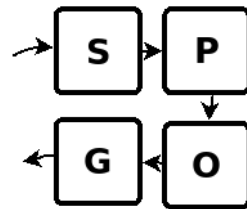


Parsing and Context Free Grammars (Chapter 4)

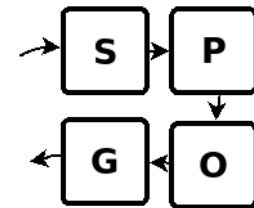


Recall: now we have a scanner that does lexical analysis for us, giving us a token sequence rather than a character sequence.

Now, we need to match the token sequence to our programming language's grammar.

Or, we need to **parse** it.

As we parse it, we create data that represents our program, and then use that in the back end to generate our output language (e.g., object code or JVM bytecode or assembly code or...)

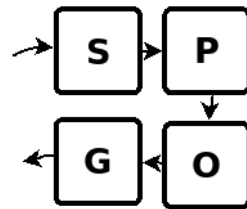


Recall: for top-down recursive descent parsing, we used a **context-free grammar** (CFG) to specify our input language.

Some form of CFGs are **the canonical way** to specify programming languages...

Even though the CFG cannot capture everything that specifies a compile-able program!

E.g.: A CFG cannot encode that a variable must be declared before it is used.



Formalities of Context Free Grammars (textbook sec 4.2.1)

A CFG is a 4-tuple $\langle T, N, s, R \rangle$ where:

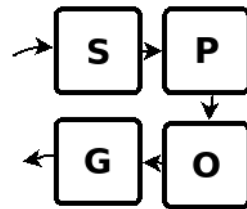
T is the set of terminals (token types)

N is the set of non-terminals

$s \in N$ is the starting non-terminal

R is the set of production rules

$r \in R, r = n \rightarrow v, n \in N, v \in (T \cup N)^*$



CFG:

$S \rightarrow A B A$

$A \rightarrow a$

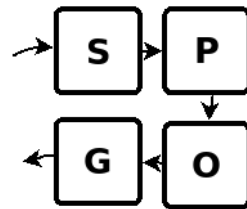
$B \rightarrow \epsilon$

$B \rightarrow b B$

This grammar accepts strings like aba, abba, ...

IOW, its language is the same as the regular expression ab^*a

Indeed, CFGs can encode any language that REs can...and more!



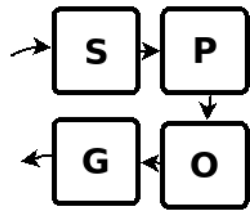
REs define *regular languages*

CFGs define *context-free languages*

So, all regular languages *are* context-free languages!

Does this mean there is a language hierarchy?

YES!

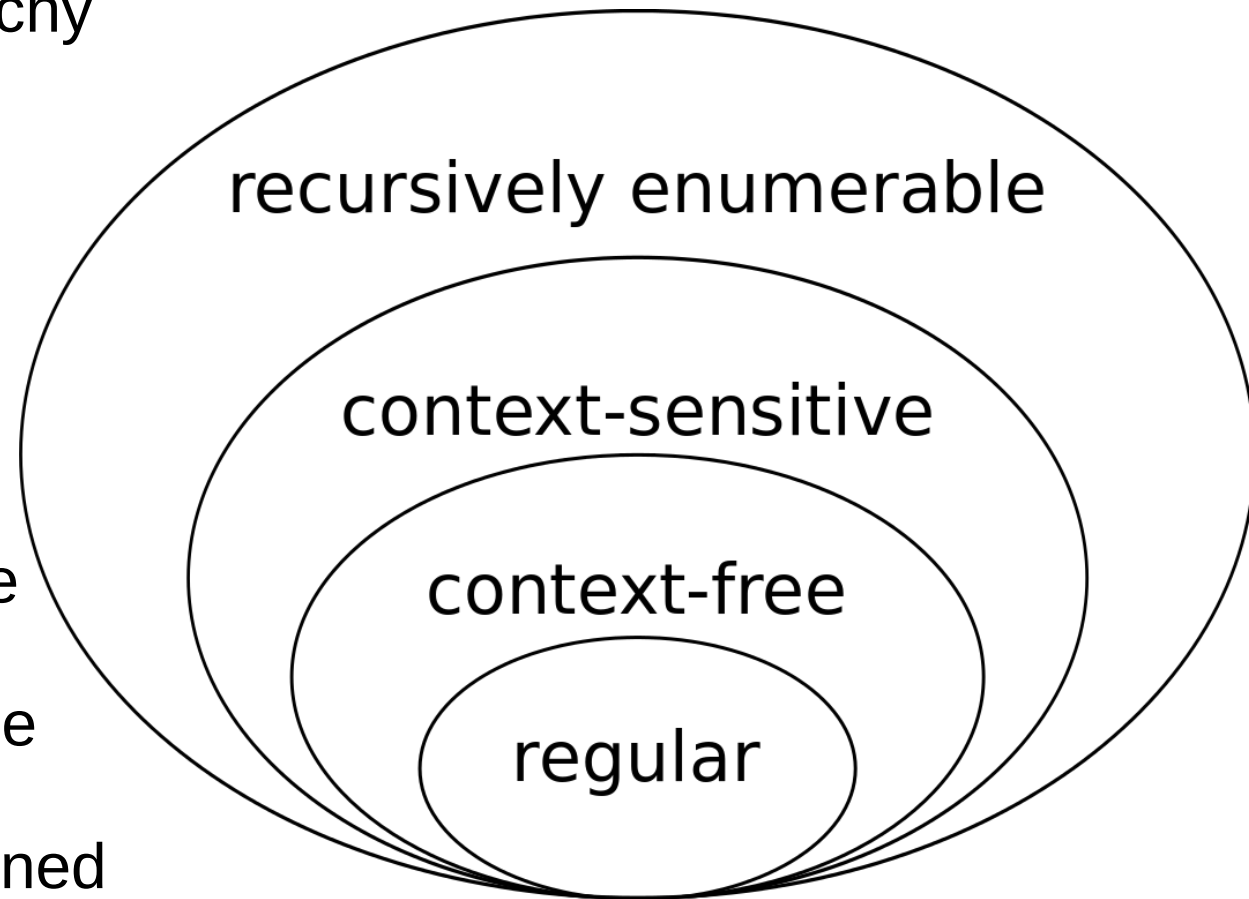


Chomsky Hierarchy of Languages

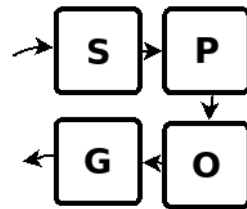
Wikipedia:Chomsky_hierarchy

Notes

- CF: the nonterminal head of the production rule has no context
- CS: head allows more than just the nonterminal
- RecEnum: think “can write a no-input function to enumerate all strings in the language”
- These four have been refined over the years



<https://en.wikipedia.org/wiki/File:Chomsky-hierarchy.svg>

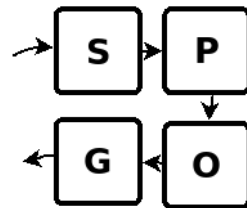


Textbook Symbol Usage

Uppercase letters are non-terminals: A,B,...

Lowercase letters are terminals: a,b,...

Lowercase Greek letters are sequences of terminals and nonterminals: α , β ,

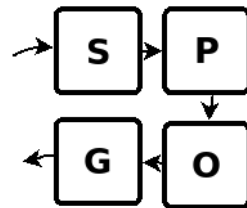


Example Expression Grammar in textbook p198

$\text{Expr} \rightarrow \text{Expr} \text{ '+' } \text{Term}$
 $\text{Expr} \rightarrow \text{Expr} \text{ '-' } \text{Term}$
 $\text{Expr} \rightarrow \text{Term}$
 $\text{Term} \rightarrow \text{Term} \text{ '*' } \text{Factor}$
 $\text{Term} \rightarrow \text{Term} \text{ '/' } \text{Factor}$
 $\text{Term} \rightarrow \text{Factor}$
 $\text{Factor} \rightarrow \text{'(' Expr ')'}$
 $\text{Factor} \rightarrow \text{id}$
 $\text{Factor} \rightarrow \text{number}$

or

$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow (E) \mid \text{id} \mid \text{number}$



$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow (E) \mid \text{id} \mid \text{number}$

string: "i*42-j" (i,j are id's, 42 is a number)

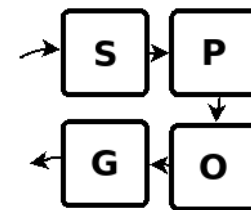
derivation:

$E \Rightarrow E-T \Rightarrow T-T \Rightarrow T * F-T \Rightarrow F * F-T \Rightarrow \text{id} * F-T \Rightarrow$
 $\text{id} * \text{number}-T \Rightarrow \text{id} * \text{number}-F \Rightarrow \text{id} * \text{number}-\text{id}$

Each derivation step gives a **sentential form** of G

If $S \Rightarrow^* \alpha$, where S is the starting nonterminal and α is a sequence of nonterminals and terminals, then α is a sentential form

A sentential form without any nonterminals is a **sentence** of G, or a sentence in $L(G)$



$$E \rightarrow E+T \mid E-T \mid T$$

$$T \rightarrow T * F \mid T / F \mid F$$

$$F \rightarrow (E) \mid \text{id} \mid \text{number}$$

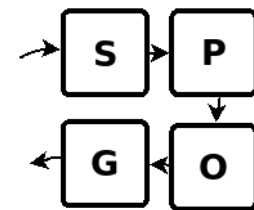
string: “i*42-j” (i,j are id’s, 42 is a number)

$$E \Rightarrow E-T \Rightarrow T-T \Rightarrow T * F-T \Rightarrow F * F-T \Rightarrow \text{id} * F-T \Rightarrow \text{id} * \text{number}-T \Rightarrow \text{id} * \text{number}-F \Rightarrow \text{id} * \text{number}-\text{id}$$

leftmost derivation: the leftmost nonterminal in each sentential form is always chosen for the next derivation; the above is a leftmost derivation; each step is a **left sentential form**

rightmost derivation: rightmost nonterminal is chosen; see below

$$E \Rightarrow E-T \Rightarrow E-F \Rightarrow E-\text{id} \Rightarrow T-\text{id} \Rightarrow T * F-\text{id} \Rightarrow T * \text{number}-\text{id} \Rightarrow F * \text{number}-\text{id} \Rightarrow \text{id} * \text{number}-\text{id}$$

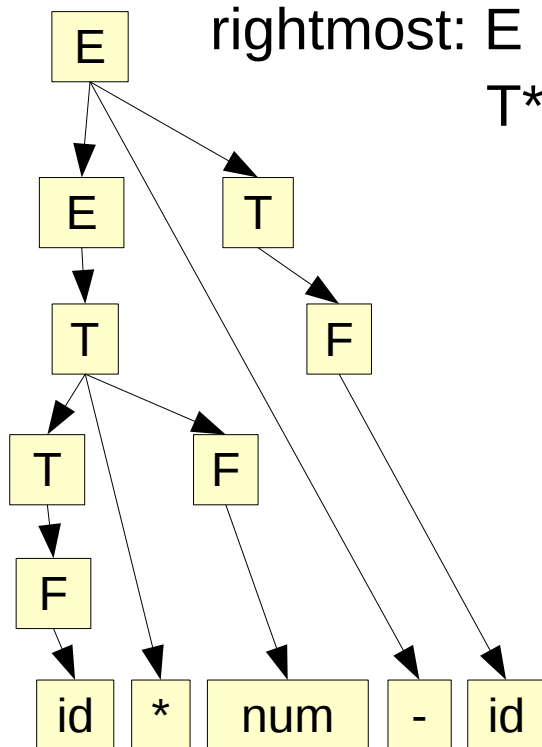


$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T*F \mid T/F \mid F$
 $F \rightarrow (E) \mid \text{id} \mid \text{number}$

string: "i*42-j" (i,j are id's, 42 is a number)

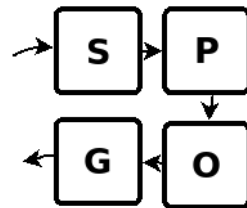
leftmost: $E \Rightarrow E-T \Rightarrow T-T \Rightarrow T*F-T \Rightarrow F*F-T \Rightarrow \text{id}*F-T \Rightarrow$
 $\text{id*number-T} \Rightarrow \text{id*number-F} \Rightarrow \text{id*number-id}$

rightmost: $E \Rightarrow E-T \Rightarrow E-F \Rightarrow E-\text{id} \Rightarrow T-\text{id} \Rightarrow T*F-\text{id} \Rightarrow$
 $T*\text{number-id} \Rightarrow F*\text{number-id} \Rightarrow \text{id*number-id}$



Notes:

- both derivations produce the same parse tree!
- this is because this grammar is unambiguous!



Eliminating *immediate* left recursion

with production rules like

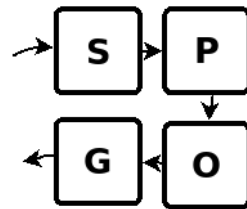
$$A \rightarrow Aa \mid Ab \mid Ac \mid q \mid r \mid s$$

then rewrite with new nonterminal B like

$$A \rightarrow qB \mid rB \mid sB$$
$$B \rightarrow aB \mid bB \mid cB \mid \text{empty}$$

(a-c) and (q-s) are *examples*; these could be *sequences* of terminals and nonterminals

bottom of page 212 of textbook



Eliminating *ALL* left recursion

rules might have left recursion that refers to other nonterminals, not just itself; this is non-immediate left recursion. How to remove? Algorithm 4.19 p 213 of tbook

Informally, given recursion of the form

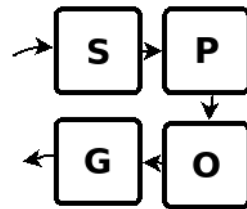
$A \rightarrow Bq$ (r, s, q are sequences of terms and non-terms)
 $B \rightarrow r \mid s$ (r and s have A somewhere in them)

Replace rules for A by expanding B rules into them, as:

$A \rightarrow rq \mid sq$

And then eliminate any immediate left recursion that might exist for A .

Repeat for all production rules (real alg imposes ordering and applies above only to previously-seen nonterminals)



Eliminating *ALL* left recursion

Example grammar:

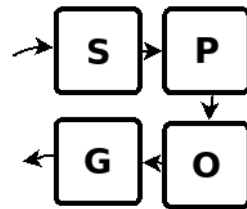
$$S \rightarrow A a \mid b$$
$$A \rightarrow A c \mid S d \mid \varepsilon$$

order: S,A; nothing to do for S (no previous nt's); at A, we replace S with its productions, so that we get

$$S \rightarrow A a \mid b$$
$$A \rightarrow A c \mid A a d \mid b d \mid \varepsilon$$

Now we eliminate the immediate left recursion to get

$$S \rightarrow A a \mid b$$
$$A \rightarrow b d B \mid B$$
$$B \rightarrow c B \mid a d B \mid \varepsilon$$



Left Factoring

Left factoring is a simple transformation that just removes equal production rule prefixes from the rules so that top down parsing can handle them

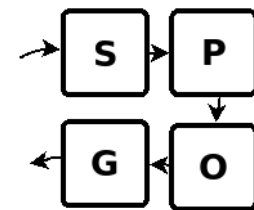
Example grammar:

$$S \rightarrow A b \mid A c$$
$$A \rightarrow a A \mid b A \mid \epsilon$$

Rules for S both begin with A, so cannot immediately choose. Left factoring just removes the suffixes into another nonterm:

$$S \rightarrow A B$$
$$B = b \mid c$$
$$A \rightarrow a A \mid b A \mid \epsilon$$

Now only one rule for S that begins with prefix A. In general, any prefix (string of terminals and nonterminals) can be factored.

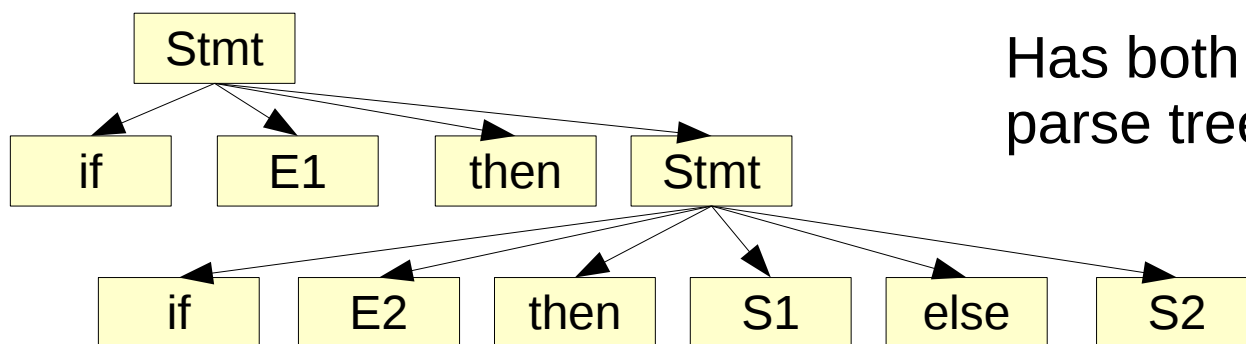


Ambiguity: if-then-else example

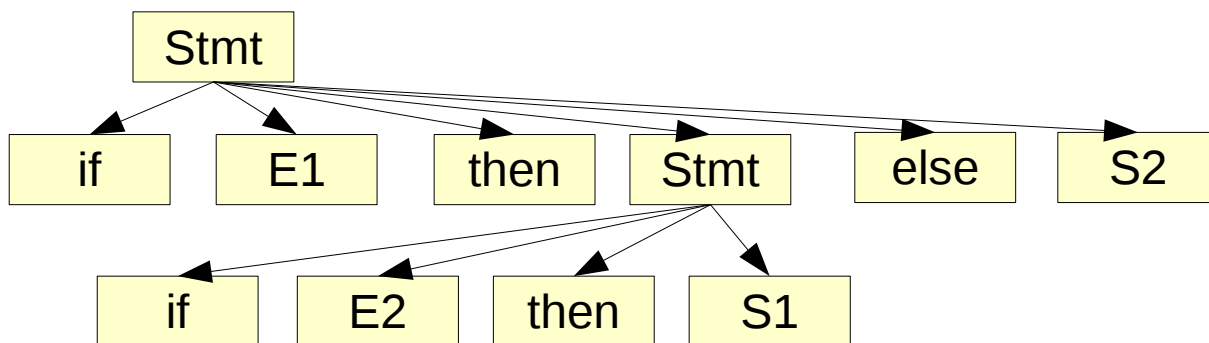
Example grammar:

Stmt \rightarrow 'if' Expr 'then' Stmt
| 'if' Expr 'then' Stmt 'else' Stmt
| // other statements

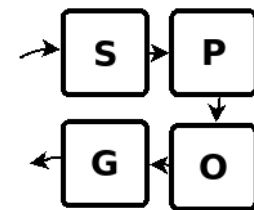
Input token sequence: if E1 then if E2 then S1 else S2



Has both of these possible
parse trees! We want the first!



Real C code:
if (x>2)
 if (x>5)
 x = 11;
else
 x = 42;

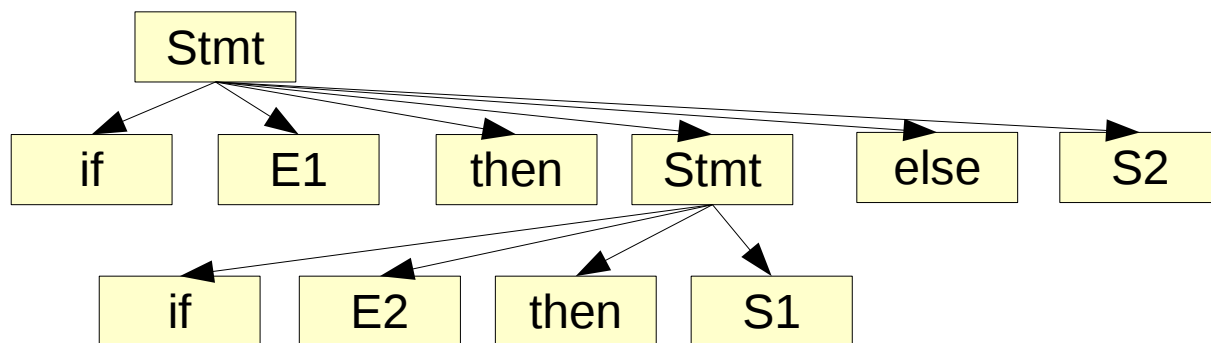
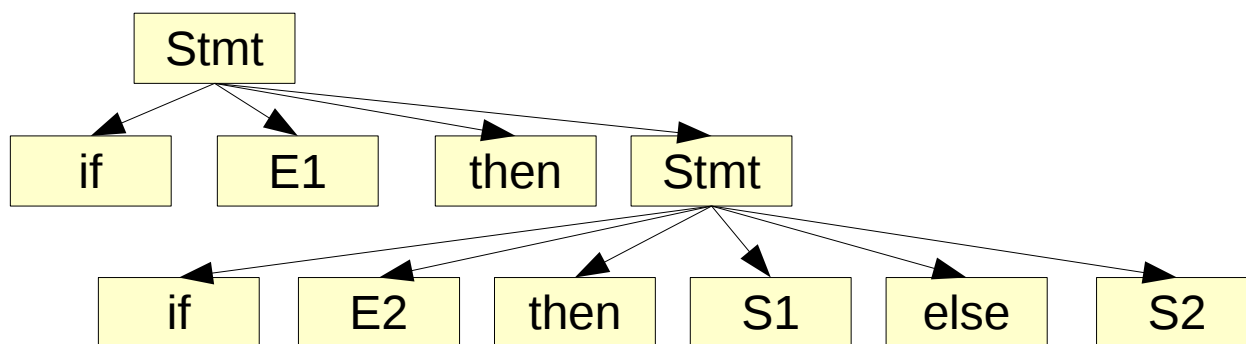


Ambiguity: if-then-else example

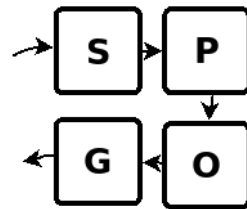
Example grammar:

Stmt \rightarrow 'if' Expr 'then' Stmt
 | 'if' Expr 'then' Stmt 'else' Stmt
 | <other statements>

“if E1 then if E2 then S1 else S2”



Textbook, p211: “In all programming languages with conditional statements of this form, the first parse tree is preferred. The general rule is, “Match each *else* with the closest unmatched *if*. This disambiguation rule can theoretically be incorporated directly into the grammar, but in practice it is rarely built into productions.”



Top Down Parsing (TB Sec 4.4)

We looked at simple top-down parsing in Chapter 2, using recursive descent parsing

- and saw that left recursion can be a problem

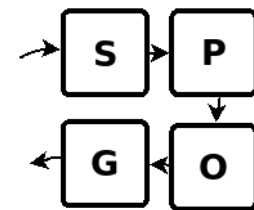
But we didn't think too hard about what works and what doesn't for top-down parsing

Example grammar:

$$S \rightarrow A S \mid B S$$
$$A \rightarrow a B$$
$$B \rightarrow b$$

In top-down, we start at S, but with no terminal symbol in any production rule, we don't have an immediate choice

But, if we look ahead at the waiting token, it will clearly tell us whether to recurse into `nontermA()` or `nontermB()`.



Top Down Parsing (TB Sec 4.4)

We can determine if a grammar is suitable for top-down.

Helper definitions:

$\text{FIRST}(\alpha)$ = the set of terminals that begin all strings derivable from α , where α is a sequence of grammar symbols (terminals and nonterminals).

$\text{FOLLOW}(A)$ = the set of terminals that can immediately follow the nonterminal A in any sentential form

Example grammar:

$S \rightarrow A S \mid B S \mid \epsilon$

$A \rightarrow a B$

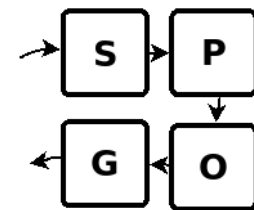
$B \rightarrow b$

$\text{FIRST}(A S) = \{a\}$, $\text{FIRST}(B S) = \{b\}$, $\text{FIRST}(\epsilon) = \{\epsilon\}$

$\text{FOLLOW}(S) = \{\$ \}$

$\text{FIRST}(a B) = \{a\}$, $\text{FOLLOW}(A) = \{a, b, \$ \}$

$\text{FIRST}(b) = \{b\}$, $\text{FOLLOW}(B) = \{a, b, \$ \}$



LL(1) Grammars (TB Sec 4.4)

If a grammar is parse-able top-down (no backtracking) by only looking at the next input token, then it is an **LL(1) grammar**:

L == scanning input Left-to-right

L == producing a Leftmost derivation

1 == looking ahead only 1 input token

All LL(1) grammars are context-free, but not all CFGs are LL(1)

Example grammar:

$S \rightarrow A S \mid B S \mid \epsilon$

$A \rightarrow a B$

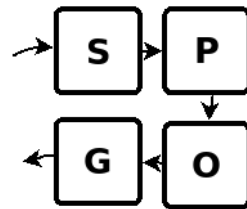
$B \rightarrow b$

$\text{FIRST}(A S) = \{a\}$, $\text{FIRST}(B S) = \{b\}$, $\text{FIRST}(\epsilon) = \{\epsilon\}$

$\text{FOLLOW}(S) = \{\$ \}$

$\text{FIRST}(a B) = \{a\}$, $\text{FOLLOW}(A) = \{a, b, \$ \}$

$\text{FIRST}(b) = \{b\}$, $\text{FOLLOW}(B) = \{a, b, \$ \}$



LL(1) Grammars (TB Sec 4.4)

Requirements (p224):

- $\text{FIRST}(\alpha)$ is disjoint for all production rules of a nonterminal S
- if one production rule of a nonterminal S produces ϵ , then $\text{FIRST}()$'s of all S rules are disjoint from $\text{FOLLOW}(S)$

All LL(1) grammars are context-free; not all CFGs are LL(1)

Example grammar:

Yes, this grammar is LL(1)

$S \rightarrow A S \mid B S \mid \epsilon$

$A \rightarrow a B$

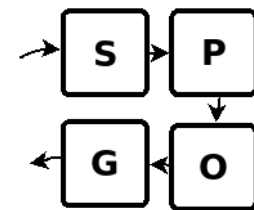
$B \rightarrow b$

$\text{FIRST}(A S) = \{a\}$, $\text{FIRST}(B S) = \{b\}$, $\text{FIRST}(\epsilon) = \{\epsilon\}$

$\text{FOLLOW}(S) = \{\$ \}$

$\text{FIRST}(a B) = \{a\}$, $\text{FOLLOW}(A) = \{a, b, \$ \}$

$\text{FIRST}(b) = \{b\}$, $\text{FOLLOW}(B) = \{a, b, \$ \}$



LL(1) Grammars (TB Sec 4.4)

Top-down parsing does not **have** to be done recursively!

Algorithm 4.31 (p224) builds a parsing table for an LL(1) grammar, and parsing can be table-driven (Alg 4.34)

$$S \rightarrow A S \mid B S \mid \epsilon$$

$$A \rightarrow a B$$

$$B \rightarrow b$$

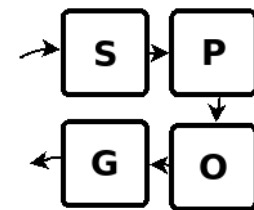
$$\text{FIRST}(A S) = \{a\}, \text{FIRST}(B S) = \{b\}, \text{FIRST}(\epsilon) = \{\epsilon\}$$

$$\text{FOLLOW}(S) = \{\$ \}$$

$$\text{FIRST}(a B) = \{a\}, \text{FOLLOW}(A) = \{a, b, \$ \}$$

$$\text{FIRST}(b) = \{b\}, \text{FOLLOW}(B) = \{a, b, \$ \}$$

NT\LA	a	b	\$
S	$S \rightarrow A S$	$S \rightarrow B S$	$S \rightarrow \epsilon$
A	$A \rightarrow a B$		
B		$B \rightarrow b$	



LL(1) Grammars (TB Sec 4.4)

$S \rightarrow A S \mid B S \mid \epsilon$

$A \rightarrow a B$

$B \rightarrow b$

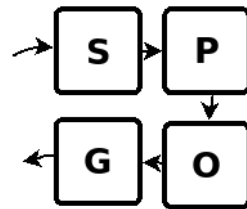
NT\LA	a	b	\$
S	$S \rightarrow A S$	$S \rightarrow B S$	$S \rightarrow \epsilon$
A	$A \rightarrow a B$		
B		$B \rightarrow b$	

Algorithm --

Repeat until match \$:
 if top is terminal, match
 and consume, or err
 if top is nonterminal,
 use lookahead to
 select table rule,
 replace on stack with
 production rule RHS

Parsing string "ab":

1. stack = [S \$], lookahead = 'a' (replace S with A S on stack)
2. stack = [A S \$], lookahead = 'a' (replace A with a B)
3. stack = [a B S \$], lookahead = 'a' (now match & consume)
4. stack = [B S \$], lookahead = 'b' (replace B with b)
5. stack = [b S \$], lookahead = 'b' (now match & consume)
6. stack = [S \$], lookahead = '\$' (replace S with (empty))
7. stack = [\$], lookahead = '\$' (match and finish)



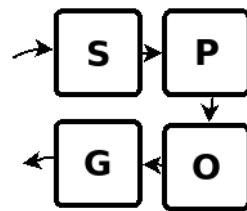
Bottom Up Parsing (TB Sec 4.5)

LL(1) grammars can be parsed with top down parsing with one-token lookahead. Great! Note that **this is really used** in the **ANTLR** parser generator tool (which is Java based). ANTLR actually uses “LL(*)” parsing.

However, **yacc** does bottom up parsing. So we are going to spend time learning bottom up parsing, in more detail than we did top down parsing!

TB p234: “We can think of bottom up parsing as the process of ‘reducing’ a string w to the start symbol of the grammar. At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of that production.”

Some of you may have seen a “shift/reduce” conflict warning in your yacc grammars; the “reduce” word in this warning is the “reduction step” in the textbook quote above.



Bottom Up Parsing (TB Sec 4.5)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

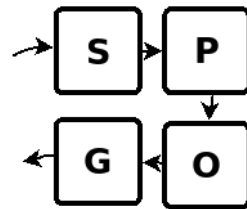
string: "i*42-j" (i,j are id's, 42 is a number)

bottom-up: $\text{id}*\text{number-id} \Rightarrow F*\text{number-id} \Rightarrow T*\text{number-id} \Rightarrow T*F\text{-id} \Rightarrow T\text{-id} \Rightarrow E\text{-id} \Rightarrow E\text{-F} \Rightarrow E\text{-T} \Rightarrow E$

If you reverse the sequence above, you get the rightmost derivation! (leftmost and rightmost shown below from earlier)

leftmost: $E \Rightarrow E\text{-T} \Rightarrow T\text{-T} \Rightarrow T*F\text{-T} \Rightarrow F*F\text{-T} \Rightarrow \text{id}*F\text{-T} \Rightarrow \text{id}*\text{number-T} \Rightarrow \text{id}*\text{number-F} \Rightarrow \text{id}*\text{number-id}$

rightmost: $E \Rightarrow E\text{-T} \Rightarrow E\text{-F} \Rightarrow E\text{-id} \Rightarrow T\text{-id} \Rightarrow T*F\text{-id} \Rightarrow T*\text{number-id} \Rightarrow F*\text{number-id} \Rightarrow \text{id}*\text{number-id}$



Bottom Up Parsing (TB Sec 4.5)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T * F \mid T / F \mid F$

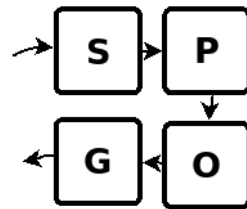
$F \rightarrow (E) \mid \text{id} \mid \text{number}$

string: "i*42-j" (i,j are id's, 42 is a number)

bottom-up: $\text{id} * \text{number} - \text{id} \Rightarrow F * \text{number} - \text{id} \Rightarrow T * \text{number} - \text{id} \Rightarrow$
 $T * F - \text{id} \Rightarrow T - \text{id} \Rightarrow E - \text{id} \Rightarrow E - F \Rightarrow E - T \Rightarrow E$

Note:

- scanning of input still proceeds Left to right
- but we produced a Rightmost derivation
- so we can call this **LR** parsing, and just like LL(k) grammars there are LR(k) grammars



Bottom Up Parsing (TB Sec 4.5)

How to do bottom up parsing? First, definitions

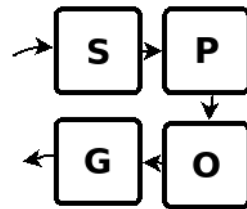
Handle: “a substring that matches the body of a production, and whose reduction represents one step along the reverse of a rightmost derivation.” (TB p235)

If $S \rightarrow^* \alpha A \omega \rightarrow \alpha \beta \omega$ in a rightmost derivation, then β is a handle in position following α (TB p235) (ω must have only terminals)

Note: Handles are contextual, not generic (see slide bottom)

Note: Not all leftmost strings that match a production rule must be handles – only those that would be in a rightmost derivation!

“If a grammar is unambiguous, then every right sentential form of the grammar has exactly one handle” (TB, p235)



Bottom Up Parsing (TB Sec 4.5)

How to do bottom up parsing? **Shift-Reduce** parsing

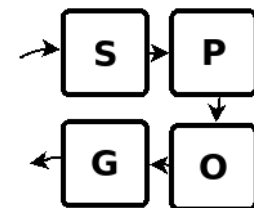
Again (as in LL parsing), use a stack.

Shift (push) input symbols (tokens) onto stack until the top of the stack contains a **handle**.

Reduce the handle to its production rule head.

Go back to shifting input.

If at end of input stack contains only the starting nonterminal, accept; else report error



Bottom Up Parsing (TB Sec 4.5)

$E \rightarrow E+T \mid E-T \mid T$ string: "i*42-j"

$T \rightarrow T*F \mid T/F \mid F$

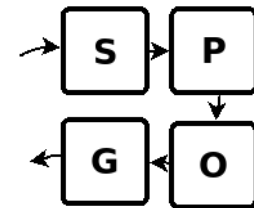
$F \rightarrow (E) \mid \text{id} \mid \text{number}$

bottom-up: $\text{id} * \text{number} - \text{id} \Rightarrow F * \text{number} - \text{id} \Rightarrow T * \text{number} - \text{id} \Rightarrow T * F - \text{id} \Rightarrow T - \text{id} \Rightarrow E - \text{id} \Rightarrow E - F \Rightarrow E - T \Rightarrow E$

Stack (top is right) Input

Action

\$	id*number-id	shift
\$ id	*number-id	reduce by $F \rightarrow \text{id}$
\$ F	*number-id	reduce by $T \rightarrow F$
\$ T	*number-id	shift
\$ T *	number-id	shift
\$ T * number	-id	reduce by $F \rightarrow \text{number}$
\$ T * F	-id	reduce by $T \rightarrow T * F$
\$ T	-id	reduce by $E \rightarrow T$
\$ E	-id	shift
\$ E -	id	shift
\$ E - id	\$	reduce by $F \rightarrow \text{id}$
\$ E - F	\$	reduce by $T \rightarrow F$
\$ E - T	\$	reduce by $E \rightarrow E - T$
\$ E	\$	accept



Bottom Up Parsing (TB Sec 4.5)

“Grammars used in compiling usually fall into the LR(1) class” (TB, p239)...but:

“An ambiguous grammar can never be LR” (TB, p239)

Two kinds of non-LR conflict in shift-reduce parsing:

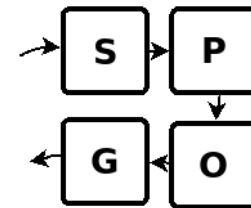
- 1) **Shift/reduce conflict**: choose whether to shift or reduce
- 2) **Reduce/reduce conflict**: two+ reduction possibilities

These occur when **generating** the parser, not in parsing!
Both signal that the grammar is NOT LR(x)!

Parser generator could give up and signal an error...or, can simply choose how to resolve the conflict, and warn you!

Yacc: chooses shift in shift/reduce conflict, chooses topmost rule in yacc source file to reduce in reduce/reduce conflict

Important: the resulting parser is NOT ambiguous!



Simple LR (SLR) Parsing (TB Sec 4.6)

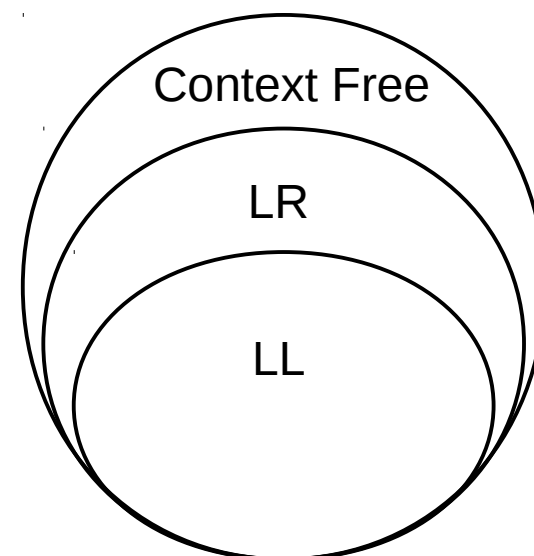
Sec 4.5 describes shift-reduce parsing, but did not detail any specific mechanism for identifying handles, generating a parser, or parsing

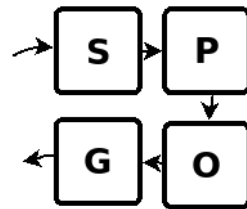
Sec 4.6 presents “the easiest method for constructing shift-reduce parsers” (TB, p241), called Simple LR, or SLR

Sec 4.7 presents more complicated methods, one of which is LALR, which is what yacc uses; we won’t do this section

As in LL, basic goal is to construct a parsing table which encodes our rules and actions

“The class of grammars that can be parsed using LR methods is a proper superset of [those] that can be parsed with [.] LL methods” (TB, p242)





Simple LR (SLR) Parsing (TB Sec 4.6)

An LR(0) **item**, of a grammar G is a production rule of G with a **dot** at some position of the body of the rule; the dot represents how far the rule has been matched and what is hoped to be seen next on the input.

$S \rightarrow AB \mid \varepsilon$

is shorthand for two rules; all dot positions are:

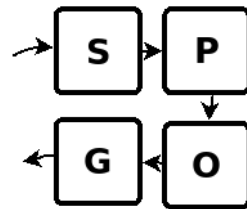
$S \rightarrow .AB$

$S \rightarrow A.B$

$S \rightarrow AB.$

$S \rightarrow .$

The $S \rightarrow \varepsilon$ rule is special: an empty rule generates only one dot position.



Simple LR (SLR) Parsing (TB Sec 4.6)

parser generator goal: construct “a deterministic finite automaton that is used to make parsing decisions. Such an automaton is called an LR(0) automaton. In particular, each state of the LR(0) automaton represents a **set of items** in the **canonical LR(0) collection.**” (TB p243)

Two functions over set I of items: Closure and Goto

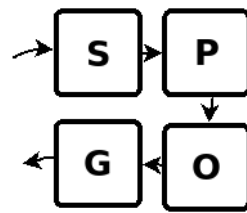
Closure(I):

1. add all items in I to Closure(I)
2. For item $A \rightarrow \alpha.B\beta$ in I and a production rule $B \rightarrow \gamma$, add item $B \rightarrow .\gamma$ to Closure(I) if not already in it. Repeat (2) until no new items added

Goto(I,X), where X is a grammar symbol:

1. Form set N of items $A \rightarrow \alpha X.\beta$ such that $A \rightarrow \alpha.X\beta$ is in I
2. Calculate and return Closure(N)

Simple LR (SLR) Parsing (TB Sec 4.6)

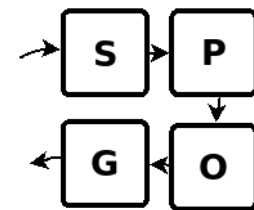


First: augment grammar with starting symbol S to have a new initial rule $S' \rightarrow S$; we accept when just about to reduce this rule

Figure 4.33: Algorithm for computing canonical LR(0) item sets

```
C = {Closure( $S' \rightarrow .S$ )}
repeat:
  foreach (unprocessed I in C):
    foreach (grammar symbol X):
      J = Goto(I,X)
      if (J not empty and not in C):
        add J to C
    mark I as processed
until nothing new added to C
```

Our parsing automaton has a state for each set in the canonical LR(0) item sets, and the Goto() function defines the transitions (which occur on grammar symbols).



Simple LR (SLR) Parsing (TB Sec 4.6)

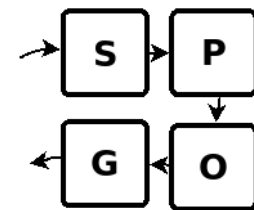
$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

(line #s below are just for reference)

1. Augment with $E' \rightarrow E$
2. Closure($E' \rightarrow .E$) is **S1** = $\{[E' \rightarrow .E], [E \rightarrow .E+T], [E \rightarrow .E-T], [E \rightarrow .T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
3. Goto(**S1**, E) = Closure($\{[E' \rightarrow E.], [E \rightarrow E.+T], [E \rightarrow E.-T]\}$) = **S2** = $\{[E' \rightarrow E.], [E \rightarrow E.+T], [E \rightarrow E.-T]\}$
4. Goto(**S1**, T) = Closure($\{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$) = **S3** = $\{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$
5. Goto(**S1**, F) = Closure($\{[T \rightarrow F.]\}$) = **S4** = $\{[T \rightarrow F.]\}$
6. Goto(**S1**, [+ , - , * , / ,)]) = $\{\}$
7. Goto(**S1**, '(') = Closure($\{[F \rightarrow (.E)]\}$) = **S5** = $\{[F \rightarrow (.E)], [E \rightarrow .E+T], [E \rightarrow .E-T], [E \rightarrow .T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
8. Goto(**S1**, id) = Closure($\{[F \rightarrow \text{id}.]\}$) = **S6** = $\{[F \rightarrow \text{id}.]\}$
9. Goto(**S1**, number) = Closure($\{[F \rightarrow \text{number}.]\}$) = **S7** = $\{[F \rightarrow \text{number}.]\}$



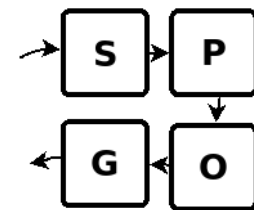
Simple LR (SLR) Parsing (TB Sec 4.6)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

1. $S2 = \{[E' \rightarrow E.], [E \rightarrow E.+T], [E \rightarrow E.-T]\}$
2. $\text{Goto}(S2, [E, T, F, *, /, (,), \text{id}, \text{number}]) = \{\}$
3. $\text{Goto}(S2, +) = \text{Closure}(\{[E \rightarrow E+.T]\}) = \mathbf{S8} = \{[E \rightarrow E+.T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
4. $\text{Goto}(S2, -) = \text{Closure}(\{[E \rightarrow E-.T]\}) = \mathbf{S9} = \{[E \rightarrow E-.T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
5. $S3 = \{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$
6. $\text{Goto}(S3, [E, T, F, +, -, (,), \text{id}, \text{number}]) = \{\}$
7. $\text{Goto}(S3, *) = \text{Closure}(\{[T \rightarrow T*.F]\}) = \mathbf{S10} = \{[T \rightarrow T*.F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
8. $\text{Goto}(S3, /) = \text{Closure}(\{[T \rightarrow T/.F]\}) = \mathbf{S11} = \{[T \rightarrow T/.F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$



Simple LR (SLR) Parsing (TB Sec 4.6)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

1. $S4 = \{[T \rightarrow F.]\}$

2. $\text{Goto}(S4, [\text{all}]) = \{\}$

3. $S5 = \{[F \rightarrow (.E)], [E \rightarrow .E+T], [E \rightarrow .E-T], [E \rightarrow .T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$

4. $\text{Goto}(S5, E) = \text{Closure}(\{[F \rightarrow (E.)], [E \rightarrow E.+T], [E \rightarrow E.-T]\}) = \mathbf{S12} = \{[F \rightarrow (E.)], [E \rightarrow E.+T], [E \rightarrow E.-T]\}$

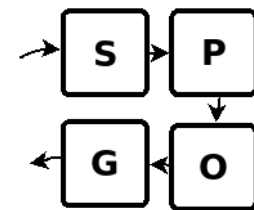
5. $\text{Goto}(S5, T) = \text{Closure}(\{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}) = \mathbf{S13} = \{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$

6. $\text{Goto}(S5, F) = \text{Closure}(\{[T \rightarrow F.]\}) = \mathbf{S4}$

7. $\text{Goto}(S5, '(') = \mathbf{S5}$ $\text{Goto}(S5, \text{id}) = \mathbf{S6}$ $\text{Goto}(S5, \text{num}) = \mathbf{S7}$

8. $S6 = \{[F \rightarrow \text{id}.]\}$, $\text{Goto}(S6, \text{all}) = \{\}$

9. $S7 = \{[F \rightarrow \text{number}.]\}$, $\text{Goto}(S7, \text{all}) = \{\}$



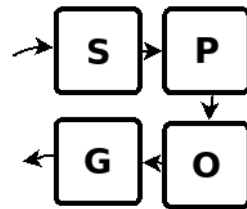
Simple LR (SLR) Parsing (TB Sec 4.6)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

1. $S8 = \{[E \rightarrow E+.T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
2. $\text{Goto}(S8, T) = \text{Closure}(\{[E \rightarrow E+T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}) = \mathbf{S14} = \{[E \rightarrow E+T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$
3. $\text{Goto}(S8, F) = \text{Closure}(\{[T \rightarrow F.]\}) = \mathbf{S4} = \{[T \rightarrow F.]\}$
4. $\text{Goto}(S8, '(') = \text{Closure}(\{[F \rightarrow (.E)]\}) = \mathbf{S5}$
5. $\text{Goto}(S8, \text{id}) = \mathbf{S6}$ $\text{Goto}(S8, \text{number}) = \mathbf{S7}$
6. $\text{Goto}(S8, [E, +, -, *, /, ,)]) = \{\}$
7. $S9 = \{[E \rightarrow E-.T], [T \rightarrow .T*F], [T \rightarrow .T/F], [T \rightarrow .F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
8. $\text{Goto}(S9, T) = \text{Closure}(\{[E \rightarrow E-T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}) = \mathbf{S15} = \{[E \rightarrow E-T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$
9. $\text{Goto}(S9, F) = \mathbf{S4}$ $\text{Goto}(S9, '(') = \mathbf{S5}$
10. $\text{Goto}(S9, \text{id}) = \mathbf{S6}$ $\text{Goto}(S9, \text{number}) = \mathbf{S7}$
11. $\text{Goto}(S9, [E, +, -, *, /, ,)]) = \{\}$



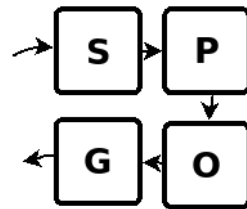
Simple LR (SLR) Parsing (TB Sec 4.6)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

1. $S_{10} = \{[T \rightarrow T*.F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
2. $\text{Goto}(S_{10}, F) = \text{Closure}(\{[T \rightarrow T*F.]\}) = \mathbf{S_{16}} = \{[T \rightarrow T*F.]\}$
3. $\text{Goto}(S_{10}, '(') = \mathbf{S_5}$ $\text{Goto}(S_{10}, \text{id}) = \mathbf{S_6}$ $\text{Goto}(S_{10}, \text{number}) = \mathbf{S_7}$
4. $\text{Goto}(S_{10}, [E, T, +, -, *, /, ,)]) = \{\}$
5. $S_{11} = \{[T \rightarrow T/.F], [F \rightarrow .(E)], [F \rightarrow .\text{id}], [F \rightarrow .\text{number}]\}$
6. $\text{Goto}(S_{11}, F) = \text{Closure}(\{[T \rightarrow T/F.]\}) = \mathbf{S_{17}} = \{[T \rightarrow T/F.]\}$
7. $\text{Goto}(S_{11}, '(') = \mathbf{S_5}$ $\text{Goto}(S_{11}, \text{id}) = \mathbf{S_6}$ $\text{Goto}(S_{11}, \text{number}) = \mathbf{S_7}$
8. $\text{Goto}(S_{11}, [E, T, +, -, *, /, ,)]) = \{\}$
9. $S_{12} = \{[F \rightarrow (E.)], [E \rightarrow E.+T], [E \rightarrow E.-T]\}$
10. $\text{Goto}(S_{12}, +) = \mathbf{S_8}$ $\text{Goto}(S_{12}, -) = \mathbf{S_9}$ $\text{Goto}(S_{12}, \text{else}) = \{\}$
11. $\text{Goto}(S_{12}, ')') = \text{Closure}(\{[F \rightarrow (E).]\}) = \{[F \rightarrow (E).]\} = \mathbf{S_{18}}$



Simple LR (SLR) Parsing (TB Sec 4.6)

$E \rightarrow E+T \mid E-T \mid T$

$T \rightarrow T*F \mid T/F \mid F$

$F \rightarrow (E) \mid \text{id} \mid \text{number}$

1. $S13 = \{[E \rightarrow T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$

2. $\text{Goto}(S13, *) = \mathbf{S10}$ $\text{Goto}(S13, /) = \mathbf{S11}$ $\text{Goto}(S13, \text{else}) = \{\}$

3. $S14 = \{[E \rightarrow E+T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$

4. $\text{Goto}(S14, *) = \mathbf{S10}$ $\text{Goto}(S14, /) = \mathbf{S11}$ $\text{Goto}(S14, \text{else}) = \{\}$

5. $S15 = \{[E \rightarrow E-T.], [T \rightarrow T.*F], [T \rightarrow T./F]\}$

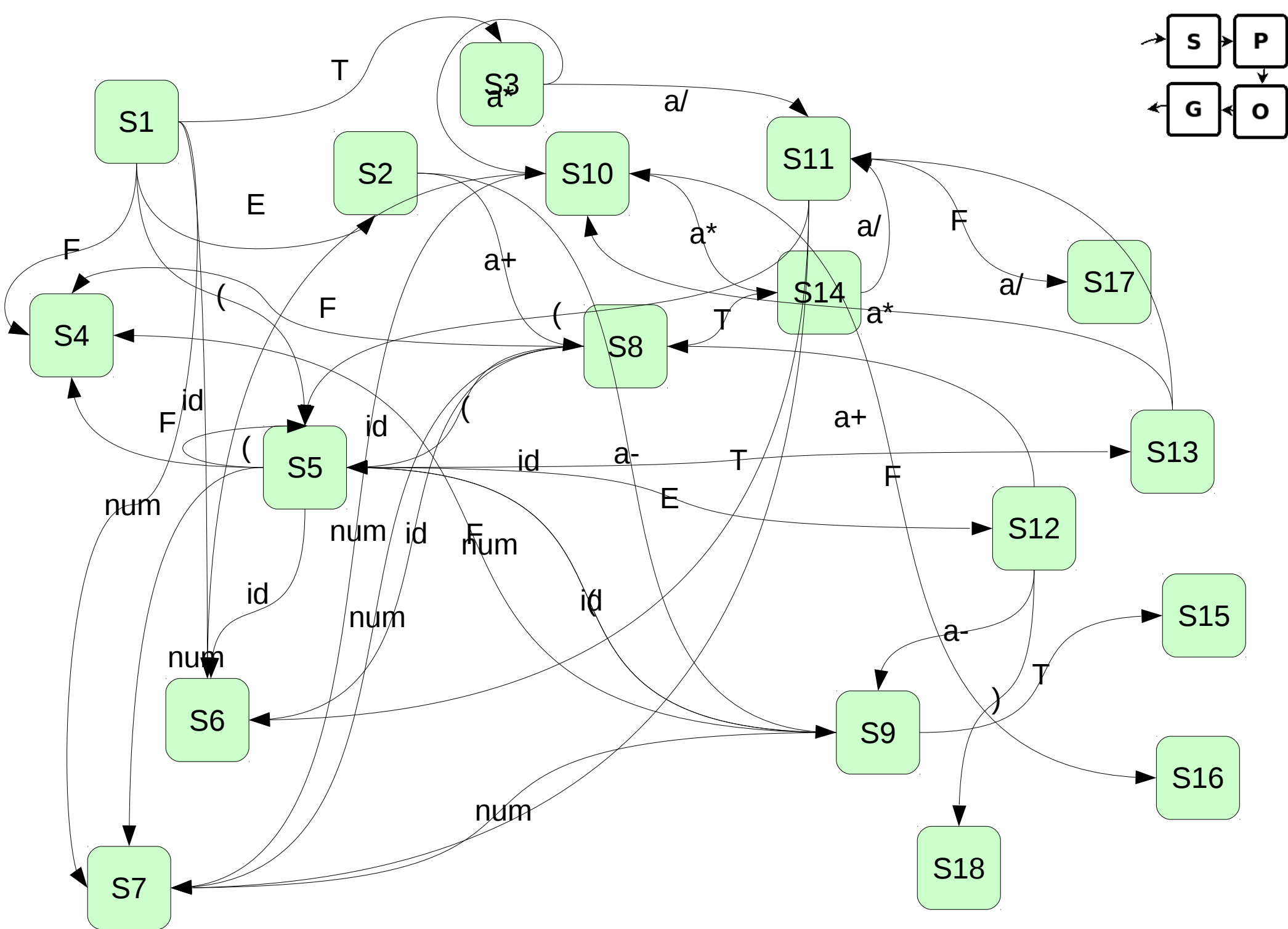
6. $\text{Goto}(S15, *) = \mathbf{S10}$ $\text{Goto}(S15, /) = \mathbf{S11}$ $\text{Goto}(S15, \text{else}) = \{\}$

7. $S16 = \{[T \rightarrow T*F.]\}$ $\text{Goto}(S16, \text{all}) = \{\}$

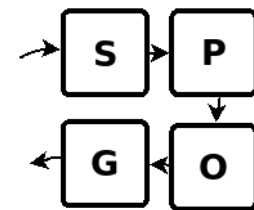
8. $S17 = \{[T \rightarrow T/F.]\}$ $\text{Goto}(S17, \text{all}) = \{\}$

9. $S18 = \{[F \rightarrow (E).]\}$ $\text{Goto}(S18, \text{all}) = \{\}$

We're done!, Now draw it!



	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13	S14	S15	S16	S17	P
S1		E	T	F	(id	nm											↓
S2								+	-									O
S3										*	/							
S4																		
S5				F	(id	nm					E	T					
S6																		
S7																		
S8				F	(id	nm							T				
S9				F	(id	nm								T			
S10					(id	nm									F		
S11					(id	nm										F	
S12								+	-									
S13										*	/							
S14										*	/							
S15										*	/							
S16																		
S17																		
S18)						



Creating the Parsing Table(s) (Alg 4.46, p253)

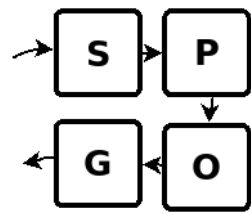
Two tables: Action and GoTo (but we can put it in one)

For each Item set I_i :

1. If $[A \rightarrow \alpha.a\beta]$ is in I_i and $\text{Goto}(I_i, a) = I_j$, then $\text{Action}[i, a] = \text{"shift } j\text{"}$
2. if $[A \rightarrow \alpha.]$ is in I_i , then $\text{Action}[i, a] = \text{"reduce } A \rightarrow \alpha\text{"}$ for all a in $\text{FOLLOW}(A)$
3. if $[S' \rightarrow S.]$ in I_i then $\text{Action}[i, \$] = \text{"accept"}$
4. All other Action table entries are "error"

GoTo table:

5. If $\text{Goto}(I_i, A) = I_j$ for nonterminal A , then $\text{GoTo}(i, A) = j$
6. Initial state is from item set containing $[S' \rightarrow .S]$



Creating the Parsing Table(s) (Alg 4.46, p253)

Follow(E,T,F) = {+, -, *, /,), \$}

S1: [1, (] = shift 5, [1, id] = shift 6, [1, num] = shift 7
Goto(1, E) = 2 Goto(1, T) = 3 Goto(1, F) = 4

S2: [2, +] = shift 8, [2, -] = shift 9; [2, \$] = accept

S3: [3, *] = shift 10, [3, /] = shift 11, [3, {+, -, *, /,), \$}] = reduce $E \rightarrow T$ (E3)

S4: [4, {+, -, *, /,), \$}] = reduce $T \rightarrow F$ (T3)

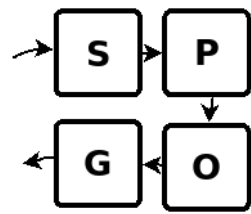
S5: [5, (] = shift 5, [5, id] = shift 6, [5, num] = shift 7
Goto(5, E) = 12, Goto(5, T) = 13 Goto(5, F) = 4

S6: [6, {+, -, *, /,), \$}] = reduce $F \rightarrow \text{id}$ (F2)

S7: [7, {+, -, *, /,), \$}] = reduce $F \rightarrow \text{num}$ (F3)

S8: [8, (] = shift 5, [8, id] = shift 6, [8, num] = shift 7, Goto(8, F)=4, Goto(8, T)=14

S9: [9, (] = shift 5, [9, id] = shift 6, [9, num] = shift 7, Goto(9, F)=4, Goto(9, T)=15



Creating the Parsing Table(s) (Alg 4.46, p253)

S10: [10,(] = shift 5, [10,id] = shift 6, [10,num] = shift 7, goto(10,F)=16

S11: [11,(] = shift 5, [11,id] = shift 6, [11,num] = shift 7, goto(11,F)=17

S12: [12,+] = shift 8, [12,-] = shift 9

S13: [13,*] = shift 10, [13,/] = shift 11 [13,{+,-,*,/,),,\$}] = reduce $E \rightarrow T$ (E3)

S14: [14,*] = shift 10, [14,/] = shift 11 [14,{+,-,*,/,),,\$}] = reduce $E \rightarrow E+T$ (E1)

S15: [15,*] = shift 10, [15,/] = shift 11 [15,{+,-,*,/,),,\$}] = reduce $E \rightarrow E-T$ (E2)

S16: [16,{+,-,*,/,),,\$}] = reduce $T \rightarrow T * F$ (T1)

S17: [17,{+,-,*,/,),,\$}] = reduce $T \rightarrow T / F$ (T2)

S18: [18,{+,-,*,/,),,\$}] = reduce $F \rightarrow (E)$ (F1)

	E	T	F	id	nm	+	-	*	/	()	\$
S1	G 2	G 3	G 4	S 6	S 7					S 5		
S2						S 8	S 9					accept
S3						R E3	R E3	S 10	S 11		R E3	R E3
S4						R T3	R T3	R T3	R T3		R T3	R T3
S5	G 12	G 13	G 4	S 6	S 7							
S6						R F2	R F2	R F2	R F2		R F2	R F2
S7						R F3	R F3	R F3	R F3		R F3	R F3
S8		G 14	G 4	S 6	S 7					S 5		
S9		G 15	G 4	S 6	S 7					S 5		
S10			G 16	S 6	S 7					S 5		
S11			G 17	S 6	S 7					S 5		
S12						S 8	S 9					
S13						R E3	R E3	S 10	S 11		R E3	R E3
S14						R E1	R E1	S 10	S 11		R E1	R E1
S15						R E2	R E2	S 10	S 11		R E2	R E2
S16						R T1	R T1	R T1	R T1		R T1	R T1
S17						R T2	R T2	R T2	R T2		R T2	R T2
S18						R F1	R F1	R F1	R F1		R F1	R F1

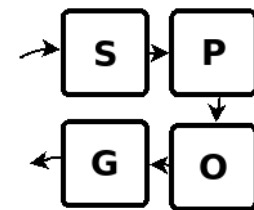


Table-Driven Parsing Algorithm (Fig 4.36, p251)

- NOTE: stack is a stack of states, not of grammar symbols!

push starting state onto stack

$a \leftarrow$ first symbol of $w\$$

while (1) {

$s \leftarrow$ top state of stack (no pop)

 if ($\text{Action}[s,a] = \text{shift } t$) {

 push t onto stack

$a \leftarrow$ next symbol of $w\$$

 } else if ($\text{Action}[s,a] = \text{reduce } A \rightarrow \beta$) {

 pop $|\beta|$ symbols off the stack

$t \leftarrow$ top state of stack

 push $\text{Goto}(t,A)$ onto stack

 output production $A \rightarrow \beta$

 } else if ($\text{Action}[s,a] = \text{accept}$) {

 break;

 } else {

 print syntax error and try to recover

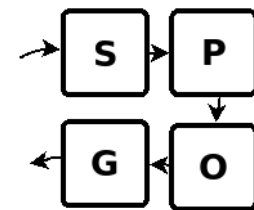
 }

}

print success; output rules are RM derivation

Notes:

- states represent symbols *in context*
- reduction still happens when a *handle* is seen
- reduction still consumes N stack items ($N ==$ size of RHS of rule)

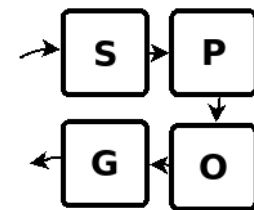


Bottom Up Parsing Using above Table

string: "i*42-j" == id * number - id

Stack (top is right) Input

1	id*number-id\$	shift 6
1 6	*number-id\$	reduce by F → id, goto 4
1 4	*number-id\$	reduce by T → F, goto 3
1 3	*number-id\$	shift 10
1 3 10	number-id\$	shift 7
1 3 10 7	-id\$	reduce by F → number, goto 16
1 3 10 16	-id\$	reduce by T → T * F, goto 3
1 3	-id\$	reduce by E → T, goto 2
1 2	-id\$	shift 9
1 2 9	id\$	shift 6
1 2 9 6	\$	reduce by F → id, goto 4
1 2 9 4	\$	reduce by T → F, goto 15
1 2 9 15	\$	reduce by E → E-T, goto 2
1 2	\$	accept!



Bottom Up Parsing: Example #2

Grammar:

$S \rightarrow A S \mid \text{empty}$

$A \rightarrow a B$

$B \rightarrow b B \mid b$

augment with: $S' \rightarrow S$

$\text{FOLLOW}(S) = \$$

$\text{FOLLOW}(A) = a, \$$

$\text{FOLLOW}(B) = a, \$$

Item sets:

$S1 = \text{Closure}(S' \rightarrow .S) = \{S' \rightarrow .S, S \rightarrow .AS, S \rightarrow ., A \rightarrow .aB\}$

$\text{Goto}(S1, S) = \text{Closure}(S' \rightarrow S.) = S2 = \{S' \rightarrow S.\}$

$\text{Goto}(S1, A) = \text{Closure}(S \rightarrow A.S) = S3 = \{S \rightarrow A.S, S \rightarrow .AS, S \rightarrow ., A \rightarrow .aB\}$

$\text{Goto}(S1, a) = \text{Closure}(A \rightarrow a.B) = S4 = \{A \rightarrow a.B, B \rightarrow .bB, B \rightarrow .b\}$

$\text{Goto}(S3, S) = \text{Closure}(S \rightarrow AS.) = S5 = \{S \rightarrow AS.\}$

$\text{Goto}(S3, A) = \text{Closure}(S \rightarrow A.S) = S3$

$\text{Goto}(S3, a) = \text{Closure}(A \rightarrow a.B) = S4$

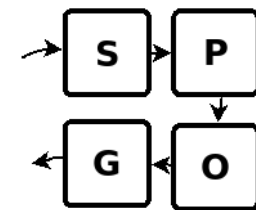
$\text{Goto}(S4, B) = \text{Closure}(A \rightarrow aB.) = S7 = \{A \rightarrow aB.\}$

$\text{Goto}(S4, b) = \text{Closure}(\{B \rightarrow b.B, B \rightarrow b.\}) = S6 = \{B \rightarrow b.B, B \rightarrow b., B \rightarrow .bB, B \rightarrow .b\}$

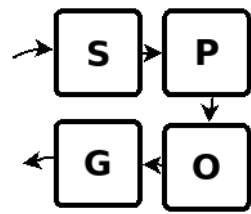
$\text{Goto}(S6, B) = \text{Closure}(\{B \rightarrow bB.\}) = S8 = \{B \rightarrow bB.\}$

$\text{Goto}(S6, b) = \text{Closure}(B \rightarrow b.B, B \rightarrow b.) = S6$

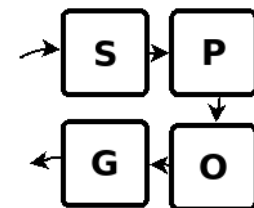
all other goto's are empty; S2 is accept



	S1	S2	S3	S4	S5	S6	S7	S8
S1		S	A	a				
S2								
S3			A	a	S			
S4						b	B	
S5								
S6						b		B
S7								
S8								



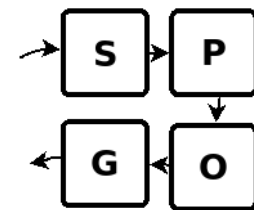
	S	A	B	a	b	\$
S1	goto 2	goto 3		shift 4		reduce $S \rightarrow e$
S2						accept
S3	goto 5	goto 3		shift 4		reduce $S \rightarrow e$
S4			goto 7		shift 6	
S5						reduce $S \rightarrow AS$
S6			goto 8	reduce $B \rightarrow b$	shift 6	reduce $B \rightarrow b$
S7				reduce $A \rightarrow aB$		reduce $A \rightarrow aB$
S8				reduce $B \rightarrow bB$		reduce $B \rightarrow bB$



Bottom Up Parsing: Example #2

string: "abbab"

Stack (top is right)	Input	Action
1	abbab\$	shift 4
1 4	bbab\$	shift 6
1 4 6	bab\$	shift 6
1 4 6 6	ab\$	reduce $B \rightarrow b$, goto 8
1 4 6 8	ab\$	reduce $B \rightarrow bB$, goto 7
1 4 7	ab\$	reduce $A \rightarrow aB$, goto 3
1 3	ab\$	shift 4
1 3 4	b\$	shift 6
1 3 4 6	\$	reduce $B \rightarrow b$, goto 7
1 3 4 7	\$	reduce $A \rightarrow aB$, goto 3
1 3 3	\$	reduce $S \rightarrow e$, goto 5
1 3 3 5	\$	reduce $S \rightarrow AS$, goto 5
1 3 5	\$	reduce $S \rightarrow AS$, goto 2
1 2	\$	accept!



More Complicated LR Parsing (TB 4.7)

4.6 was just “simple LR”, SLR

4.7 introduces more complicated (and powerful) LR parsing

LR, or “canonical” LR: idea, add a terminal symbol to the definition of SLR items; this terminal will constrain the reduce actions to choose which rule to reduce based on the terminal (SLR applies the same reduce to all following terminals). This LR(1) method ends up having *many* items!

LALR, or “lookahead LR”: LR(1) generates many item sets with common “cores”; these can be merged (carefully)

LALR produces a number of item sets on the order of SLR, while LR(1) can be an order of magnitude bigger