bullet\gamma, b]$  to  $I$ , where  $b \in \text{FIRST}(\beta a)$ 
until no more items can be added to  $I$ 
return  $I$ 
```



goto1(I)

Let I be a set of LR(1) items and X be a grammar symbol.

Then, $\text{GOTO}(I, X)$ is the closure of the set of all items

$[A \rightarrow \alpha X \bullet \beta, a]$ such that $[A \rightarrow \alpha \bullet X \beta, a] \in I$

If I is the set of valid items for some viable prefix γ , then $\text{GOTO}(I, X)$ is the set of valid items for the viable prefix γX .

$\text{goto}(I, X)$ represents state after recognizing X in state I .

function $\text{goto1}(I, X)$

let J be the set of items $[A \rightarrow \alpha X \bullet \beta, a]$

such that $[A \rightarrow \alpha \bullet X \beta, a] \in I$

return $\text{closure1}(J)$



Building the LR(1) item sets for grammar G

We start the construction with the item $[S' \rightarrow \bullet S, \$]$, where

S' is the start symbol of the augmented grammar G'

S is the start symbol of G

$\$$ represents EOF

To compute the collection of sets of LR(1) items

```
function items( $G'$ )
     $s_0 \leftarrow \text{closure1}(\{[S' \rightarrow \bullet S, \$]\})$ 
     $C \leftarrow \{s_0\}$ 
    repeat
        for each set of items  $s \in C$ 
            for each grammar symbol  $X$ 
                if  $\text{goto1}(s, X) \neq \emptyset$  and  $\text{goto1}(s, X) \notin C$ 
                    add  $\text{goto1}(s, X)$  to  $C$ 
    until no more item sets can be added to  $C$ 
    return  $C$ 
```



Constructing the LR(1) parsing table

Build lookahead into the DFA to begin with

- 1 construct the collection of sets of LR(1) items for G'
- 2 state i of the LR(1) machine is constructed from I_i
 - 1 $[A \rightarrow \alpha \bullet a\beta, b] \in I_i$ and $\text{goto}_1(I_i, a) = I_j$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"shift } j\text{"}$
 - 2 $[A \rightarrow \alpha \bullet, a] \in I_i, A \neq S'$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"reduce } A \rightarrow \alpha\text{"}$
 - 3 $[S' \rightarrow S \bullet, \$] \in I_i$
 $\Rightarrow \text{ACTION}[i, \$] \leftarrow \text{"accept"}$
- 3 $\text{goto}_1(I_i, A) = I_j$
 $\Rightarrow \text{GOTO}[i, A] \leftarrow j$
- 4 set undefined entries in ACTION and GOTO to "error"
- 5 initial state of parser s_0 is $\text{closure}_1([S' \rightarrow \bullet S, \$])$



Back to previous example (\notin SLR(1))

| | | |
|---|---|---|
| $S \rightarrow L = R$ | $I_0 : S' \rightarrow \bullet S, \$$ | $I_5 : L \rightarrow \text{id}\bullet, \$$ |
| R | $S \rightarrow \bullet L = R, \$$ | $I_6 : S \rightarrow L = \bullet R, \$$ |
| $L \rightarrow *R$ | $S \rightarrow \bullet R, \$$ | $R \rightarrow \bullet L, \$$ |
| id | $L \rightarrow \bullet *R, =$ | $L \rightarrow \bullet *R, \$$ |
| $R \rightarrow L$ | $L \rightarrow \bullet \text{id}, =$ | $L \rightarrow \bullet \text{id}, \$$ |
| | $R \rightarrow \bullet L, \$$ | $I_7 : L \rightarrow *R\bullet, \$$ |
| | $L \rightarrow \bullet *R, \$$ | $I_8 : R \rightarrow L\bullet, \$$ |
| | $L \rightarrow \bullet \text{id}, \$$ | $I_9 : S \rightarrow L = R\bullet, \$$ |
| $I_1 : S' \rightarrow S\bullet, \$$ | | $I_{10} : R \rightarrow L\bullet, \$$ |
| $I_2 : S \rightarrow L\bullet = R, \$$ | | $I_{11} : L \rightarrow * \bullet R, \$$ |
| | $R \rightarrow L\bullet, \$$ | $R \rightarrow \bullet L, \$$ |
| $I_3 : S \rightarrow R\bullet, \$$ | | $L \rightarrow \bullet *R, \$$ |
| $I_4 : L \rightarrow * \bullet R, = \$$ | | $L \rightarrow \bullet \text{id}, \$$ |
| | $R \rightarrow \bullet L, = \$$ | $I_{12} : L \rightarrow \text{id}\bullet, \$$ |
| | $L \rightarrow \bullet *R, = \$$ | $I_{13} : L \rightarrow *R\bullet, \$$ |
| | $L \rightarrow \bullet \text{id}, = \$$ | |

I_2 no longer has shift-reduce conflict: reduce on $\$$, shift on $=$



Example: back to SLR(1) expression grammar

In general, LR(1) has many more states than LR(0)/SLR(1):

| | | | |
|---|-----------------------|---|---------------------------|
| 1 | $S \rightarrow E$ | 4 | $T \rightarrow T * F$ |
| 2 | $E \rightarrow E + T$ | 5 | $ F$ |
| 3 | $ T$ | 6 | $F \rightarrow \text{id}$ |

7 | (E)

LR(1) item sets:

I_0 :

$$\begin{aligned} S &\rightarrow \bullet E, \$ \\ E &\rightarrow \bullet E + T, +\$ \\ E &\rightarrow \bullet T, +\$ \\ T &\rightarrow \bullet T * F, *+ \$ \\ T &\rightarrow \bullet F, *+ \$ \\ F &\rightarrow \bullet \text{id}, *+ \$ \\ F &\rightarrow \bullet (E), *+ \$ \end{aligned}$$

I'_0 : shifting (

$$\begin{aligned} F &\rightarrow (\bullet E), *+ \$ \\ E &\rightarrow \bullet E + T, +) \\ E &\rightarrow \bullet T, +) \\ T &\rightarrow \bullet T * F, *+) \\ T &\rightarrow \bullet F, *+) \\ F &\rightarrow \bullet \text{id}, *+) \\ F &\rightarrow \bullet (E), *+) \end{aligned}$$

I''_0 : shifting (

$$\begin{aligned} F &\rightarrow (\bullet E), *+) \\ E &\rightarrow \bullet E + T, +) \\ E &\rightarrow \bullet T, +) \\ T &\rightarrow \bullet T * F, *+) \\ T &\rightarrow \bullet F, *+) \\ F &\rightarrow \bullet \text{id}, *+) \\ F &\rightarrow \bullet (E), *+) \end{aligned}$$



Another example

Consider:

| | |
|---|--------------------|
| 0 | $S' \rightarrow S$ |
| 1 | $S \rightarrow CC$ |
| 2 | $C \rightarrow cC$ |
| 3 | d |

| state | ACTION | GOTO | |
|-------|---------|------|--|
| | c d \$ | S C | |
| 0 | s3 s4 - | 1 2 | |
| 1 | - - acc | - - | |
| 2 | s6 s7 - | - 5 | |
| 3 | s3 s4 - | - 8 | |
| 4 | r3 r3 - | - - | |
| 5 | - - r1 | - - | |
| 6 | s6 s7 - | - 9 | |
| 7 | - - r3 | - - | |
| 8 | r2 r2 - | - - | |
| 9 | - - r2 | - - | |

LR(1) item sets:

$$I_0 : S' \rightarrow \bullet S, \$$$

$$S \rightarrow \bullet CC, \$$$

$$C \rightarrow \bullet cC, cd$$

$$C \rightarrow \bullet d, cd$$

$$I_1 : S' \rightarrow S \bullet, \$$$

$$I_2 : S \rightarrow C \bullet C, \$$$

$$C \rightarrow \bullet cC, \$$$

$$C \rightarrow \bullet d, \$$$

$$I_3 : C \rightarrow c \bullet C, cd$$

$$C \rightarrow \bullet cC, cd$$

$$C \rightarrow \bullet d, cd$$

$$I_4 : C \rightarrow d \bullet, cd$$

$$I_5 : S \rightarrow CC \bullet, \$$$

$$I_6 : C \rightarrow c \bullet C, \$$$

$$C \rightarrow \bullet cC, \$$$

$$C \rightarrow \bullet d, \$$$

$$I_7 : C \rightarrow d \bullet, \$$$

$$I_8 : C \rightarrow cC \bullet, cd$$

$$I_9 : C \rightarrow cC \bullet, \$$$



LALR(1) parsing

Define the *core* of a set of LR(1) items to be the set of LR(0) items derived by ignoring the lookahead symbols.

Thus, the two sets

- $\{[A \rightarrow \alpha \bullet \beta, a], [A \rightarrow \alpha \bullet \beta, b]\}$, and
- $\{[A \rightarrow \alpha \bullet \beta, c], [A \rightarrow \alpha \bullet \beta, d]\}$

have the same core.

Key idea:

If two sets of LR(1) items, I_i and I_j , have the same core, we can merge the states that represent them in the ACTION and GOTO tables.



LALR(1) table construction

To construct LALR(1) parsing tables, we can insert a single step into the LR(1) algorithm

- (1.5) For each core present among the set of LR(1) items, find all sets having that core and replace these sets by their union.

The goto function must be updated to reflect the replacement sets.

The resulting algorithm has large space requirements, as we still are required to build the full set of LR(1) items.



LALR(1) table construction

The revised (*and renumbered*) algorithm

- ① construct the collection of sets of LR(1) items for G'
- ② for each core present among the set of LR(1) items, find all sets having that core and replace these sets by their union (update the `goto1` function incrementally)
- ③ state i of the LALR(1) machine is constructed from I_i .
 - ① $[A \rightarrow \alpha \bullet a\beta, b] \in I_i$ and $\text{goto1}(I_i, a) = I_j$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"shift } j\text{"}$
 - ② $[A \rightarrow \alpha \bullet, a] \in I_i, A \neq S'$
 $\Rightarrow \text{ACTION}[i, a] \leftarrow \text{"reduce } A \rightarrow \alpha\text{"}$
 - ③ $[S' \rightarrow S \bullet, \$] \in I_i \Rightarrow \text{ACTION}[i, \$] \leftarrow \text{"accept"}$
- ④ $\text{goto1}(I_i, A) = I_j \Rightarrow \text{GOTO}[i, A] \leftarrow j$
- ⑤ set undefined entries in ACTION and GOTO to "error"
- ⑥ initial state of parser s_0 is $\text{closure1}([S' \rightarrow \bullet S, \$])$



Example

Reconsider:

| | |
|---|-----------------------------|
| | $S' \rightarrow S$ |
| 0 | $S \rightarrow CC$ |
| 1 | $CC \rightarrow \bullet cC$ |
| 2 | $C \rightarrow cC$ |
| 3 | $ $ d |

$$\begin{array}{ll}
 I_0 : S' \rightarrow \bullet S, \$ & \\
 S \rightarrow \bullet CC, \$ & \\
 C \rightarrow \bullet cC, cd & \\
 C \rightarrow \bullet d, cd & \\
 I_1 : S' \rightarrow S\bullet, \$ & \\
 I_2 : S \rightarrow C\bullet C, \$ & \\
 C \rightarrow \bullet cC, \$ & \\
 C \rightarrow \bullet d, \$ &
 \end{array}$$

$$\begin{array}{ll}
 I_3 : C \rightarrow c\bullet C, cd & I_6 : C \rightarrow c\bullet C, \$ \\
 C \rightarrow \bullet cC, cd & C \rightarrow \bullet cC, \$ \\
 C \rightarrow \bullet d, cd & C \rightarrow \bullet d, \$ \\
 I_4 : C \rightarrow d\bullet, cd & I_7 : C \rightarrow d\bullet, \$ \\
 I_5 : S \rightarrow CC\bullet, \$ & I_8 : C \rightarrow cC\bullet, cd \\
 & I_9 : C \rightarrow cC\bullet, \$ \\
 &
 \end{array}$$

Merged states:

$$\begin{array}{l}
 I_{36} : C \rightarrow c\bullet C, cd\$ \\
 C \rightarrow \bullet cC, cd\$ \\
 C \rightarrow \bullet d, cd\$ \\
 I_{47} : C \rightarrow d\bullet, cd\$ \\
 I_{89} : C \rightarrow cC\bullet, cd\$
 \end{array}$$

| state | ACTION | | | GOTO | |
|-------|--------|-----|-----|------|---|
| | c | d | \$ | S | C |
| 0 | s36 | s47 | - | 1 | 2 |
| 1 | - | - | acc | - | - |
| 2 | s36 | s47 | - | - | 5 |
| 36 | s36 | s47 | - | - | 8 |
| 47 | r3 | r3 | r3 | - | - |
| 5 | - | - | r1 | - | - |
| 89 | r2 | r2 | r2 | - | - |



More efficient LALR(1) construction

Observe that we can:

- represent I_i by its *basis* or *kernel*:
items that are either $[S' \rightarrow \bullet S, \$]$
or do not have \bullet at the left of the RHS
- compute *shift*, *reduce* and *goto* actions for state derived from I_i directly from its kernel

This leads to a method that avoids building the complete canonical collection of sets of LR(1) items

Self reading: Section 4.7.5 Dragon book



The role of precedence

Precedence and associativity can be used to resolve shift/reduce conflicts in ambiguous grammars.

- lookahead with higher precedence \Rightarrow *shift*
- same precedence, left associative \Rightarrow *reduce*

Advantages:

- more concise, albeit ambiguous, grammars
- shallower parse trees \Rightarrow fewer reductions

Classic application: expression grammars



The role of precedence

With precedence and associativity, we can use:

$$\begin{array}{lcl} E & \rightarrow & E * E \\ & | & \\ & E / E & \\ & | & \\ & E + E & \\ & | & \\ & E - E & \\ & | & \\ & (E) & \\ & | & \\ & -E & \\ & | & \\ & \text{id} & \\ & | & \\ & \text{num} & \end{array}$$

This eliminates useless reductions (*single productions*)



Error recovery in shift-reduce parsers

The problem

- encounter an invalid token
- bad pieces of tree hanging from stack
- incorrect entries in symbol table

We want to *parse* the rest of the file

Restarting the parser

- find a restartable state on the stack
- move to a consistent place in the input
- print an informative message to `stderr`

(line number)



Error recovery in yacc/bison/Java CUP

The error mechanism

- designated token `error`
- valid in any production
- `error` shows synchronization points

When an error is discovered

- pops the stack until `error` is legal
- skips input tokens until it successfully shifts 3 (some default value)
- `error` productions can have actions

This mechanism is fairly general

Read the section on Error Recovery of the on-line CUP manual



Example

Using error

```
stmt_list : stmt
           | stmt_list ; stmt
```

can be augmented with error

```
stmt_list : stmt
           | error
           | stmt_list ; stmt
```

This should

- throw out the erroneous statement
- synchronize at ";" or "end"
- invoke yyerror ("syntax error")

Other "natural" places for errors

- all the "lists": FieldList, CaseList
- missing parentheses or brackets
- extra operator or missing operator

(yychar)



Left versus right recursion

Right Recursion:

- needed for termination in predictive parsers
- requires more stack space
- right associative operators

Left Recursion:

- works fine in bottom-up parsers
- limits required stack space
- left associative operators

Rule of thumb:

- right recursion for top-down parsers
- left recursion for bottom-up parsers

Left recursive grammar:

$$\begin{aligned} E &\rightarrow E + T | E \\ T &\rightarrow T * F | F \\ F &\rightarrow (E) + Int \end{aligned}$$

After left recursion removal

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' | \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' | \epsilon \\ F &\rightarrow (E) + Int \end{aligned}$$

Parse the string $3 + 4 + 5$



- *Recursive descent*

A hand coded recursive descent parser directly encodes a grammar (typically an LL(1) grammar) into a series of mutually recursive procedures. It has most of the linguistic limitations of LL(1).

- $\text{LL}(k)$

An $\text{LL}(k)$ parser must be able to recognize the use of a production after seeing only the first k symbols of its right hand side.

- $\text{LR}(k)$

An $\text{LR}(k)$ parser must be able to recognize the occurrence of the right hand side of a production after having seen all that is derived from that right hand side with k symbols of lookahead.



Grammar hierarchy

- $\text{LR}(k) > \text{LR}(1) > \text{LALR}(1) > \text{SLR}(1) > \text{LR}(0)$
- $\text{LL}(k) > \text{LL}(1) > \text{LL}(0)$
- $\text{LR}(0) > \text{LL}(0)$
- $\text{LR}(1) > \text{LL}(1)$
- $\text{LR}(k) > \text{LL}(k)$



