

# Building a Predictive Parser

I.e., How to build the parse table for a  
recursive-descent parser

# Last Time: Intro LL(1) Predictive Parser

*Predict* the parse tree  
top-down

Parser structure

- 1 token of lookahead
- A stack tracking the current parse tree's frontier
- Selector/parse table

Necessary conditions

- Left-factored
- Free of left-recursion



# Today: Building the Parse Table

Review grammar transformations

- Why they are necessary
- How they work

Build the parse table

- $\text{FIRST}(X)$ : Set of terminals that can begin at a subtree rooted at  $X$
- $\text{FOLLOW}(X)$ : Set of terminals that can appear after  $X$

# Review of LL(1) Grammar Transformations

Necessary (but not sufficient conditions) for LL(1) parsing:

- Free of left recursion
  - “No left-recursive rules”
  - Why? Need to look past the list to know when to cap it
- Left-factored
  - “No rules with a common prefix, for any nonterminal”
  - Why? We would need to look past the prefix to pick the production

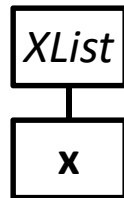
# Why Left Recursion is a Problem (Blackbox View)

CFG snippet:  $XList \rightarrow XList\ x \mid x$

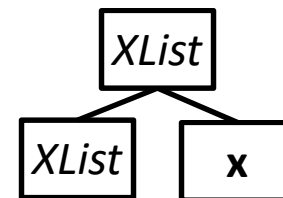
Current parse tree: *XList*

Current token: **x**

How should we grow the tree top-down?



**(OR)**



Correct if there are no more **xs**

Correct if there are more **xs**

**We don't know which to choose without more lookahead**

# Why Left Recursion is a Problem (Whitebox View)

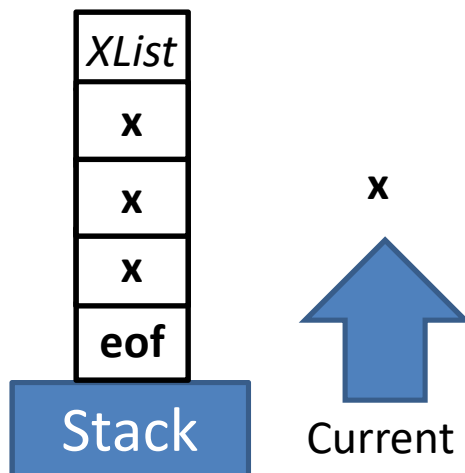
CFG snippet:  $XList \rightarrow XList \ x \mid x$

Current parse tree:  $XList$

Current token:  $x$

Parse table:

$XList$	$x$	$\epsilon$
---------	-----	------------



# Left-Recursion Elimination: Review

Replace  $A \rightarrow A \alpha \mid \beta$   
With  $A \rightarrow \beta A'$   
 $A' \rightarrow \alpha A' \mid \varepsilon$

Head of the list

Where  $\beta$  does not start with  $A$ , *or* may not be present

Preserves the language (a list of  $\alpha$ s, starting with a  $\beta$ ),  
but uses right recursion

# Left-Recursion Elimination: Ex1

$$A \rightarrow A \alpha \mid \beta \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}$$

$$\begin{array}{l} E \rightarrow E \text{ cross id} \mid \text{id} \\ \quad \underbrace{\hspace{1.5cm}}_{\alpha} \quad \underbrace{\hspace{1cm}}_{\beta} \end{array} \quad \Rightarrow \quad \begin{array}{l} E \rightarrow \text{id } E' \\ E' \rightarrow \underbrace{\text{cross id}}_{\alpha} E' \mid \varepsilon \end{array}$$



# Left-Recursion Elimination: Ex2

$$A \rightarrow A \alpha \mid \beta \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow ( E ) \mid \text{id} \end{array} \quad \Rightarrow \quad \begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \mid \epsilon \\ T \rightarrow F T' \\ T' \rightarrow * F T' \mid \epsilon \\ F \rightarrow ( E ) \mid \text{id} \end{array}$$

# Left-Recursion Elimination: Ex3

$$A \rightarrow A \alpha \mid \beta \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

$DList \rightarrow DList D \mid \epsilon$   
 $D \rightarrow Type \text{ id semi}$   
 $Type \rightarrow \text{bool} \mid \text{int}$

$DList \rightarrow \epsilon DList'$   
 $DList' \rightarrow D DList' \mid \epsilon$   
 $D \rightarrow Type \text{ id semi}$   
 $Type \rightarrow \text{bool} \mid \text{int}$

$DList \rightarrow D DList \mid \epsilon$   
 $D \rightarrow Type \text{ id semi}$   
 $Type \rightarrow \text{bool} \mid \text{int}$

# Left Factoring: Review

Removing a common prefix from a grammar

Replace  $A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_m \mid \gamma_1 \mid \dots \mid \gamma_n$

With  $A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_n$   
 $A' \rightarrow \beta_1 \mid \dots \mid \beta_m$

Where  $\beta_i$  and  $\gamma_i$  are sequence of symbols with no common prefix

Note:  $\gamma_i$  may not be present, and one of the  $\beta$  may be  $\epsilon$

Combine all “problematic” rules that start with  $\alpha$  into one rule  $\alpha A'$   
Now  $A'$  represents the suffix of the “problematic” rules

# Left Factoring: Example 1

$$A \rightarrow \alpha \beta_1 \mid \dots \mid \alpha \beta_m \mid \gamma_1 \mid \dots \mid \gamma_n \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_n \\ A' \rightarrow \beta_1 \mid \dots \mid \beta_m \end{array}$$

$$X \rightarrow \overset{\alpha}{\underbrace{\hspace{1cm}}} \overset{\beta_1}{\underbrace{\hspace{1cm}}} \mid \overset{\alpha}{\underbrace{\hspace{1cm}}} \overset{\beta_2}{\underbrace{\hspace{1cm}}} \mid \overset{\alpha}{\underbrace{\hspace{1cm}}} \overset{\beta_3}{\underbrace{\hspace{1cm}}} \mid \overset{\gamma_1}{\underbrace{\hspace{1cm}}} d$$

$X \rightarrow \langle a \rangle \mid \langle b \rangle \mid \langle c \rangle \mid d$

---


$$X \rightarrow \overset{\alpha}{\underbrace{\hspace{1cm}}} X' \mid \overset{\gamma_1}{\underbrace{\hspace{1cm}}} d$$

$$X' \rightarrow \underbrace{a}_{\beta_1} \mid \underbrace{b}_{\beta_2} \mid \underbrace{c}_{\beta_3}$$

# Left Factoring: Example 2

$$A \rightarrow \alpha \beta_1 \mid \dots \mid \alpha \beta_m \mid \gamma_1 \mid \dots \mid \gamma_n \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_n \\ A' \rightarrow \beta_1 \mid \dots \mid \beta_m \end{array}$$

$\beta_1$   $\beta_2$

$Stmt \rightarrow \text{id assign } E \mid \text{id } ( EList ) \mid \text{return}$

$E \rightarrow \text{intlit} \mid \text{id}$

$EList \rightarrow E \mid E \text{ comma } EList$

---

$Stmt \rightarrow \text{id } Stmt' \mid \text{return}$

$Stmt' \rightarrow \text{assign } E \mid ( EList )$

$E \rightarrow \text{intlit} \mid \text{id}$

$EList \rightarrow E \mid E \text{ comma } EList$

# Left Factoring: Example 3

$$A \rightarrow \alpha \beta_1 \mid \dots \mid \alpha \beta_m \mid \gamma_1 \mid \dots \mid \gamma_n \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_n \\ A' \rightarrow \beta_1 \mid \dots \mid \beta_m \end{array}$$

$\alpha$        $\beta_1 = \epsilon$        $\alpha$        $\beta_2$   
 $S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{semi}$   
 $E \rightarrow \text{boollit}$

---

$S \rightarrow \text{if } E \text{ then } S S' \mid \text{semi}$

$S' \rightarrow \text{else } S \mid \epsilon$

$E \rightarrow \text{boollit}$

# Left Factoring: Not Always Immediate

$$A \rightarrow \alpha \beta_1 \mid \dots \mid \alpha \beta_m \mid \gamma_1 \mid \dots \mid \gamma_n \quad \Rightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_n \\ A' \rightarrow \beta_1 \mid \dots \mid \beta_m \end{array}$$

This snippet yearns for left factoring

```
S → A | C | return
A → id assign E
C → id ( Elist )
```

but we cannot! At least without *inlining*

```
S → id assign E | id ( Elist ) | return
```

# Let's be more constructive

So far, we have only talked about what precludes us from building a predictive parser

It is time to actually build the parse table



# Building the Parse Table

What do we actually need to ensure that production  $A \rightarrow \alpha$  is the correct one to apply?

Assume  $\alpha$  is an arbitrary sequence of symbols

1. What terminals could  $\alpha$  possibly start with  
→ we call this the FIRST set
2. What terminal could possibly come after  $A$   
→ we call this the FOLLOW set

# Why is FIRST Important?

Assume the top-of-stack symbol is  $A$  and current token is  $a$

- Production 1:  $A \rightarrow \alpha$
- Production 2:  $A \rightarrow \beta$

FIRST lets us disambiguate:

- If  $a$  is in  $\text{FIRST}(\alpha)$ , we know Production 1 is a viable choice
- If  $a$  is in  $\text{FIRST}(\beta)$ , we know Production 2 is a viable choice
- If  $a$  is only in one of  $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$ , we can predict the production we need

# FIRST Sets

$\text{FIRST}(\alpha)$  is the set of terminals that begin the strings derivable from  $\alpha$ , and also, if  $\alpha$  can derive  $\varepsilon$ , then  $\varepsilon$  is in  $\text{FIRST}(\alpha)$ .

Formally, let's write it together

$\text{FIRST}(\alpha) =$

# FIRST Sets

$\text{FIRST}(\alpha)$  is the set of terminals that begin the strings derivable from  $\alpha$ , and also, if  $\alpha$  can derive  $\epsilon$ , then  $\epsilon$  is in  $\text{FIRST}(\alpha)$ .

Formally, let's write it together

$$\text{FIRST}(\alpha) = \{t \mid (t \in \Sigma \wedge \alpha \Rightarrow^* t\beta) \vee (t = \epsilon \wedge \alpha \Rightarrow^* \epsilon)\}$$

# FIRST Construction: Single Symbol

We begin by doing FIRST sets for a single, arbitrary symbol  $X$

- If  $X$  is a terminal:  $\text{FIRST}(X) = \{ X \}$
- If  $X$  is  $\epsilon$ :  $\text{FIRST}(\epsilon) = \{ \epsilon \}$
- If  $X$  is a nonterminal, for each  $X \rightarrow Y_1 Y_2 \dots Y_k$ 
  - Put  $\text{FIRST}(Y_1) - \{\epsilon\}$  into  $\text{FIRST}(X)$
  - If  $\epsilon$  is in  $\text{FIRST}(Y_1)$ , put  $\text{FIRST}(Y_2) - \{\epsilon\}$  into  $\text{FIRST}(X)$
  - If  $\epsilon$  is also in  $\text{FIRST}(Y_2)$ , put  $\text{FIRST}(Y_3) - \{\epsilon\}$  into  $\text{FIRST}(X)$
  - ...
  - If  $\epsilon$  is in FIRST of all  $Y_i$  symbols, put  $\epsilon$  into  $\text{FIRST}(X)$

Repeat this step until there are no changes to any nonterminal's FIRST set

# FIRST(X) Example

## Building FIRST(X) for nonterm X

for each  $X \rightarrow Y_1 Y_2 \dots Y_k$

- Add  $\text{FIRST}(Y_1) - \{\epsilon\}$
- If  $\epsilon$  is in  $\text{FIRST}(Y_{1 \text{ to } i-1})$ : add  $\text{FIRST}(Y_i) - \{\epsilon\}$
- If  $\epsilon$  is in all RHS symbols, add  $\epsilon$

$\text{Exp} \rightarrow \text{Term Exp}'$

$\text{Exp}' \rightarrow \text{minus Term Exp}' \mid \epsilon$

$\text{Term} \rightarrow \text{Factor Term}'$

$\text{Term}' \rightarrow \text{divide Factor Term}' \mid \epsilon$

$\text{Factor} \rightarrow \text{intlit} \mid \text{lparen Exp rparen}$

$\text{FIRST}(\text{Factor}) = \{ \text{intlit}, \text{lparen} \}$

$\text{FIRST}(\text{Term}') = \{ \text{divide}, \epsilon \}$

$\text{FIRST}(\text{Term}) = \{ \text{intlit}, \text{lparen} \}$

$\text{FIRST}(\text{Exp}') = \{ \text{minus}, \epsilon \}$

$\text{FIRST}(\text{Exp}) = \{ \text{intlit}, \text{lparen} \}$

# FIRST( $\alpha$ )

We now extend FIRST to strings of symbols  $\alpha$

- We want to define FIRST for all RHS

Looks very similar to the procedure for single symbols

Let  $\alpha = Y_1 Y_2 \dots Y_k$

- Put  $\text{FIRST}(Y_1) - \{\epsilon\}$  in  $\text{FIRST}(\alpha)$ 
  - If  $\epsilon$  is in  $\text{FIRST}(Y_1)$ : add  $\text{FIRST}(Y_2) - \{\epsilon\}$  to  $\text{FIRST}(\alpha)$
  - If  $\epsilon$  is in  $\text{FIRST}(Y_2)$ : add  $\text{FIRST}(Y_3) - \{\epsilon\}$  to  $\text{FIRST}(\alpha)$
  - ...
  - If  $\epsilon$  is in FIRST of all  $Y_i$  symbols, put  $\epsilon$  into  $\text{FIRST}(\alpha)$

# Building $\text{FIRST}(\alpha)$ from $\text{FIRST}(X)$

## Building $\text{FIRST}(X)$ for nonterm $X$

for each  $X \rightarrow Y_1 Y_2 \dots Y_k$

- Add  $\text{FIRST}(Y_1) - \{\epsilon\}$
- If  $\epsilon$  is in  $\text{FIRST}(Y_{1 \text{ to } i-1})$ : add  $\text{FIRST}(Y_i) - \{\epsilon\}$
- If  $\epsilon$  is in all RHS symbols, add  $\epsilon$

## Building $\text{FIRST}(\alpha)$

Let  $\alpha = Y_1 Y_2 \dots Y_k$

- Add  $\text{FIRST}(Y_1) - \{\epsilon\}$
- If  $\epsilon$  is in  $\text{FIRST}(Y_{1 \text{ to } i-1})$ : add  $\text{FIRST}(Y_i) - \{\epsilon\}$
- If  $\epsilon$  is in all RHS symbols, add  $\epsilon$



# FIRST( $\alpha$ ) Example

## Building FIRST( $\alpha$ )

Let  $\alpha = Y_1 Y_2 \dots Y_k$

- Add FIRST( $Y_1$ ) -  $\{\epsilon\}$
- If  $\epsilon$  is in FIRST( $Y_{1 \text{ to } i-1}$ ): add FIRST( $Y_i$ ) -  $\{\epsilon\}$
- If, for all RHS symbols  $Y_j$ ,  $\epsilon$  is in FIRST( $Y_j$ ), add  $\epsilon$

$E \rightarrow T X$

$X \rightarrow + T X \mid \epsilon$

$T \rightarrow F Y$

$Y \rightarrow * F Y \mid \epsilon$

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

$\text{FIRST}(E) = \{ (, \text{id} \}$

$\text{FIRST}(T) = \{ (, \text{id} \}$

$\text{FIRST}(F) = \{ (, \text{id} \}$

$\text{FIRST}(X) = \{ +, \epsilon \}$

$\text{FIRST}(Y) = \{ *, \epsilon \}$

$\text{FIRST}(T X) = \{ (, \text{id} \}$

$\text{FIRST}(+ T X) = \{ + \}$

$\text{FIRST}(F Y) = \{ (, \text{id} \}$

$\text{FIRST}( * F Y ) = \{ * \}$

$\text{FIRST}( ( E ) ) = \{ ( \}$

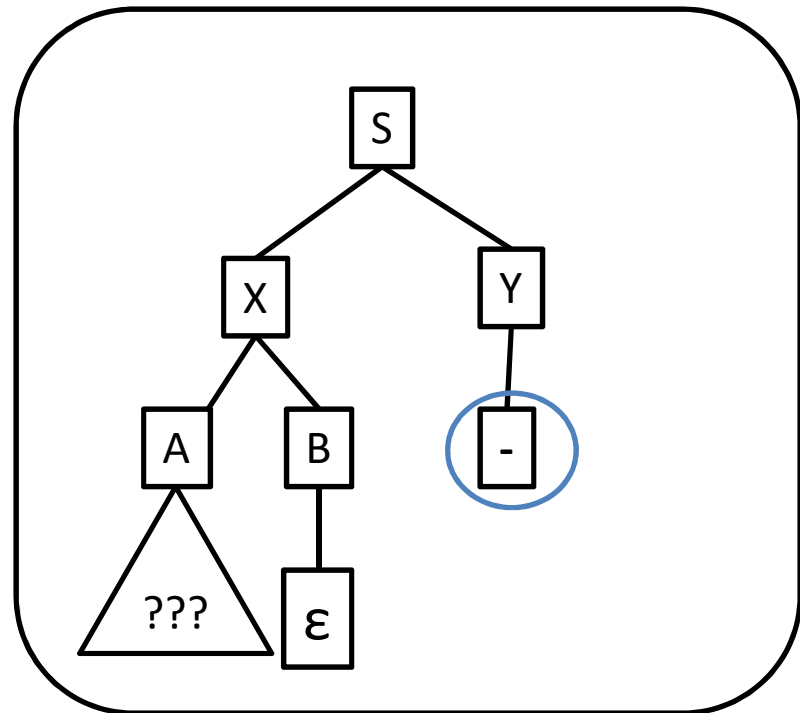
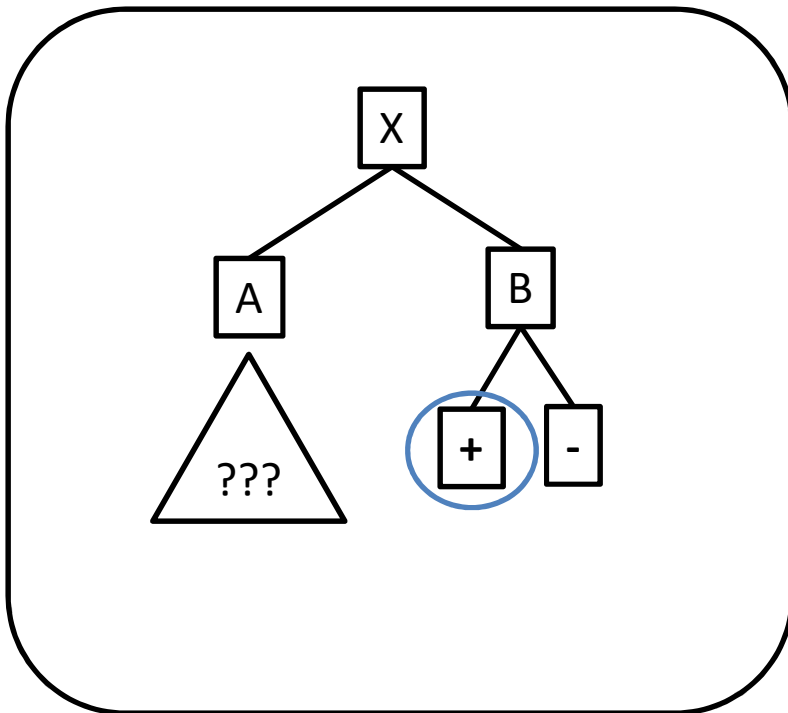
$\text{FIRST}( \text{id} ) = \{ \text{id} \}$

FIRST sets alone do not provide enough information to construct a parse table

If a rule  $R$  can derive  $\varepsilon$ , we need to know what terminals can come just after  $R$

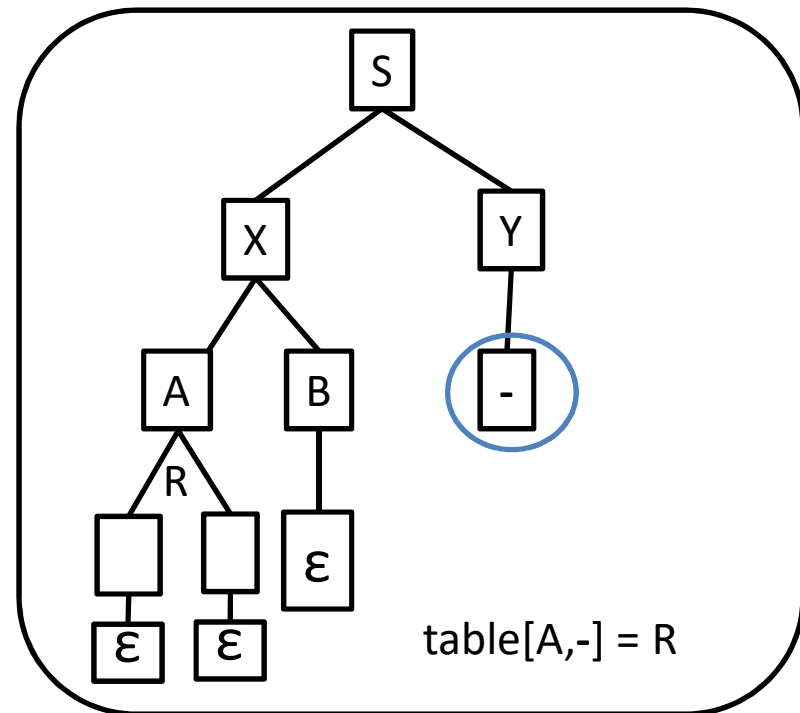
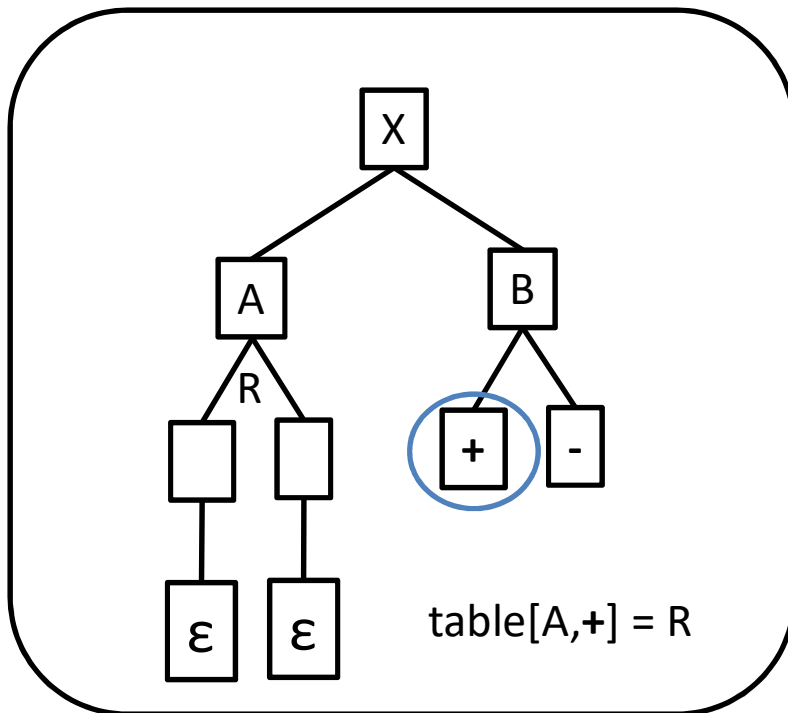
# FOLLOW Sets: Pictorially

For nonterminal  $A$ ,  $\text{FOLLOW}(A)$  is the set of terminals that can appear immediately to the right of  $A$



# FOLLOW Sets: Pictorially

For nonterminal  $A$ ,  $\text{FOLLOW}(A)$  is the set of terminals that can appear immediately to the right of  $A$



# FOLLOW Sets

For nonterminal  $A$ ,  $\text{FOLLOW}(A)$  is the set of terminals that can appear immediately to the right of  $A$

Let's write it together,

$\text{FOLLOW}(A) =$

# FOLLOW Sets

For nonterminal  $A$ ,  $\text{FOLLOW}(A)$  is the set of terminals that can appear immediately to the right of  $A$

Let's write it together,

$\text{FOLLOW}(A) =$

$$\{t \mid (t \in \Sigma \wedge S \Rightarrow^+ \alpha A t \beta) \vee (t = EOF \wedge S \Rightarrow^* \alpha A)\}$$

# FOLLOW Sets: Construction

To build FOLLOW(A)

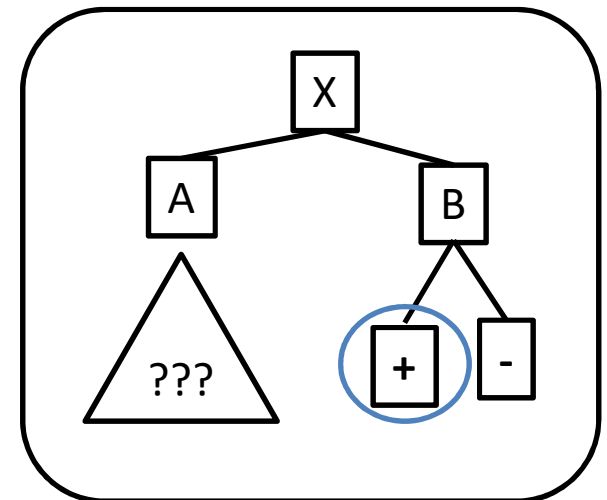
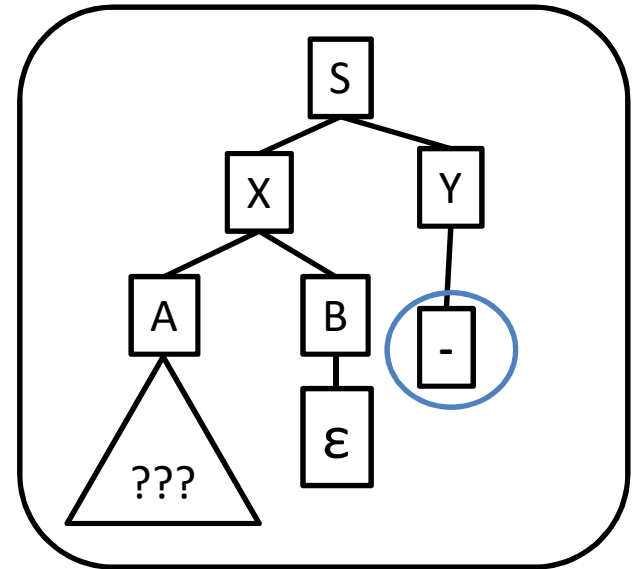
– If A is the start nonterminal, add eof

Where  $\alpha, \beta$  may be empty

– For rules  $X \rightarrow \alpha A \beta$

- Add  $\text{FIRST}(\beta) - \{\epsilon\}$
- If  $\epsilon$  is in  $\text{FIRST}(\beta)$  or  $\beta$  is empty, add  $\text{FOLLOW}(X)$

Continue building FOLLOW sets until reach a fixed point (i.e., no more symbols can be added)



# FOLLOW Sets Example

FOLLOW(A) for  $X \rightarrow \alpha A \beta$

If A is the start, add **eof**

Add FIRST( $\beta$ ) –  $\{\epsilon\}$

Add FOLLOW(X) if  $\epsilon$  in FIRST( $\beta$ ) or  $\beta$  is empty

$S \rightarrow Bc \mid DB$

$B \rightarrow ab \mid cS$

$D \rightarrow d \mid \epsilon$

FIRST (S) = { **a, c, d** }

FIRST (B) = { **a, c** }

FIRST (D) = { **d,  $\epsilon$**  }

FIRST (B c) = { **a, c** }

FIRST (D B) = { **d, a, c** }

FIRST (**a b**) = { **a** }

FIRST (**c S**) = { **c** }

FOLLOW (S) = { **eof** }

FOLLOW (B) = { **c, eof** }

FOLLOW (D) = { **a, c** }

---

FOLLOW (S) = { **eof, c** }

FOLLOW (B) = { **c, eof** }

FOLLOW (D) = { **a, c** }

---

FOLLOW (S) = { **eof, c** }

FOLLOW (B) = { **c, eof** }

FOLLOW (D) = { **a, c** }

---



# Building the Parse Table

```
for each production  $X \rightarrow \alpha$  {  
  for each terminal t in FIRST( $\alpha$ ) {  
    put  $\alpha$  in Table[X][t]  
  }  
  if  $\epsilon$  is in FIRST( $\alpha$ ) {  
    for each terminal t in FOLLOW( $X$ ) {  
      put  $\alpha$  in Table[X][t]  
    }  
  }  
}
```

Table collision  $\Leftrightarrow$  Grammar is not in LL(1)

# Putting it all together

Build FIRST sets for each nonterminal

Build FIRST sets for each production's RHS

Build FOLLOW sets for each nonterminal

Use FIRST and FOLLOW to fill parse table for each production

# Tips n' Tricks

## FIRST sets

- Only contain alphabet terminals and  $\varepsilon$
- Defined for arbitrary RHS and nonterminals
- Constructed by starting at the beginning of a production

## FOLLOW sets

- Only contain alphabet terminals and eof
- Defined for nonterminals only
- Constructed by jumping into production

FIRST( $\alpha$ ) for  $\alpha = Y_1 Y_2 \dots Y_k$

Add FIRST( $Y_1$ ) -  $\{\epsilon\}$

If  $\epsilon$  is in FIRST( $Y_{1 \text{ to } i-1}$ ): add FIRST( $Y_i$ ) -  $\{\epsilon\}$

If  $\epsilon$  is in all RHS symbols, add  $\epsilon$

FOLLOW(A) for  $X \rightarrow \alpha A \beta$

If A is the start, add **eof**

Add FIRST( $\beta$ ) -  $\{\epsilon\}$

Add FOLLOW(X) if  $\epsilon$  in FIRST( $\beta$ ) or  $\beta$  empty

Table[X][t]

for each production  $X \rightarrow \alpha$

for each terminal **t** in FIRST( $\alpha$ )

put  $\alpha$  in Table[X][**t**]

if  $\epsilon$  is in FIRST( $\alpha$ ) {

for each terminal **t** in FOLLOW(X) {

put  $\alpha$  in Table[X][**t**]

FIRST (S) = { **a, c, d** }

FIRST (B) = { **a, c** }

FIRST (D) = { **d,  $\epsilon$**  }

FIRST (B c) = { **a, c** }

FIRST (D B) = { **d, a, c** }

FIRST (a b) = { **a** }

FIRST (c S) = { **c** }

FIRST (d) = { **d** }

FIRST ( $\epsilon$ ) = {  **$\epsilon$**  }

FOLLOW (S) = { **eof, c** }

FOLLOW (B) = { **c, eof** }

FOLLOW (D) = { **a, c** }

CFG

S  $\rightarrow$  B c | D B

B  $\rightarrow$  a b | c S

D  $\rightarrow$  d |  $\epsilon$



	a	b	c	d	eof
S	B c D B		B c D B	D B	
B	a b		c S		
D	$\epsilon$		$\epsilon$	d	

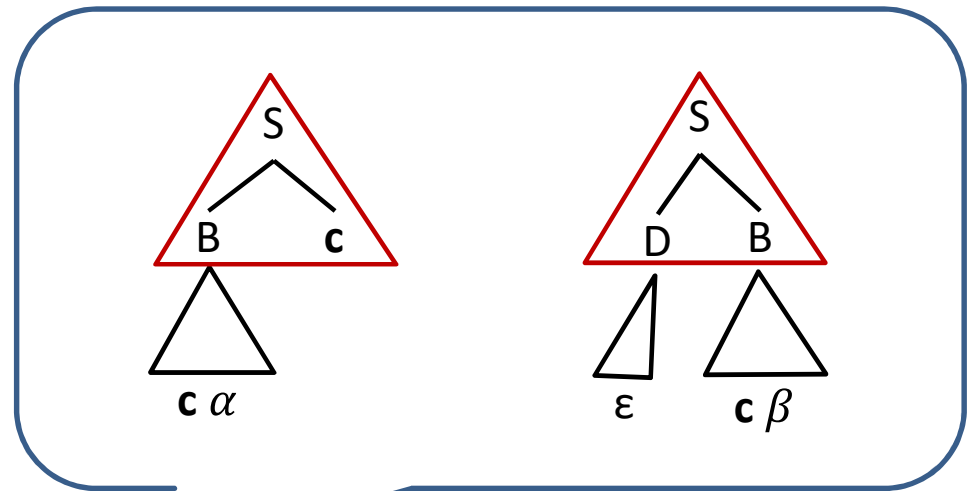
# Why is a Table Collision a Problem?

CFG

$S \rightarrow Bc \mid DB$

$B \rightarrow ab \mid cS$

$D \rightarrow d \mid \epsilon$



	a	b	c	d	$\epsilon$
S	Bc DB		Bc DB	DB	
B	ab		cS		
D	$\epsilon$		$\epsilon$	d	

current  
token

