



Top-down Parsing Recursive Descent & LL(1)



Roadmap (Where are we?)

- Predictive top-down parsing
 - The LL(1) Property
 - First and Follow sets
 - Simple recursive descent parsers
 - Table-driven LL(1) parsers



LL(1) Parser

- L = scan input left to right
- L = Leftmost derivation
- 1 = lookahead is enough to pick right production rule to use
- No Backtracking
- No Left Recursion



Predictive Parsing

Given production rules

$$A \rightarrow \alpha$$

$$A \rightarrow \beta$$

*the parser should be able to choose
between α or β using one lookahead*

Predictive Parser is a top-down parser free
of backtracking



First Sets

For some rhs $\alpha \in G$

FIRST(α) is set of tokens (terminals) that appear as first symbol in some string deriving from α

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

Some number of derivations gets us x at the beginning

Goal \rightarrow SheepNoise

SheepNoise \rightarrow SheepNoise baa

| baa

For SheepNoise:

FIRST(Goal) = { baa }

FIRST(SN) = { baa }

FIRST(baa) = { baa }



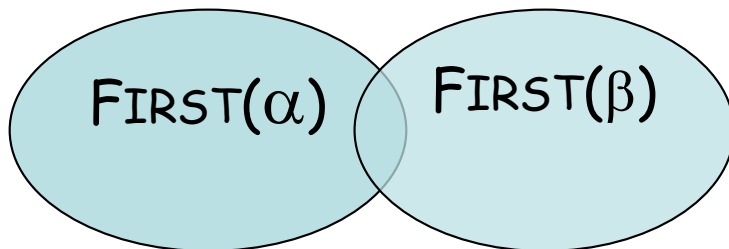
LL(1) Property

If $A \rightarrow \alpha$ 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 exactly one symbol!

Almost correct! See the next slide



Does not have LL(1) Property



What about ε -productions?

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\varepsilon \in \text{FIRST}(\alpha)$, then we need to ensure

$$\text{FOLLOW}(A) \cap \text{FIRST}(\beta) = \emptyset$$

where,

$\text{FOLLOW}(A)$ = the set of terminal symbols that can immediately follow A in a sentential form

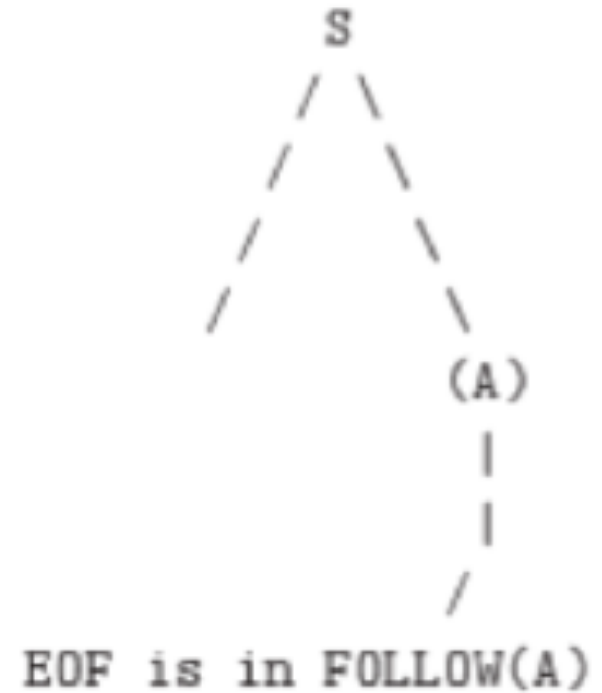
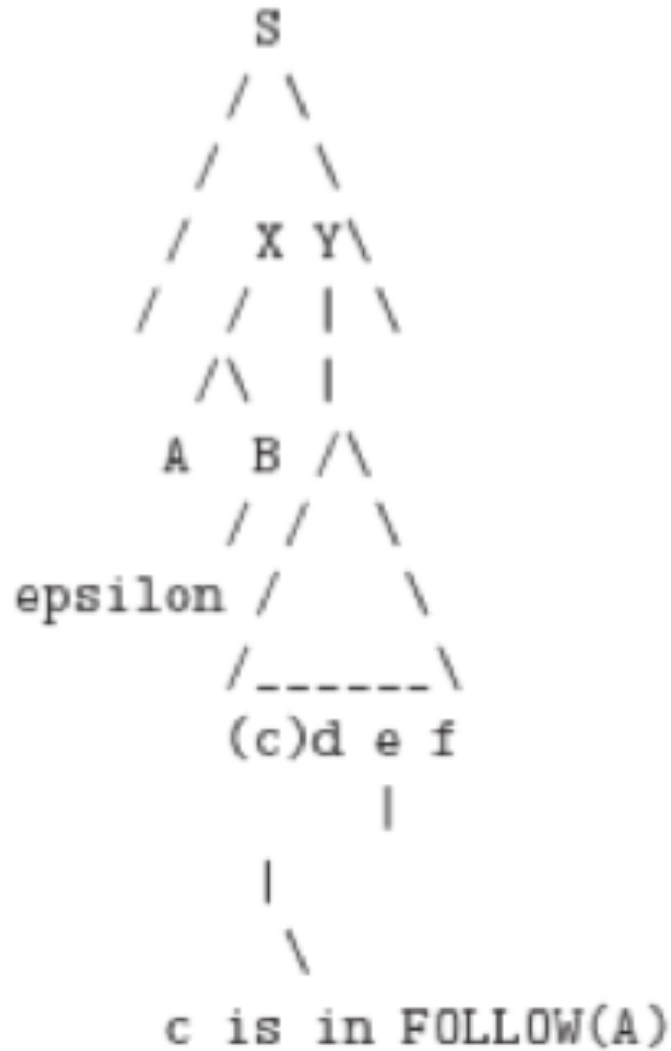
Formally,

$\text{Follow}(A) = \{t \mid (t \text{ is a terminal and } G \Rightarrow^* \alpha A \underline{t} \beta) \text{ or } (t \text{ is eof and } G \Rightarrow^* \alpha A)\}$

Note: eof if A is at the end of the derived sentence



Follow Sets Intuition



FIRST⁺ sets

Definition of FIRST⁺(A → α)

if $\varepsilon \in \text{FIRST}(\alpha)$ then

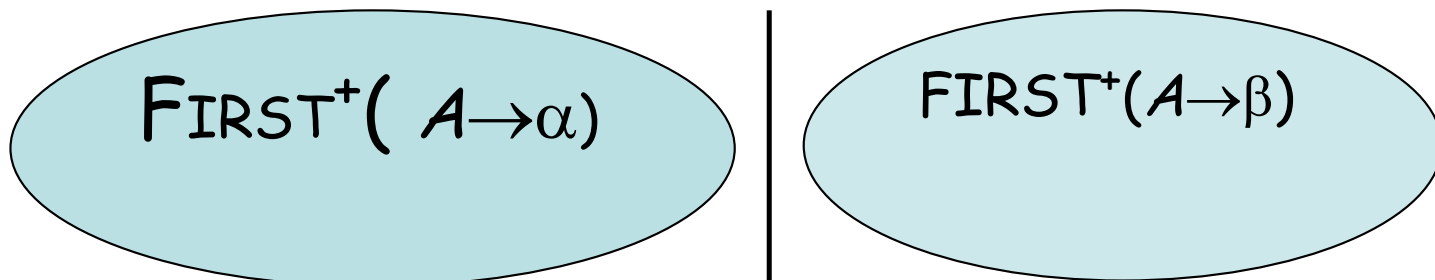
$$\text{FIRST}^+(A \rightarrow \alpha) = \text{FIRST}(\alpha) \cup \text{FOLLOW}(A)$$

else

$$\text{FIRST}^+(A \rightarrow \alpha) = \text{FIRST}(\alpha)$$

Grammar is *LL(1)* iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies

$$\text{FIRST}^+(A \rightarrow \alpha) \cap \text{FIRST}^+(A \rightarrow \beta) = \emptyset$$





What If My Grammar Is Not LL(1) ?

Can we transform a non-LL(1) grammar into an LL(1) grammar?

- In general, the answer is no
- In some cases, however, the answer is yes
- Perform:
 - Eliminate left-recursion **Previously**
 - Perform left factoring **today**



What If My Grammar Is Not LL(1) ?

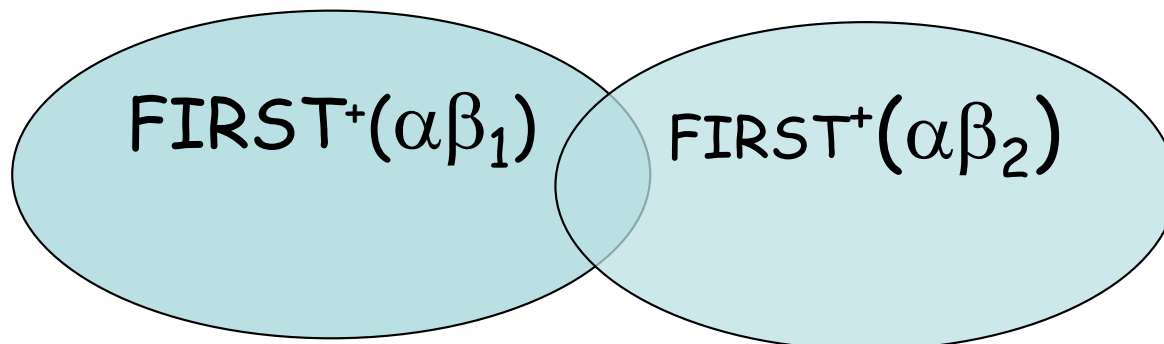
Given grammar G with productions

$$A \rightarrow \alpha \beta_1$$

$$A \rightarrow \alpha \beta_2$$

if α derives anything other than ε and

$$\text{FIRST}^+(A \rightarrow \alpha \beta_1) \cap \text{FIRST}^+(A \rightarrow \alpha \beta_2) \neq \emptyset$$



This grammar is not LL(1)



Left Factoring

If we pull the common prefix, α , into a separate production, we may make the grammar LL(1).

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1$$

$$| \beta_2$$

Create a new Nonterminal

Now, if $\text{FIRST}^+(A' \rightarrow \beta_1) \cap \text{FIRST}^+(A' \rightarrow \beta_2) = \emptyset$,
 G may be LL(1)

Left Factoring

For each nonterminal A
find the longest prefix α common to 2 or more
alternatives for A

if $\alpha \neq \epsilon$ then

replace all of the productions

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3 \mid \dots \mid \alpha \beta_n \mid \gamma$$

with

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots \mid \beta_n$$

Repeat until no NT has rhs' with a common prefix

NT with common prefix

Left Factoring

For each nonterminal A
find the longest prefix α common to 2 or more
alternatives for A

if $\alpha \neq \epsilon$ then

replace all of the productions

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3 \mid \dots \mid \alpha \beta_n \mid \gamma$$

with

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots \mid \beta_n$$

Repeat until no NT has rhs' with a common prefix

Put common prefix α into a
separate production rule

Left Factoring

For each nonterminal A
find the longest prefix α common to 2 or more
alternatives for A

if $\alpha \neq \varepsilon$ then

replace all of the productions

$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3 \mid \dots \mid \alpha \beta_n \mid \gamma$

with

$A \rightarrow \alpha A' \mid \gamma$

$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots \mid \beta_n$

Repeat until no NT has rhs' with a common prefix

Create new Nonterminal (A')
with all unique suffixes



Left Factoring

For each nonterminal A
find the longest prefix α common to 2 or more
alternatives for A
if $\alpha \neq \varepsilon$ then
replace all of the productions
 $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3 \mid \dots \mid \alpha \beta_n \mid \gamma$
with
 $A \rightarrow \alpha A' \mid \gamma$
 $A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots \mid \beta_n$
Repeat until no NT has rhs' with a common prefix

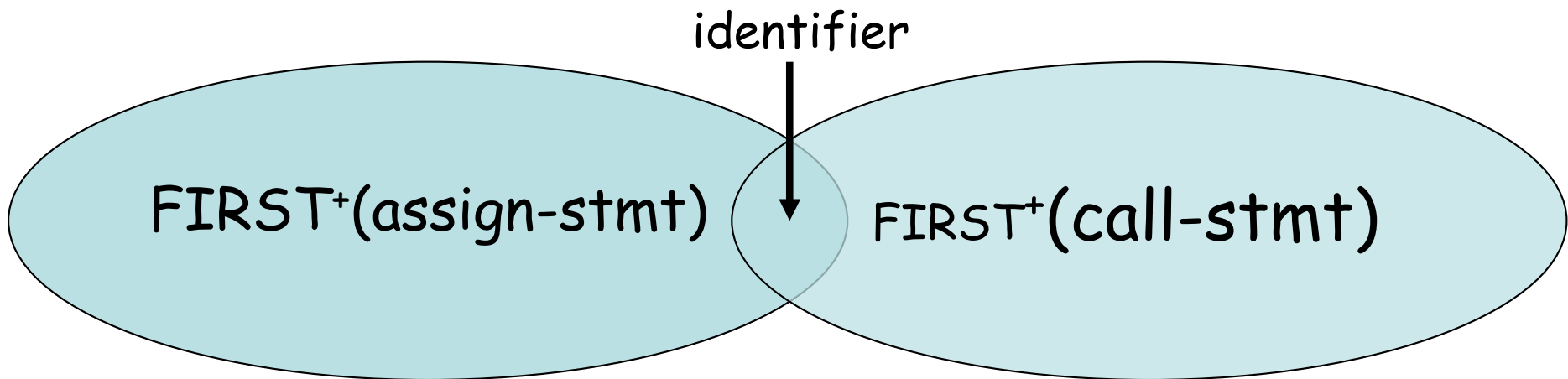
Transformation makes some grammars into LL(1) grammars
There are languages for which no LL(1) grammar exists 15



Left Factoring not possible

Here is an example where a programming language fails to be LL(1) and is not in a form that can be left factored

statement \rightarrow *assign-stmt* | *call-stmt* | **other**
assign-stmt \rightarrow **identifier** := *exp*
call-stmt \rightarrow **identifier** (*exp-list*)





Left Factoring Example

Consider a simple right-recursive expression grammar

0		<i>Goal</i>	→	<i>Expr</i>
1		<i>Expr</i>	→	<i>Term + Expr</i>
2				<i>Term - Expr</i>
3				<i>Term</i>
4		<i>Term</i>	→	<i>Factor * Term</i>
5				<i>Factor / Term</i>
6				<i>Factor</i>
7		<i>Factor</i>	→	<u>number</u>
8				<u>id</u>

To choose between 1, 2, & 3, an LL(1) parser must look past the number or id to see the operator.

$FIRST^+(1) = FIRST^+(2) = FIRST^+(3)$

and

$FIRST^+(4) = FIRST^+(5) = FIRST^+(6)$

Let's left factor this grammar.



Left Factoring Example

After Left Factoring, we have

0		<i>Goal</i>	→	<i>Expr</i>
1		<i>Expr</i>	→	<i>Term Expr'</i>
2		<i>Expr'</i>	→	+ <i>Expr</i>
3				- <i>Expr</i>
4				ϵ
5		<i>Term</i>	→	<i>Factor Term'</i>
6		<i>Term'</i>	→	* <i>Term</i>
7				/ <i>Term</i>
8				ϵ
9		<i>Factor</i>	→	<u>number</u>
10				<u>id</u>

Clearly,

$FIRST^+(2)$, $FIRST^+(3)$, & $FIRST^+(4)$

are disjoint, as are

$FIRST^+(6)$, $FIRST^+(7)$, & $FIRST^+(8)$

The grammar now has the LL(1) property



FIRST Sets

FIRST(α)

For some $\alpha \in (T \cup NT)^*$, define **FIRST(α)** as the set of tokens that appear as the first symbol in some string that derives from α

That is, $\underline{x} \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* \underline{x} \gamma$, for some γ

Computing FIRST Sets

```
for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop
```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Set terminals

Computing FIRST Sets

```

for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop

```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Set empty set for First of nonterminals



Computing FIRST Sets

```
for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop
```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Fixed point algorithm; Monotone because we always add to First sets; never delete from sets

Computing FIRST Sets

```
for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop
```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Iterate through each production

Computing FIRST Sets

```

for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop

```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

RHS is some set of T and NT .

Computing FIRST Sets

```
for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop
```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Initialize rhs to First of first symbol minus epsilon



Computing FIRST Sets

```
for each  $x \in T$ ,  $FIRST(x) \leftarrow \{x\}$ 
for each  $A \in NT$ ,  $FIRST(A) \leftarrow \emptyset$ 
while (FIRST sets are still changing) do
  for each  $p \in P$ , of the form  $A \rightarrow \beta$  do
    if  $\beta$  is  $B_1 B_2 \dots B_k$  then begin:
       $FS \leftarrow FIRST(B_1) - \{\epsilon\}$ 
      for  $i \leftarrow 1$  to  $k-1$  by 1 while  $\epsilon \in FIRST(B_i)$  do
         $FS \leftarrow FS \cup (FIRST(B_{i+1}) - \{\epsilon\})$ 
      end // for loop
    end // if-then
    if  $i = k$  and  $\epsilon \in FIRST(B_k)$ 
      then  $FS \leftarrow FS \cup \{\epsilon\}$ 
     $FIRST(A) \leftarrow FIRST(A) \cup FS$ 
  end // for loop
end // while loop
```

Outer loop is monotone increasing for FIRST sets

→ $|T \cup NT \cup \epsilon|$ is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Iterate through symbols in production until have a symbol that does not have epsilon in First set



Expression Grammar

0	<i>Goal</i>	→	<i>Expr</i>
1	<i>Expr</i>	→	<i>Term Expr'</i>
2	<i>Expr'</i>	→	+ <i>Term Expr'</i>
3			- <i>Term Expr'</i>
4			ϵ
5	<i>Term</i>	→	<i>Factor Term'</i>
6	<i>Term'</i>	→	* <i>Factor Term'</i>
7			/ <i>Factor Term'</i>
8			ϵ
9	<i>Factor</i>	→	<u>number</u>
10			<u>id</u>
11			(<i>Expr</i>)

Symbol	FIRST
<u>num</u>	<u>num</u>
<u>id</u>	<u>id</u>
+	+
-	-
*	*
/	/
((
))
<u>eof</u>	<u>eof</u>
ϵ	ϵ
<i>Goal</i>	<u>num, id, (</u>
<i>Expr</i>	<u>num, id, (</u>
<i>Expr'</i>	+, -, ϵ
<i>Term</i>	<u>num, id, (</u>
<i>Term'</i>	*, /, ϵ
<i>Factor</i>	<u>num, id, (</u>