



CS415 Compilers

Syntax Analysis Part 3

These slides are based on slides copyrighted by
Keith Cooper, Ken Kennedy & Linda Torczon at Rice
University

- Midterm has been graded. Exams will be handed out in recitation. Please see canvas for grades.
- Third homework:
Due Thursday, March 10
- First project (local instruction scheduler) deadlines:
code: March 10 @ 11:59pm - single tar file
report: March 11 @ 11:59pm - single pdf file
Late policy:
 - Grace period: 1 hour
 - 20% penalty for every started 24 hour period after the deadline.
 - Saturday/Sunday count as a single 24 hour period.
- Warning grades will be submitted by Friday, March 11

Parsing (Syntax Analysis)

Top-Down Parsing
EAC Chapters 3.3

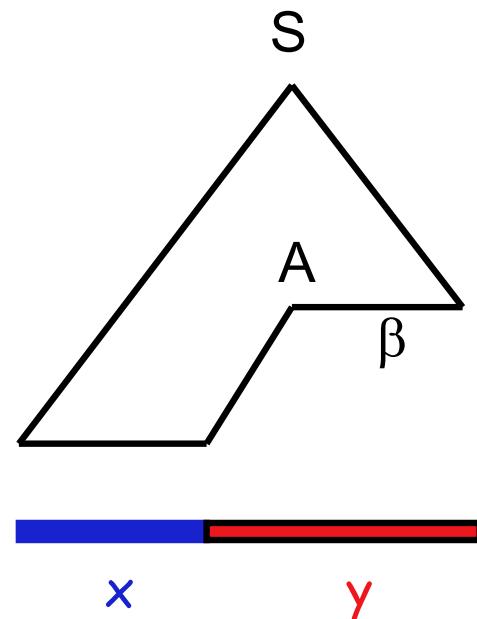
LL(1), recursive descent

1 input symbol lookahead

construct leftmost derivation (forwards)

input: read left-to-right

$$S \Rightarrow^*_{Im} x A \beta \Rightarrow_{Im} x \delta \beta \Rightarrow^*_{Im} x y$$



LL(1), recursive descent

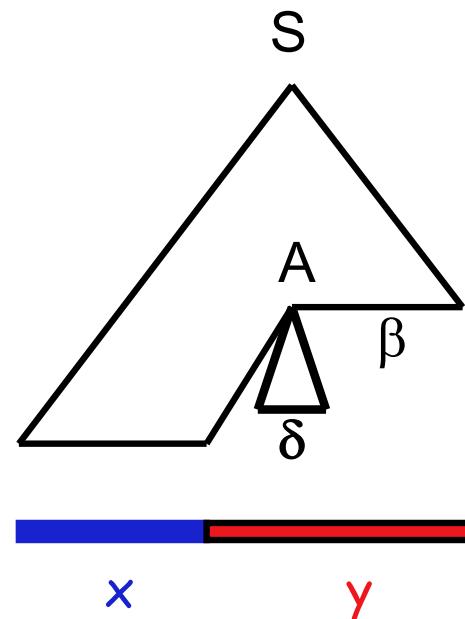
1 input symbol lookahead

construct leftmost derivation (forwards)

input: read left-to-right

$S \Rightarrow^*_{Im} x A \beta \Rightarrow_{Im} x \delta \beta \Rightarrow^*_{Im} x y$

rule $A \rightarrow \delta$



$$a \in \text{FIRST}_1(\alpha) \text{ iff } \alpha \Rightarrow^* a\gamma, \text{ for some } \gamma$$

To build $\text{FIRST}(X)$ for all grammar symbols X :

1. if X is a terminal (token), $\text{FIRST}(X) := \{ X \}$
2. if $X \rightarrow \varepsilon$, then $\varepsilon \in \text{FIRST}(X)$
3. iterate until no more terminals or ε can be added to any $\text{FIRST}(X)$:

if $X \rightarrow Y_1 Y_2 \dots Y_k$ then

$a \in \text{FIRST}(X)$ if $a \in \text{FIRST}(Y_i)$ and

$\varepsilon \in \text{FIRST}(Y_j)$ for all $1 \leq j < i$

$\varepsilon \in \text{FIRST}(X)$ if $\varepsilon \in \text{FIRST}(Y_i)$ for all $1 \leq i \leq k$

end iterate

Note: if $\varepsilon \notin \text{FIRST}(Y_1)$, then $\text{FIRST}(Y_i)$ is irrelevant, for $1 < i$

$$a \in \text{FIRST}(\alpha) \text{ iff } \alpha \Rightarrow^* a\gamma, \text{ for some } \gamma$$

To build $\text{FIRST}(\alpha)$ for $\alpha = X_1 X_2 \dots X_n$:

1. $a \in \text{FIRST}(\alpha)$ if $a \in \text{FIRST}(X_i)$ and
 $\varepsilon \in \text{FIRST}(X_j)$ for all $1 \leq j < i$
2. $\varepsilon \in \text{FIRST}(\alpha)$ if $\varepsilon \in \text{FIRST}(X_i)$ for all $1 \leq i \leq n$

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

$\text{FOLLOW}(A) := \text{the set of terminals that can appear immediately to the right of } A \text{ in some sentential form.}$

Thus, a non-terminal's FOLLOW set specifies the tokens that can legally appear after it; a terminal has no FOLLOW set

$$\text{FOLLOW}(A) = \{ a \in (T \cup \{\text{eof}\}) \mid S \text{ eof} \Rightarrow^* \alpha A a \gamma \}$$

To build $\text{FOLLOW}(X)$ for all non-terminal X :

1. Place eof in $\text{FOLLOW}(\langle \text{goal} \rangle)$

iterate until no more terminals or eof can be added
to any $\text{FOLLOW}(X)$:

2. If $A \rightarrow \alpha B \beta$ then

put $\{\text{FIRST}(\beta) - \varepsilon\}$ in $\text{FOLLOW}(B)$

3. If $A \rightarrow \alpha B$ then

put $\text{FOLLOW}(A)$ in $\text{FOLLOW}(B)$

4. If $A \rightarrow \alpha B \beta$ and $\varepsilon \in \text{FIRST}(\beta)$ then

put $\text{FOLLOW}(A)$ in $\text{FOLLOW}(B)$

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\varepsilon \in \text{FIRST}(\alpha)$, then we need to ensure that $\text{FIRST}(\beta)$ is disjoint from $\text{FOLLOW}(A)$, too

Define $\text{FIRST}^+(\delta)$ for rule $A \rightarrow \delta$ as

- $(\text{FIRST}(\delta) - \{\varepsilon\}) \cup \text{FOLLOW}(A)$, if $\varepsilon \in \text{FIRST}(\delta)$
- $\text{FIRST}(\delta)$, otherwise

The LL(1) Property

A grammar is $LL(1)$ iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies
 $FIRST^+(\alpha) \cap FIRST^+(\beta) = \emptyset$

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !

Question: Can there be two rules $A \rightarrow \alpha$ and $A \rightarrow \beta$ in a $LL(1)$ grammar such that $\varepsilon \in FIRST(\alpha)$ and $\varepsilon \in FIRST(\beta)$?

Given a grammar that has the $LL(1)$ property

- Problem: NT "A" needs to be replaced in next derivation step
- Assume $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with
 $\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_2) = \emptyset$, $\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_3) = \emptyset$, and
 $\text{FIRST}^+(\beta_2) \cap \text{FIRST}^+(\beta_3) = \emptyset$ (pair-wise disjoint sets)

```
/* find rule for A */
if (current token ∈ FIRST+ (β1))
    select A → β1
else if (current token ∈ FIRST+(β2))
    select A → β2
else if (current token ∈ FIRST+(β3))
    select A → β3
else
    report an error and return false
```

Grammars with the $LL(1)$ property are called **predictive grammars** because the parser can "predict" the correct expansion at each point in the parse.

Parsers that capitalize on the $LL(1)$ property are called **predictive parsers**.

One kind of predictive parser is the **recursive descent** parser. The other is a table-driven parser **table-driven parser**.

Is the following grammar LL(1)?

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

Is the following grammar LL(1)?

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

$$\text{First}(aSb) = \{ a \}$$

$$\text{First}(\epsilon) = \{ \epsilon \}$$

$$\text{Follow}(S). = \{ \text{eof}, b \}$$

$$\text{First}^+(aSb) = \{ a \}$$

$$\text{First}^+(\epsilon) = (\text{First}(\epsilon) - \{ \epsilon \}) \cup \text{Follow}(S) = \{ \text{eof}, b \}$$

LL(1)?

Is the following grammar LL(1)?

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

$$\text{First}(aSb) = \{ a \}$$

$$\text{First}(\epsilon) = \{ \epsilon \}$$

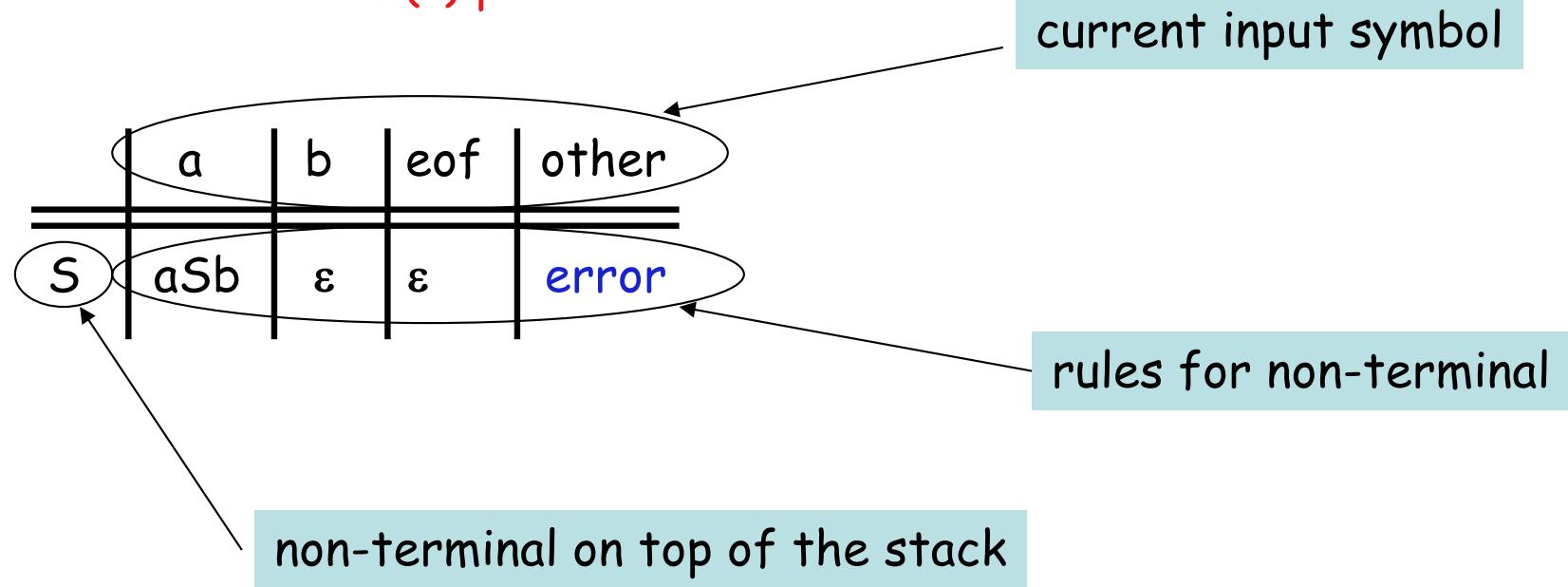
$$\text{Follow}(S). = \{ \text{eof}, b \}$$

$$\text{First}^+(aSb) = \{ a \}$$

$$\text{First}^+(\epsilon) = (\text{First}(\epsilon) - \{ \epsilon \}) \cup \text{Follow}(S) = \{ \text{eof}, b \}$$

LL(1)? YES, since $\{ a \} \cap \{ \text{eof}, b \} = \emptyset$

Table-driven LL(1) parser



Building the complete table

- Need a row for every NT & a column for every $T + \text{"eof"}$
- Need an algorithm to build the table

Filling in $\text{TABLE}[X,y]$, $X \in NT$, $y \in T \cup \{\text{eof}\}$

- entry is the rule $X \rightarrow \beta$, if $y \in \text{FIRST}^+(\beta)$
- entry is **error** otherwise

If any entry is defined multiple times, G is not $LL(1)$

This is the $LL(1)$ table construction algorithm

```
token ← next_token() // scanner call
push EOF onto Stack // bottom of Stack marker
push the start symbol,  $S$ , onto Stack
TOS ← top of Stack
loop forever
    if TOS = EOF and token = EOF then
        break & report success
    else if TOS is a terminal then
        if TOS matches token then
            pop Stack           // recognized TOS
            token ← next_token()
        else report error looking for TOS
    else                                // TOS is a non-terminal
        if TABLE[TOS,token] is  $A \rightarrow B_1 B_2 \dots B_k$  then
            pop Stack           // get rid of A
            push  $B_k, B_{k-1}, \dots, B_1$  // in that order
        else report error expanding TOS
    TOS ← top of Stack
```



Table-driven *LL(1)* parser

	a	b	eof	other
S	aSb	ϵ	ϵ	error

How to parse input a a a b b b ?

Describe action as sequence of states

(PDA stack content, remaining input, next action)

use eof as bottom-of-stack marker

PDA stack content: [X, ... Z], where Z is the TOS

next actions: rule or next input+pop or error or accept

Table-driven *LL(1)* parser

	a	b	eof	other
S	aSb	ϵ	ϵ	error

([eof, **S**], **aaabbb**, aSb) \Rightarrow
([eof, b, S, **a**], **aaabbb**, next input+pop) \Rightarrow
([eof, b, **S**], **aabbbb**, aSb) \Rightarrow
([eof, b, b, S, **a**], **aabbbb**, next input+pop) \Rightarrow
([eof, b, b, **S**], **abbb**, aSb) \Rightarrow
([eof, b, b, b, S, **a**], **abbb**, next input+pop) \Rightarrow
([eof, b, b, b, **S**], **bbb**, ϵ) \Rightarrow
([eof, b, b, **b**], **bbb**, next input+pop) \Rightarrow ([eof, b, **b**], **bb**, next input+pop) \Rightarrow
([eof, **b**], **b**, next input+pop) \Rightarrow ([eof], eof, accept)

Recursive descent LL(1) parser

	a	b	eof	other
S	aSb	ϵ	ϵ	error

1. Every NT is associated with a parsing procedure.
2. The parsing procedure for $A \in \text{NT}$, `proc A`, is responsible to parse and consume any (token) string that can be derived from A ; it may recursively call other parsing procedures.
3. The parser is invoked by calling `proc S` for start symbol S .

Recursive descent LL(1) parser

	a	b	eof	other
S	aSb	ϵ	ϵ	error

```
main () {
    token = next_token();
    if (S () and token = eof)
        print "accept"
    else
        print "error";
}
```

```
bool S () {
    switch token {
        case a: token = next_token();
                  if (not S ()) return false;
                  if token = b
                      {token = next_token(); return true;}
                  else
                      return false;
                      break;
        case b, eof: return true; break;
        default: return false;
    }
}
```

Recursive descent LL(1) parser

	a	b	eof	other
S	aSb	ϵ	ϵ	error

```
main () {
    token = next_token();
    if (S () and token = eof )
        print "accept"
    else
        print "error";
}
```

How to parse input a a a b b b ?

```
bool S () {
    switch token {
        case a: token = next_token();
            if (not S ()) return false;
            if token = b
                {token = next_token(); return true;}
            else
                return false;
                break;
        case b, eof: return true; break;
        default: return false;
    }
}
```

Top-down parsers cannot handle left-recursive grammars

Formally,

A grammar is *left recursive* if $\exists A \in NT$ such that
 \exists a derivation $A \Rightarrow^+ A\alpha$, for some string $\alpha \in (NT \cup T)^+$

Our expression grammar is left recursive

- This can lead to non-termination in a top-down parser
- For a top-down parser, any recursion must be right recursion
- We would like to convert the left recursion to right recursion

Non-termination is a bad property in any part of a compiler

What if my grammar does not have the LL(1) property?

⇒ Sometimes, we can transform the grammar

The Algorithm

$\forall A \in NT,$
*find the longest prefix α that occurs in two
or more right-hand sides of A*
if $\alpha \neq \epsilon$ then replace all of the A productions,
$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma,$$

with
$$A \rightarrow \alpha Z \mid \gamma$$

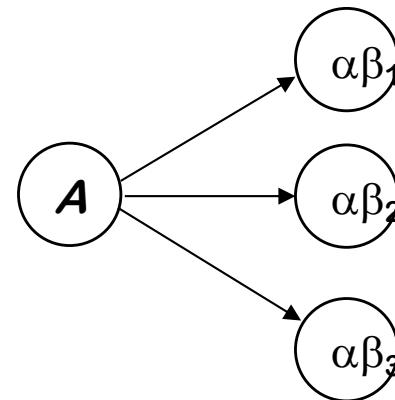
$$Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

where Z is a new element of NT
Repeat until no common prefixes remain

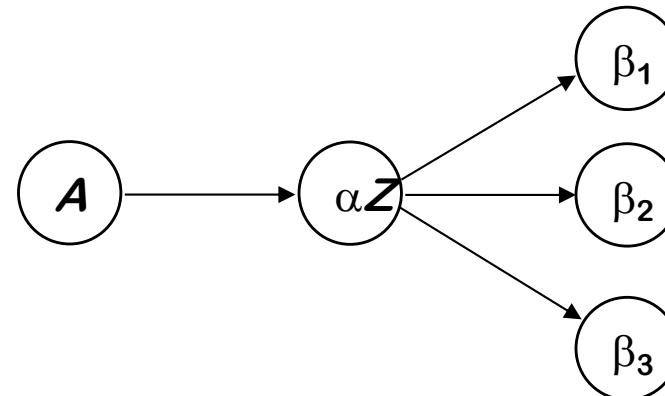
A graphical explanation for the same idea

$$\begin{array}{l} A \rightarrow \alpha\beta_1 \\ | \quad \alpha\beta_2 \\ | \quad \alpha\beta_3 \end{array}$$

becomes ...



$$\begin{array}{l} A \rightarrow \alpha Z \\ Z \rightarrow \beta_1 \\ | \quad \beta_2 \\ | \quad \beta_3 \end{array}$$



Consider the following fragment of the expression grammar

$$\begin{array}{lcl} \text{Factor} & \rightarrow & \underline{\text{Identifier}} \\ & | & \underline{\text{Identifier}} [\text{ExprList}] \\ & | & \underline{\text{Identifier}} (\text{ExprList}) \end{array}$$

$$\begin{aligned} \text{FIRST}(rhs_1) &= \{ \underline{\text{Identifier}} \} \\ \text{FIRST}(rhs_2) &= \{ \underline{\text{Identifier}} \} \\ \text{FIRST}(rhs_3) &= \{ \underline{\text{Identifier}} \} \end{aligned}$$

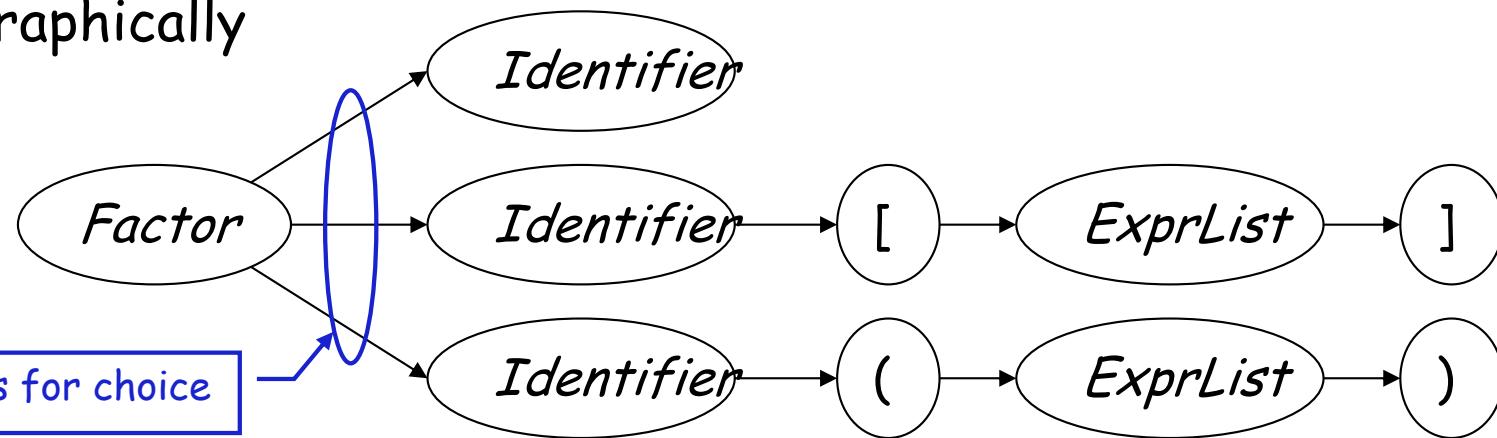
After left factoring, it becomes

$$\begin{array}{lcl} \text{Factor} & \rightarrow & \underline{\text{Identifier}} \text{ Arguments} \\ \text{Arguments} & \rightarrow & [\text{ExprList}] \\ & | & (\text{ExprList}) \\ & | & \epsilon \end{array}$$

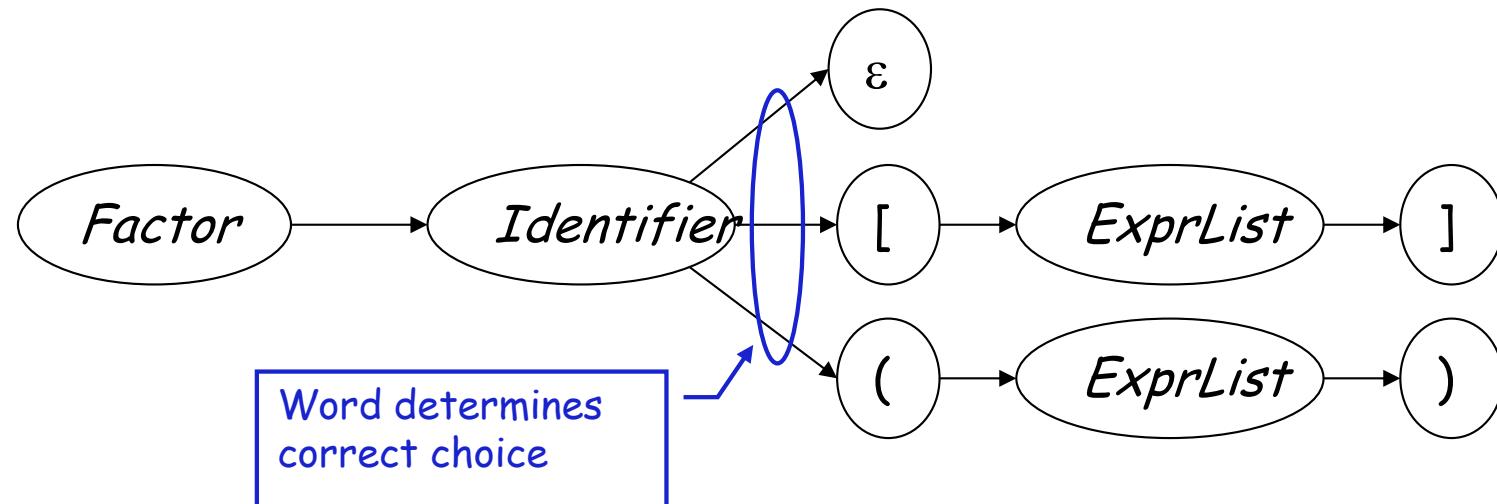
$$\begin{aligned} \text{FIRST}(rhs_1) &= \{ \underline{\text{Identifier}} \} \\ \text{FIRST}(rhs_2) &= \{ [\ } \} \\ \text{FIRST}(rhs_3) &= \{ (\ } \} \\ \text{FIRST}(rhs_4) &= \text{FOLLOW}(\text{Factor}) \\ \Rightarrow & \text{It has the } LL(1) \text{ property} \end{aligned}$$

This form has the same syntax, with the $LL(1)$ property

Graphically



becomes ...



Question

By *eliminating left recursion* and *left factoring*, can we transform an arbitrary CFG to a form where it meets the $LL(1)$ condition? (and can be parsed predictively with a single token lookahead?)

Answer

Given a CFG that doesn't meet the $LL(1)$ condition, it is undecidable whether or not an equivalent $LL(1)$ grammar exists.

Example $\{a^n 0 b^n \mid n \geq 1\} \cup \{a^n 1 b^{2n} \mid n \geq 1\}$ has no $LL(k)$ grammar

$$G \rightarrow \underline{a} A \underline{b}$$

$$\mid \underline{a} B \underline{b} b$$

$$A \rightarrow \underline{a} A \underline{b}$$

$$\mid \underline{0}$$

$$B \rightarrow \underline{a} B \underline{b} b$$

$$\mid \underline{1}$$

Problem: need an unbounded number of "a" symbols before you can determine whether you are in the A group or the B group.

More Syntax Analysis (bottom-up)

LR(1) parsing

Read EaC: Chapter 3.4