

Announcements

HW4 due today

HW5 assigned today

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 parse tree 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 selector 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)

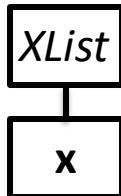
CFG snippet: $XList \rightarrow XList\ x \mid x$

Current parse tree: *XList*

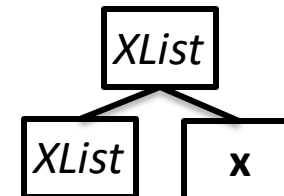
Current token: **x**

How should we grow the tree top-down?

Either prediction would be wrong in some cases.



(OR)



Correct if there are no more **xs**

Correct if there are more **xs**

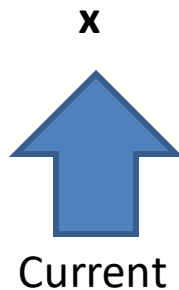
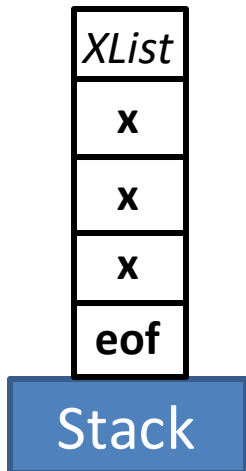
We don't know which without more lookahead

Why Left Recursion is a Problem (Whitebox View)

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

Current parse tree: $XList$ x eof Current token: x

Parse table: $XList$ $XList\ x$ ϵ



(Stack overflow)

Left Recursion Elimination: Review

Turn left recursion into left recursion.

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

Head of the list

Where β does not start with A or may not be present

Preserve order (a list of α starting with β) but use right recursion

Left Recursion Elimination: Ex1

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

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

Left Recursion Elimination: Ex2

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

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

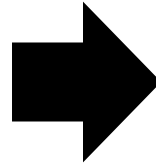
Left Recursion Elimination: Ex3

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

$SList \rightarrow SList D \mid \varepsilon$

$D \rightarrow Type \text{ id semi}$

$Type \rightarrow \text{bool} \mid \text{int}$

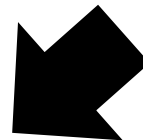


$SList \rightarrow \varepsilon SList'$

$SList' \rightarrow D SList' \mid \varepsilon$

$D \rightarrow Type \text{ id semi}$

$Type \rightarrow \text{bool} \mid \text{int}$



$SList \rightarrow D SList \mid \varepsilon$

$D \rightarrow Type \text{ id semi}$

$Type \rightarrow \text{bool} \mid \text{int}$

Left Factoring: Review

Removing common prefix from 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 β_i and γ_i are sequence of symbols with no common prefix
 γ_i May not be present, one of the β may be ε

Squash all “problem” rules starting with α together into one rule $\alpha A'$
Now A' represents the suffix of the “problem” rules

Left Factoring: Example 1

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

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

$X \rightarrow < a > \mid < b > \mid < c > \mid d$

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

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

Left Factoring: Example 2

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

β_1 β_2

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

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

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

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

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

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

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

Left Factoring: Example 3

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

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

$S \rightarrow$ **if** E **then** $S S'$ **| semi**

$S' \rightarrow$ **else** $S \mid \epsilon$

$E \rightarrow$ **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've only talked about what precludes us from building a predictive parser

It's time to actually build the parse table

Building the Parse Table

What do we actually need to ensure arbitrary production $A \rightarrow \alpha$ is the correct one to apply?
Assume α is an arbitrary sequence of symbols.

1. What **terminals** could α 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 in only in one of $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$, we can predict the production we need

Otherwise the grammar is not LL(1).

FIRST Sets

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

Notations:
N for non-terminals;
 Σ for terminals.

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 α , and also, if α can derive ϵ , then ϵ is in $\text{FIRST}(\alpha)$.

Formally, let's write it together

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

If t is epsilon and alpha can derive it.

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 ϵ : $\text{FIRST}(\epsilon) = \{ \epsilon \}$
- 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 ϵ is in $\text{FIRST}(Y_1)$, put $\text{FIRST}(Y_2) - \{\epsilon\}$ into $\text{FIRST}(X)$
 - If ϵ is also in $\text{FIRST}(Y_2)$, put $\text{FIRST}(Y_3) - \{\epsilon\}$ into $\text{FIRST}(X)$
 - ...
 - If ϵ is in FIRST of all Y_i symbols, put ϵ 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 ϵ is in $\text{FIRST}(Y_{1 \text{ to } i-1})$: add $\text{FIRST}(Y_i) - \{\epsilon\}$
- If ϵ is in all RHS symbols, add ϵ

$Exp \rightarrow Term\ Exp'$

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

$Term \rightarrow Factor\ Term'$

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

$Factor \rightarrow \text{intlitt} \mid \text{lparens}\ Exp\ \text{rparens}$

$\text{FIRST}(Factor) = \{ \text{intlitt}, \text{lparens} \}$

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

$\text{FIRST}(Term) = \{ \text{intlitt}, \text{lparens} \}$

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

$\text{FIRST}(Exp) = \{ \text{intlitt}, \text{lparens} \}$

FIRST(α)

We now extend FIRST to strings of symbols α

- 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 ϵ is in $\text{FIRST}(Y_1)$: add $\text{FIRST}(Y_2) - \{\epsilon\}$ to $\text{FIRST}(\alpha)$
 - If ϵ is in $\text{FIRST}(Y_2)$: add $\text{FIRST}(Y_3) - \{\epsilon\}$ to $\text{FIRST}(\alpha)$
 - ...
 - If ϵ is in FIRST of all Y_i symbols, put ϵ 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 ϵ is in $\text{FIRST}(Y_{1 \text{ to } i-1})$: add $\text{FIRST}(Y_i) - \{\epsilon\}$
- If ϵ is in all RHS symbols, add ϵ

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

Let $\alpha = Y_1 Y_2 \dots Y_k$

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

FIRST(α) Example

Building FIRST(α)

Let $\alpha = Y_1 Y_2 \dots Y_k$

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

$E \rightarrow TX$

$X \rightarrow +TX \mid \epsilon$

$T \rightarrow FY$

$Y \rightarrow *FY \mid \epsilon$

$F \rightarrow (E) \mid id$

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

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

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

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

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

$FIRST(TX) = \{ (, id \}$

$FIRST(+TX) = \{ + \}$

$FIRST(FY) = \{ (, id \}$

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

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

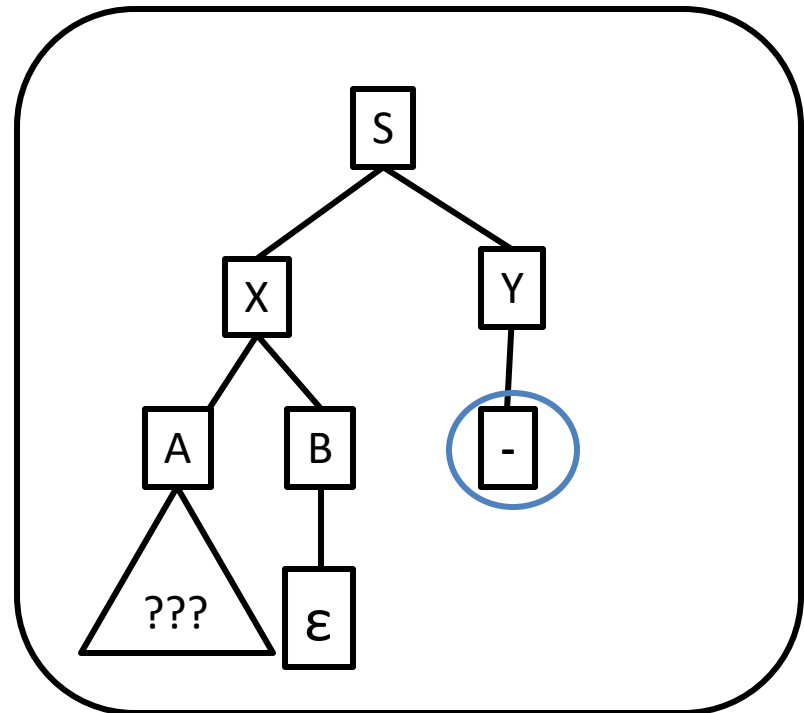
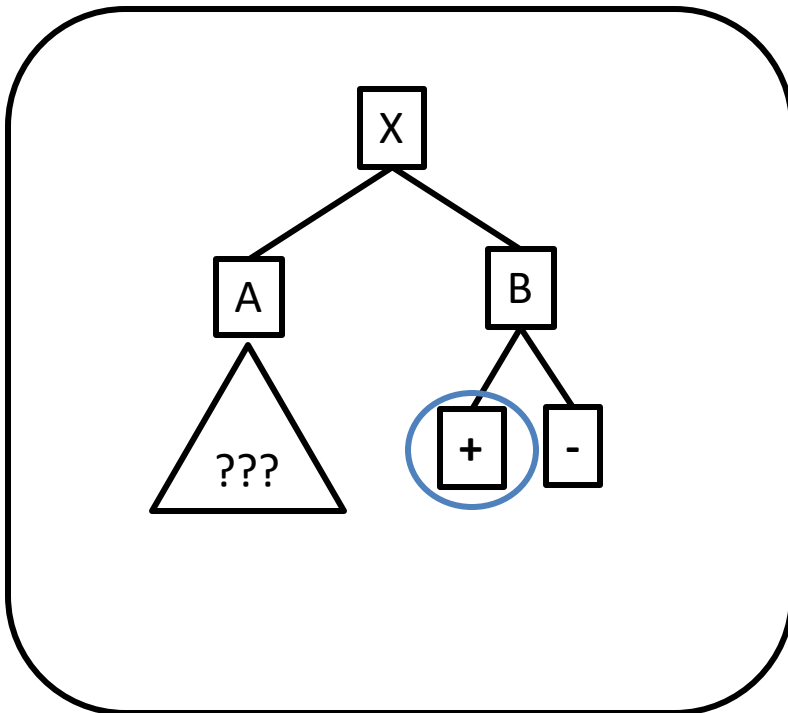
$FIRST(id) = \{ id \}$

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

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

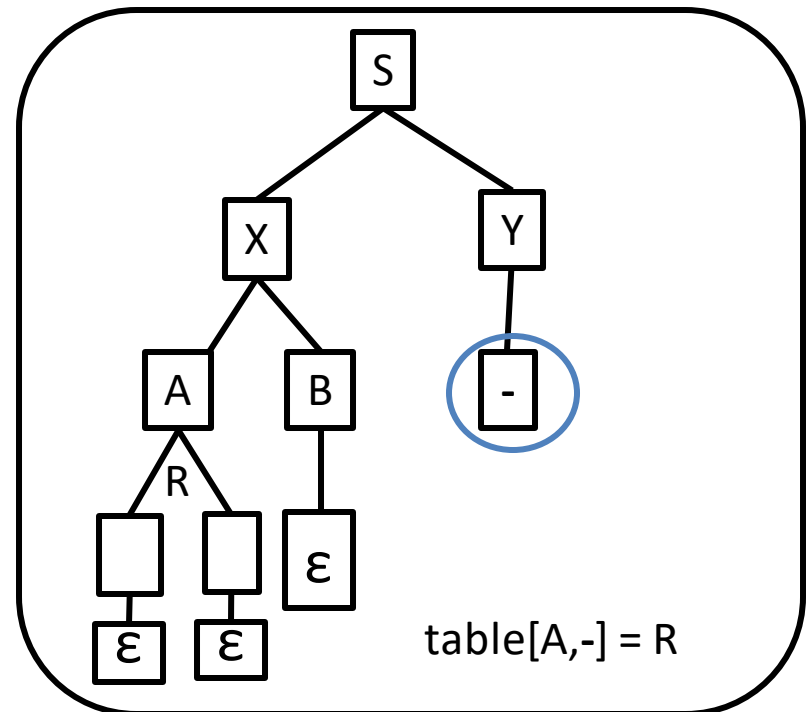
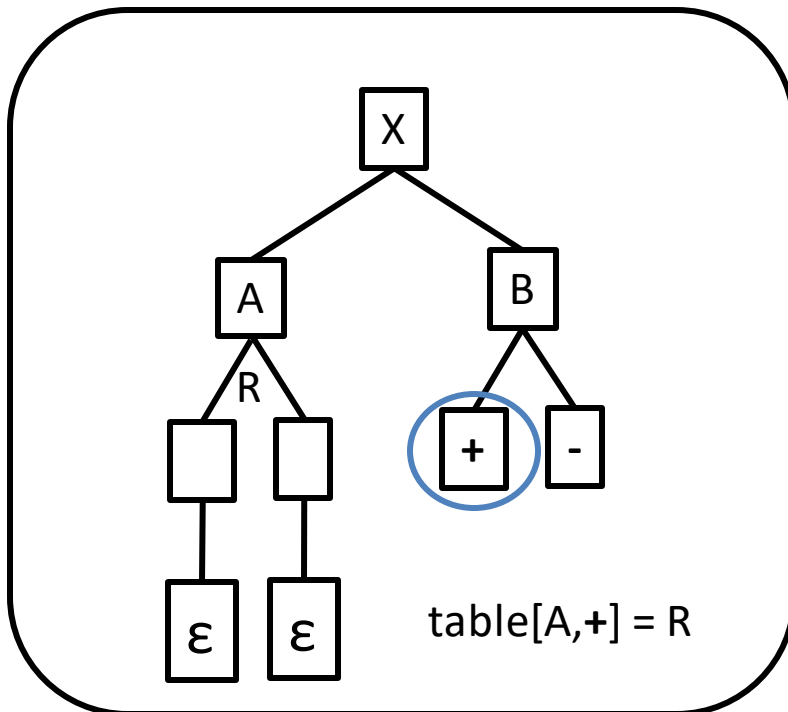
FOLLOW Sets: Pictorially

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,

$\text{FOLLOW}(A) =$

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

eof is always in $\text{FOLLOW}(S)$.

FOLLOW Sets: Construction

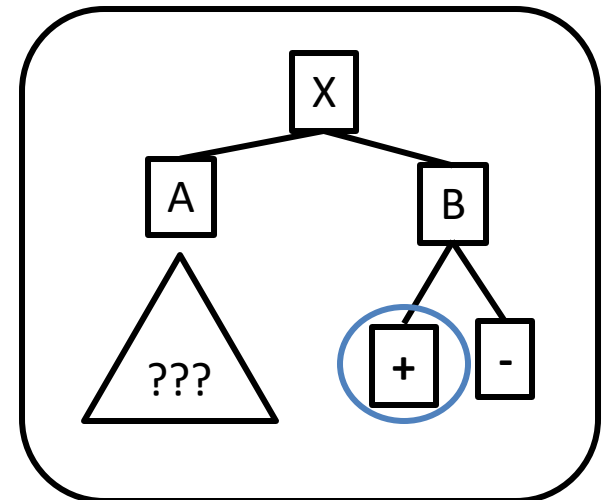
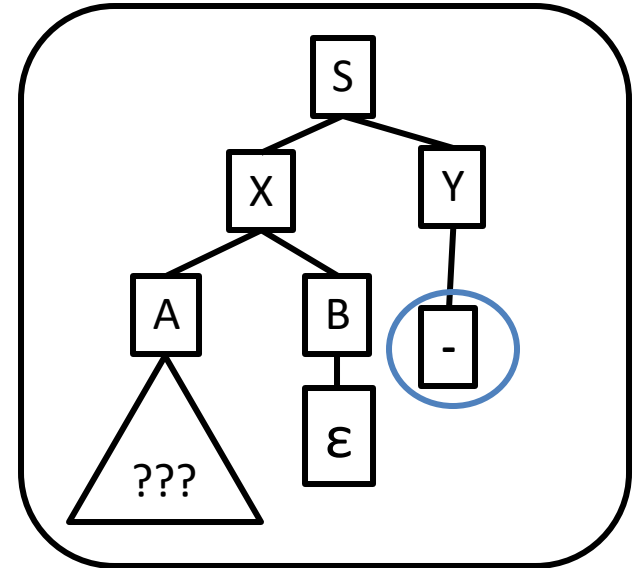
To build FOLLOW(A)

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

Where α, β may be empty

- For rules $X \rightarrow \alpha A \beta$
 - Add $\text{FIRST}(\beta) - \{\epsilon\}$
 - If ϵ is in $\text{FIRST}(\beta)$ or β is empty, add 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(β) – $\{\epsilon\}$

Add FOLLOW(X) if ϵ in FIRST(β) or β 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, ϵ** }

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

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 ε
- 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(α) for $\alpha = Y_1 Y_2 \dots Y_k$

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

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

If ϵ is in all RHS symbols, add ϵ

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

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

Add FIRST(β) - $\{\epsilon\}$

Add FOLLOW(X) if ϵ in FIRST(β) or β empty

Table[X][t]

for each production $X \rightarrow \alpha$

for each terminal **t** in FIRST(α)

put α in Table[X][**t**]

if ϵ is in FIRST(α) {

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

put α in Table[X][**t**]

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

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

FIRST (D) = { **d, ϵ** }

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

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

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

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

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

FIRST (ϵ) = { **ϵ** }

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 | ϵ



	a	b	c	d	eof
S	B c D B		B c D B	D B	
B	a b		c S		
D	ϵ		ϵ	d	

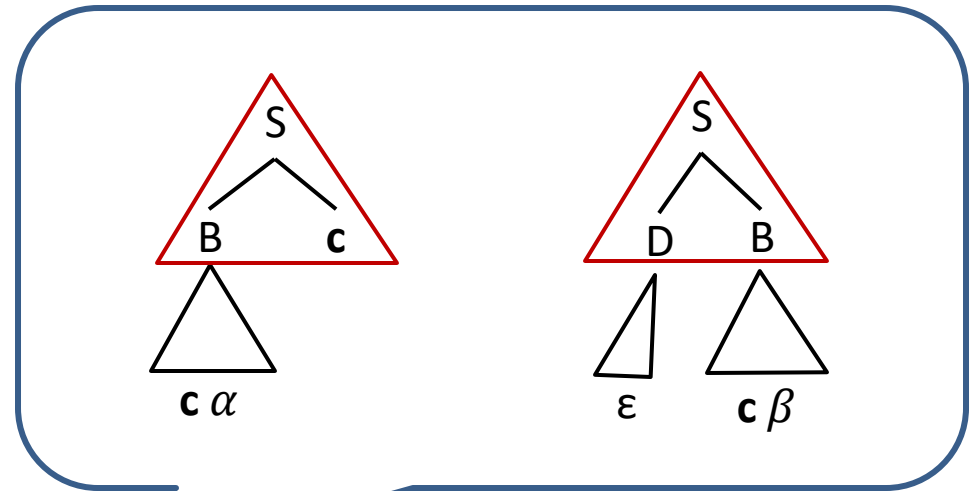
Why is a Table Collision a Problem?

CFG

$S \rightarrow Bc \mid DB$

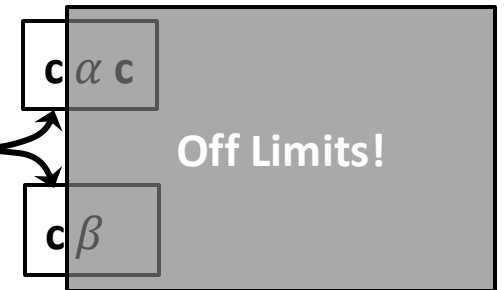
$B \rightarrow ab \mid cS$

$D \rightarrow d \mid \epsilon$



	a	b	c	d	ϵ
S	Bc DB		Bc DB	DB	
B	ab		cS		
D	ϵ		ϵ	d	

current
token



Recap

FIRST and FOLLOW sets define the parse table

If the grammar is LL(1), the table is unambiguous

- i.e., each cell has at most one entry

If the grammar is not LL(1) we can attempt a transformation sequence:

1. Remove left recursion
2. Left-factoring

Next time: Grammar transformations affect the structure of the parse tree. How does this affect syntax-directed translation (in particular, parse tree AST)?