

# Check-in

## Review: Predictive Parsing

Assume an LL(1) parser with...

this selector table:

	(	)	{	}
S	( S )	)	{	}

this syntax stack:

S
)
)
eof

and this ( lookahead token:

(

Draw the configuration of the parser after it processes the tokens ( )



# Housekeeping

Administrivia

## **Projects**

- P1 grades are in
- P2 (nominally) due Wednesday
- P3 out Friday

## **Trials**

- Trial 1 due tonight



# Housekeeping

Administrivia

## **Labs**

- Based on the confusion about abstract classes, I've decided to shift the labs a bit



University of Kansas | Drew Davidson

# *ECCS 665* **COMPILER** *CONSTRUCTION*

FIRST Sets



# Last Time

Review – Predictive Parsing

## Intro to Parsing

- Complexity

## A New Type of Language – LL(k)

- Intro
- LL(1) parsing

### You Should Know

- What parsing is
- What LL(1) languages are
- How an LL(1) parser operates



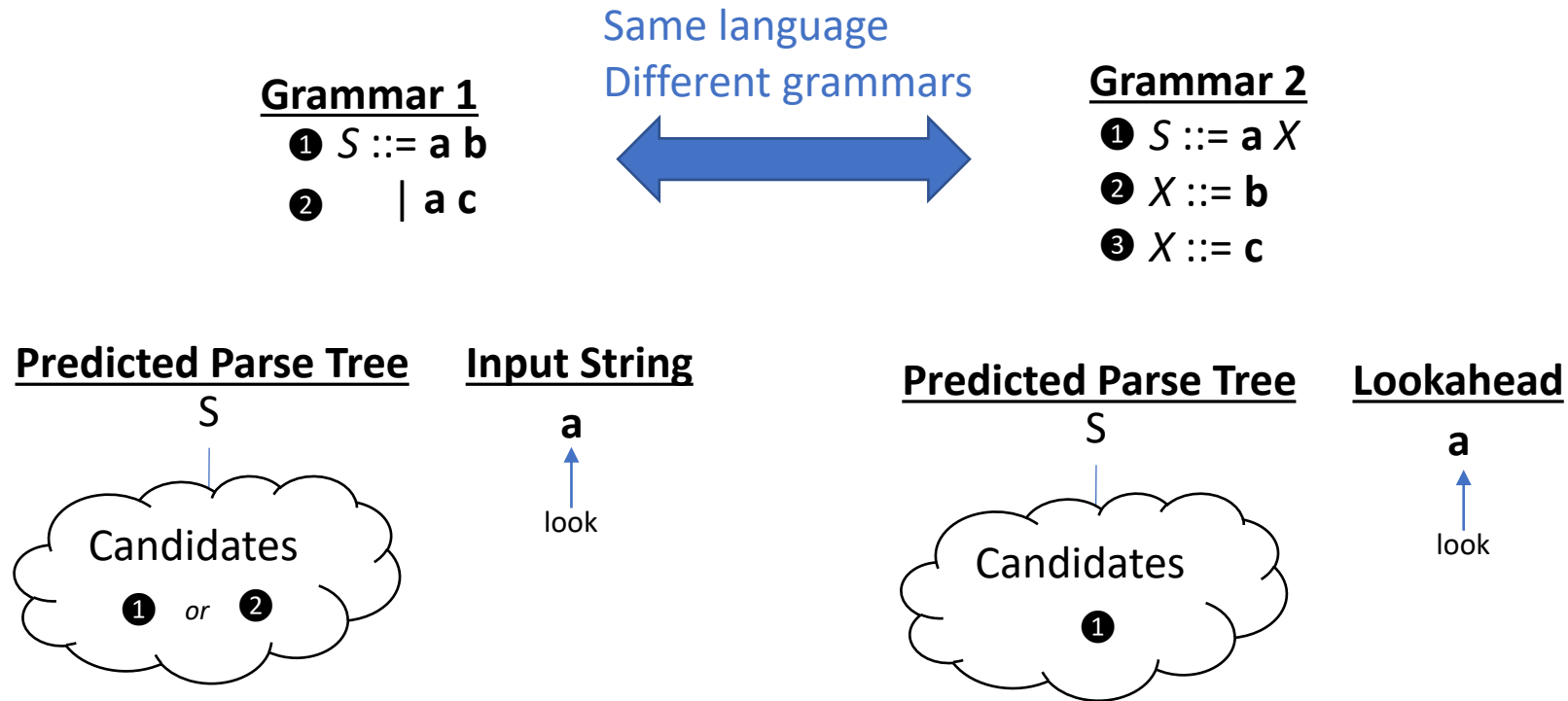
**Parsing**



# Where we Left Off

Review – Predictive Parsing

**The language might be LL(1) ... even when the grammar is not!**





# Today's Outline

Preview – FIRST Sets

## Transforming Grammars

- Fixing LL(1) “near misses”

## Building LL(1) Parsers

- What the selector table needs
- FIRST Sets



Parsing



# LL(1) Grammar Limitations

Transforming Grammars – Fixing LL(1) Near Misses

**Given a language, we can't always find an LL(1) grammar *even if one exists***

- Best we can do: simple transformations that remove “obvious” disqualifiers





# Checking if a Grammar is LL(1)

Transforming Grammars – Fixing LL(1) Near Misses

If either of the following hold, the grammar is not LL(1):

- The grammar is left-recursive
- The grammar **isn't** left-factored



We can transform *some* grammars while preserving the recognized language



# (Immediate) Left Recursion

Transforming Grammars – Fixing LL(1) Near Misses

- Recall, a grammar such that  $X \xRightarrow{+} X \alpha$  is left recursive
- A grammar is immediately left recursive if this can happen in one step:  
$$A \rightarrow A \alpha \mid \beta$$

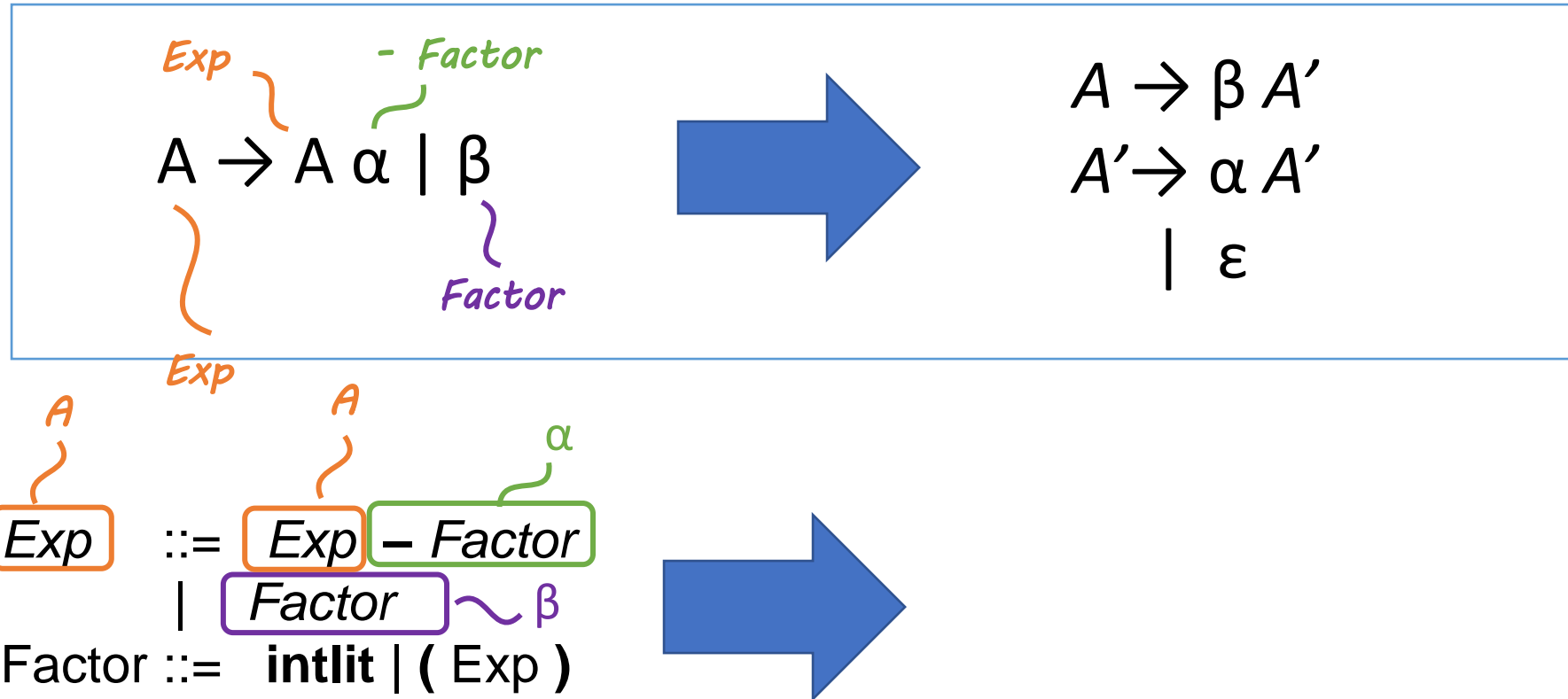






# Immediate Left Recursion Removal

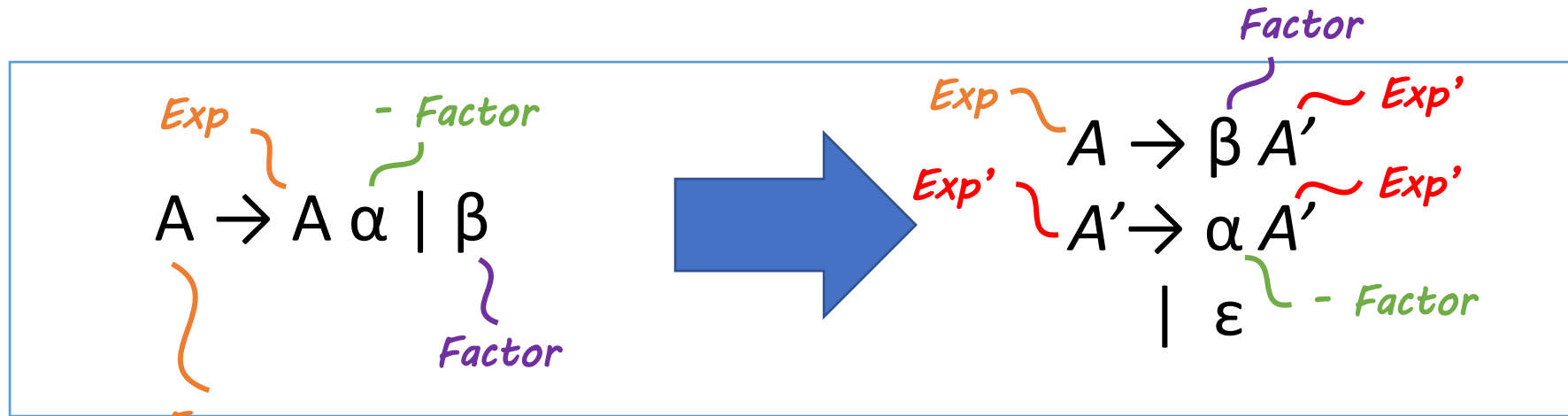
(Predictive) Parsing - LL(1) Transformations



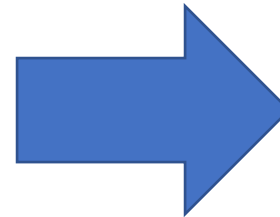


# Immediate Left Recursion Removal

(Predictive) Parsing - LL(1) Transformations



$Exp ::= Exp - Factor$   
 $Exp ::= Factor$   
 $Factor ::= \text{intlit} \mid ( Exp )$



$Exp ::= Factor Exp'$   
 $Exp' ::= - Factor Exp'$   
 $Exp' ::= \epsilon$   
 $Factor ::= \text{intlit} \mid ( Exp )$

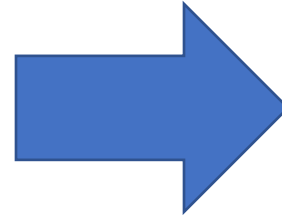


# Immediate Left Recursion Removal

*(Predictive) Parsing - LL(1) Transformations*

(general rule)

Given Productions

$$\begin{array}{l} A ::= \beta_1 \\ | \beta_2 \\ | \beta_n \\ | A \alpha_1 \\ | A \alpha_2 \\ | A \alpha_m \end{array}$$


Convert to

$$\begin{array}{l} A ::= \beta_1 A' \\ | \beta_2 A' \\ | \beta_n A' \\ A' ::= \alpha_1 A' \\ | \alpha_2 A' \\ | \alpha_m A' \\ | \varepsilon \end{array}$$



# Left Factoring Grammar

*(Predictive) Parsing - LL(1) Transformations*

- If a nonterminal has (at least) two productions whose RHS has a common prefix, the grammar is **not** left factored  
(and **not** an LL(1) grammar)

*Question: What makes this grammar not left-factored?*

$Exp ::= ( Exp )$   
|  $\{ Exp \}$   
|  $()$   
|  $ab$   
|  $bb$



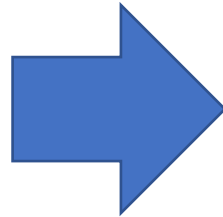
# Left Factoring: Simple Rule

*(Predictive) Parsing - LL(1) Transformations*

Given Productions

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$$

*Pull suffix into  
a new nonterminal*



Convert to

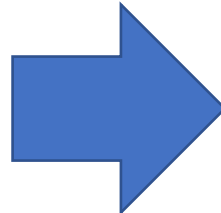
$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

*Add a new rule  
for suffixes*

$$X ::= \alpha \begin{matrix} \beta_1 \\ \beta_2 \end{matrix}$$

Diagram illustrating the transformation of a grammar rule. The left side shows a rule  $X ::= \alpha \begin{matrix} \beta_1 \\ \beta_2 \end{matrix}$  where  $\alpha$  is a common prefix and  $\beta_1, \beta_2$  are suffixes. The right side shows the transformed rule  $X ::= \alpha b X'$  where  $\alpha b$  is the new prefix and  $X'$  is a new nonterminal. The suffixes are now  $\beta_1 = cd$  and  $\beta_2 = ef$ .



$$X ::= \alpha b X'$$

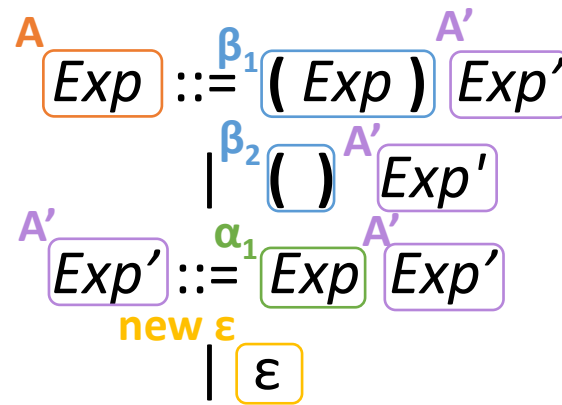
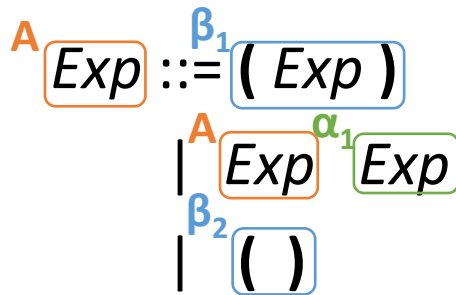
$$X' ::= \beta_1 \mid \beta_2$$



# Attempt LL(1) Conversion

(Predictive) Parsing - LL(1) Transformations

Remove immediate left-recursion



$$A \rightarrow A \alpha \mid \beta$$

becomes

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A'$$

$$\mid \epsilon$$



# Attempt LL(1) Conversion

(Predictive) Parsing - LL(1) Transformations

$Exp ::= ( Exp )$   
 $\quad | Exp \ Exp$   
 $\quad | ( )$

Remove immediate left-recursion



$Exp \overset{A}{::=} \overset{\alpha}{( Exp )} Exp' \overset{\beta_1}{}$   
 $\quad | \overset{\alpha}{( )} Exp' \overset{\beta_2}{}$   
 $Exp' ::= Exp \ Exp'$   
 $\quad | \epsilon$

Left-factored

$Exp \overset{A}{::=} \overset{\alpha}{( Exp'' } \overset{A'}{)}$   
 $Exp'' ::= Exp ) Exp' \overset{\beta_1}{}$   
 $\quad | ) Exp' \overset{\beta_2}{}$   
 $Exp' ::= Exp \ Exp'$   
 $\quad | \epsilon$

$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$  becomes

$A \rightarrow \alpha A'$   
 $A' \rightarrow \beta_1 \mid \beta_2$



# Attempt LL(1) Conversion

*(Predictive) Parsing - LL(1) Transformations*

$$\begin{aligned} \text{Exp} &::= ( \text{Exp} ) \\ &| \text{Exp} \text{ Exp} \\ &| ( ) \end{aligned}$$

Remove immediate left-recursion


$$\begin{aligned} \text{Exp} &::= ( \text{Exp} ) \text{Exp}' \\ &| ( ) \text{Exp}' \\ \text{Exp}' &::= \text{Exp} \text{Exp}' \\ &| \varepsilon \end{aligned}$$

Left-factored


$$\begin{aligned} \text{Exp} &::= ( \text{Exp}'' \\ \text{Exp}'' &::= \text{Exp} ) \text{Exp}' \\ &| ) \text{Exp}' \\ \text{Exp}' &::= \text{Exp} \text{Exp}' \\ &| \varepsilon \end{aligned}$$



# Current Status

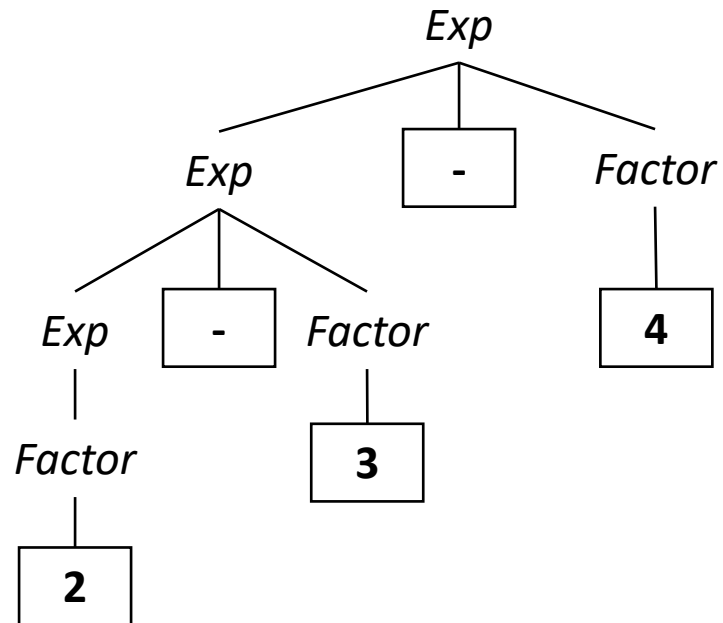
*(Predictive) Parsing - LL(1) Transformations*

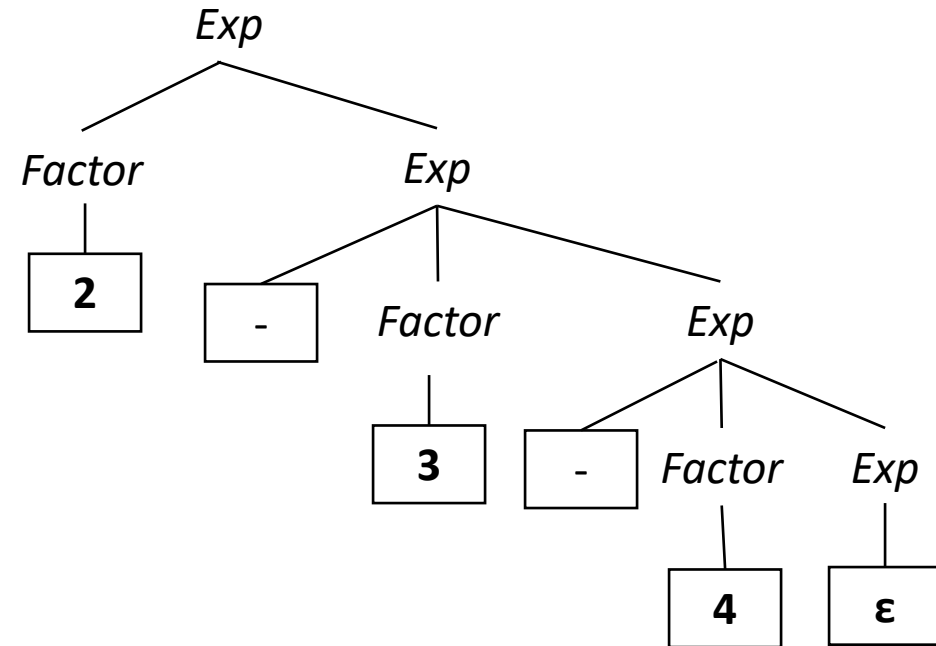
- We've removed 2 disqualifiers from LL(1)
  - Left-recursive grammar
  - **Not** Left-Factored grammar



# Let's Check on the Parse Tree

## LL(1) Grammar Transformations

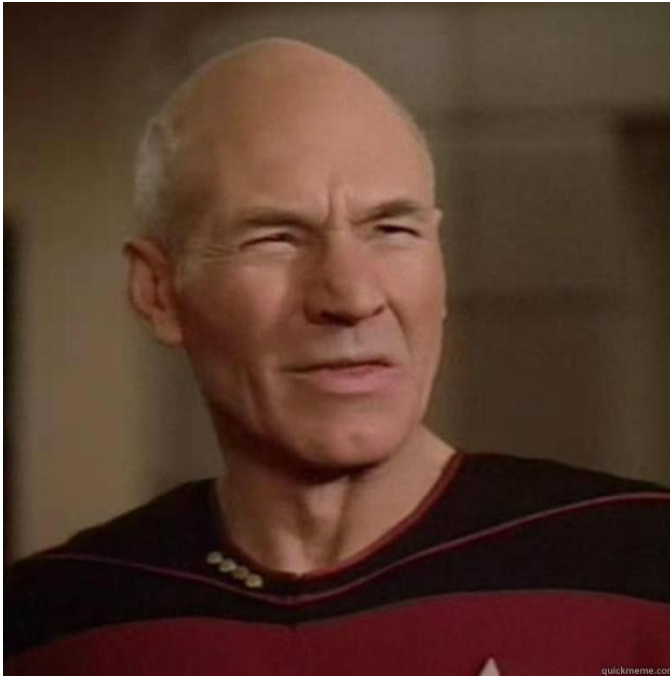
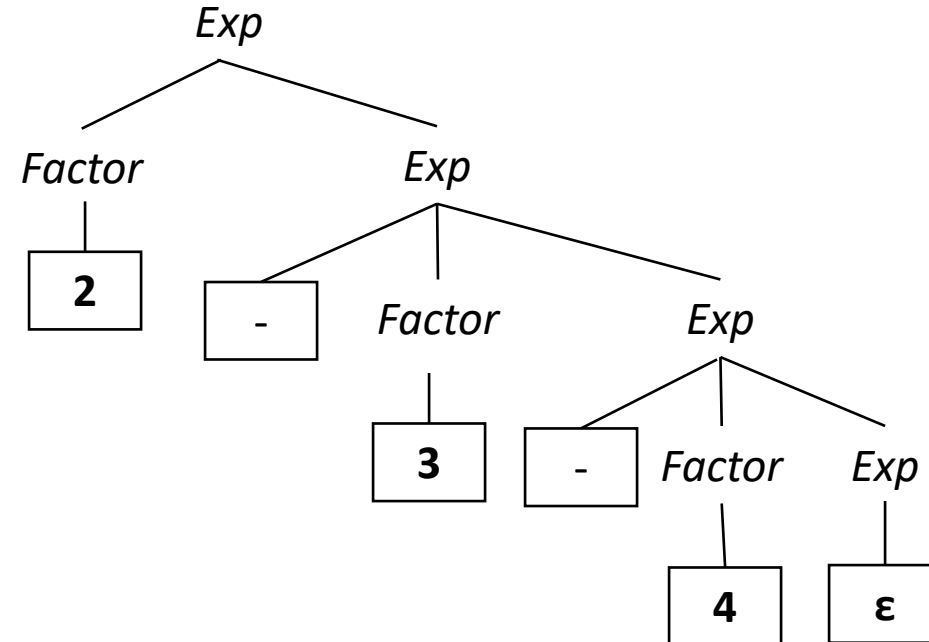
$$\begin{aligned} \text{Exp} &\rightarrow \text{Exp} - \text{Factor} \\ &\quad | \text{Factor} \\ \text{Factor} &\rightarrow \text{intlit} \mid ( \text{Exp} ) \end{aligned}$$


$$\begin{aligned} \text{Exp} &\rightarrow \text{Factor Exp}' \\ \text{Exp}' &\rightarrow - \text{Factor Exp}' \\ &\quad | \epsilon \\ \text{Factor} &\rightarrow \text{intlit} \mid ( \text{Exp} ) \end{aligned}$$




# Let's Check on the Parse Tree

## *LL(1) Grammar Transformations*


$$\begin{aligned} \text{Exp} &\rightarrow \text{Factor Exp}' \\ \text{Exp}' &\rightarrow - \text{Factor Exp}' \\ &\quad | \quad \epsilon \\ \text{Factor} &\rightarrow \text{intlit} \mid ( \text{Exp} ) \end{aligned}$$




# Nevermind, We'll Fix Parse Trees Later

*LL(1) Grammar Transformations*

- \\_ ( ツ ) \\_ /



# Today's Outline

## Lecture 9 – FIRST sets

### Transforming Grammars



- Fixing LL(1) “near misses”

### Building LL(1) Parsers

- Understanding LL(1) Selector Tables
- FIRST Sets



Parsing



# Recall the LL(1) Parser's Operation

## *Building LL(1)Selector Table*

### **LL(1)**

- Processes **Left-to-right**
- **Leftmost** derivation
- **1** token of lookahead

### **Predictive Parser: “guess & check”**

- Starts at the root, *guesses* how to unfold a nonterminal (derivation step)
- *Checks* that terminals match prediction



# Recall the LL(1) Parser's Operation

## Building LL(1) Selector Table

### Example LL(1) Grammar:

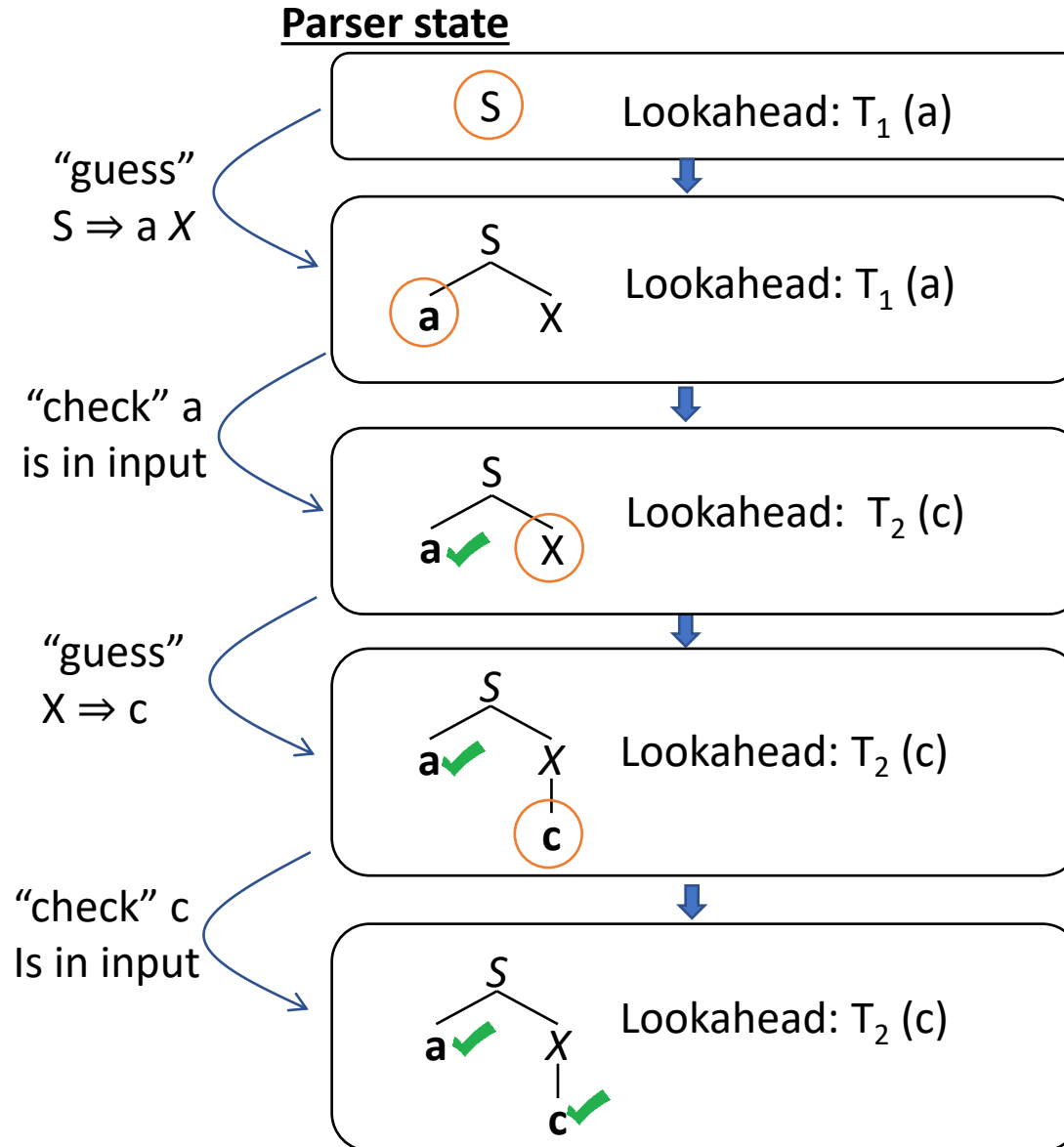
$S ::= a X$

$X ::= b a \mid c$

### Example Input:

$a \quad c$   
 $\uparrow \quad \uparrow$   
 $T_1 \quad T_2$

In practice,  
table-driven parser  
uses a stack to  
match this tree





# How does the Parser Guess?

## *Building Parser Tables*

**The intuition is a bit tricky**

- We need to get into the mindset of the parser



*Pretend your consciousness has been transported inside an LL(1) parser*



# Become the Parser

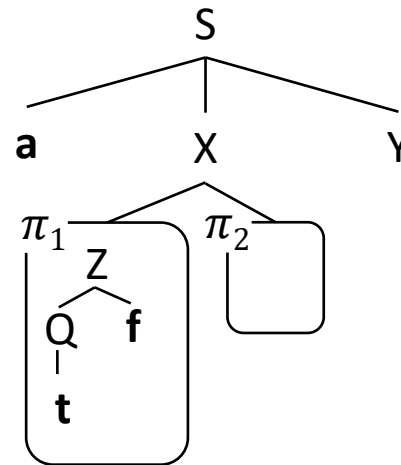
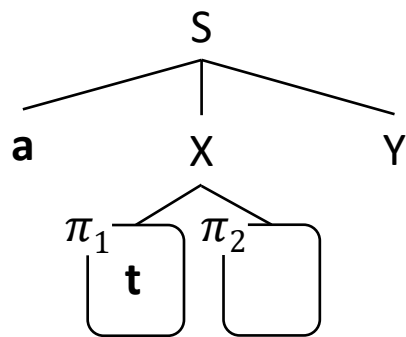
## Building Parser Tables

You need to unfold a nonterminal  $X$  with lookahead token  $t$

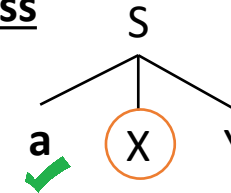
Assume there's an  $X$  production  $X ::= \pi_1 \pi_2$  (where  $\pi_1$  and  $\pi_2$  are some kind of symbol)

How do we know to guess this production?

Case 1:  $\pi_1$  subtree may start with  $t$



Parse in Progress



Lookahead:  $T_2(t)$

Grammar Fragment

...  
 $X ::= \pi_1 \pi_2$   
 ...



# Become the Parser

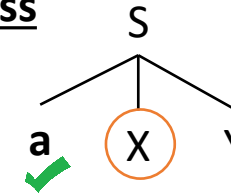
## Building Parser Tables

You need to unfold a nonterminal  $X$  with lookahead token  $t$

Assume there's an  $X$  production  $X ::= \pi_1 \pi_2$  (where  $\pi_1$  and  $\pi_2$  are some kind of symbol)

How do we know to guess this production?

Parse in Progress

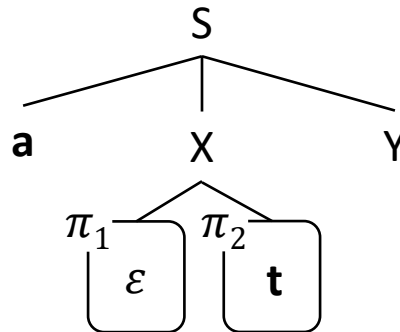


Lookahead:  $T_2(t)$

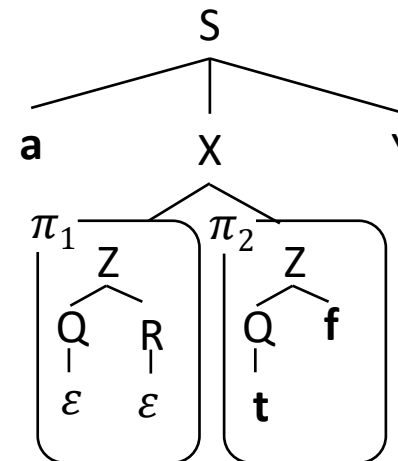
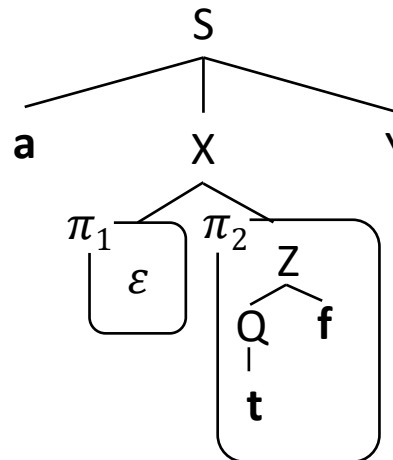
Grammar Fragment

...  
 $X ::= \pi_1 \pi_2$   
 ...

Case 1:  $\pi_1$  subtree may start with  $t$



Case 2:  $\pi_1$  subtree may be empty and  $\pi_2$  starts with  $t$





# Become the Parser

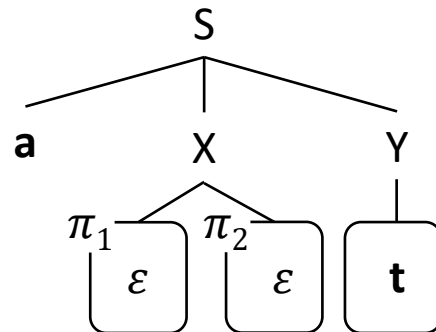
## Building Parser Tables

You need to unfold a nonterminal  $X$  with lookahead token  $\mathbf{t}$

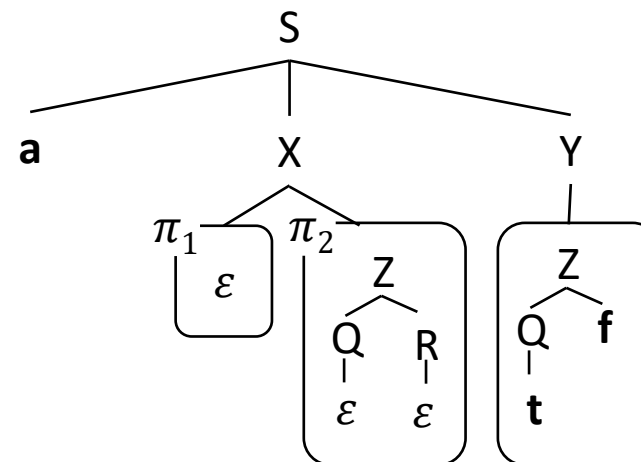
Assume there's an  $X$  production  $X ::= \pi_1 \pi_2$  (where  $\pi_1$  and  $\pi_2$  are some kind of symbol)

How do we know to guess this production?

Case 1:  $\pi_1$  subtree may start with  $\mathbf{t}$

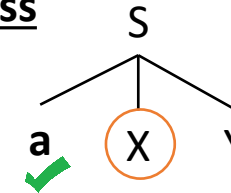


Case 2:  $\pi_1$  subtree may be empty and  $\pi_2$  starts with  $\mathbf{t}$



Case 3: both  $\pi_1$  and  $\pi_2$  may be empty and the sibling may start with  $\mathbf{t}$

Parse in Progress



Lookahead:  $T_2(\mathbf{t})$

Grammar Fragment

...  
 $X ::= \pi_1 \pi_2$   
 ...



# Become the Parser

## *Building Parser Tables*

You need to unfold a nonterminal  $X$  with lookahead token  $\mathbf{t}$

Assume there's an  $X$  production  $X ::= \pi_1 \pi_2$  (where  $\pi_1$  and  $\pi_2$  are some kind of symbol)

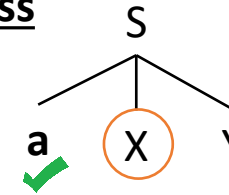
How do we know to guess this production?

Case 1:  $\pi_1$  subtree may start with  $\mathbf{t}$

Case 2:  $\pi_1$  subtree may be empty and  $\pi_2$  starts with  $\mathbf{t}$

Case 3: both  $\pi_1$  and  $\pi_2$  may be empty and the sibling may start with  $\mathbf{t}$

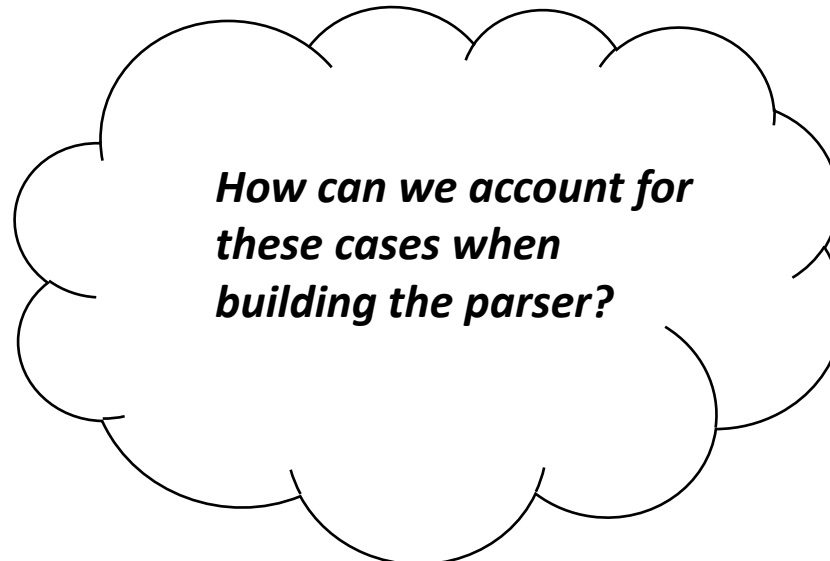
Parse in Progress



Lookahead:  $T_2(\mathbf{t})$

Grammar Fragment

...  
 $X ::= \pi_1 \pi_2$   
...





# Become the Parser

## *Building Parser Tables*

You need to unfold a nonterminal  $X$  with lookahead token  $\mathbf{t}$

Assume there's an  $X$  production  $X ::= \pi_1 \pi_2$  (where  $\pi_1$  and  $\pi_2$  are some kind of symbol)

How do we know to guess this production?

Case 1:  $\pi_1$  subtree may start with  $\mathbf{t}$

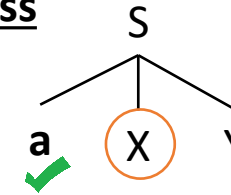
Case 2:  $\pi_1$  subtree may be empty and  $\pi_2$  starts with  $\mathbf{t}$

Case 3: both  $\pi_1$  and  $\pi_2$  may be empty and the sibling may start with  $\mathbf{t}$

**FIRST Sets**

**FOLLOW Sets**

Parse in Progress



Lookahead:  $T_2(\mathbf{t})$

Grammar Fragment

...  
 $X ::= \pi_1 \pi_2$   
...

Two sets are sufficient to capture these cases and to build the selector table



# Today's Outline

## Lecture 9 – FIRST sets

### Transforming Grammars

- ✓ Fixing LL(1) “near misses”

### Building LL(1) Parsers

- ✓ Reverse-Engineering Selector Tables

- FIRST Sets



Parsing



# An Informal Definition

*Building LL(1) Selector Table: FIRST sets, single symbol*

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



# A Formal Definition

*Building LL(1) Selector Table: FIRST sets, single symbol*

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

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

$$\left\{ \hat{\alpha} \mid \left( \hat{\alpha} \in \Sigma \wedge \alpha \xRightarrow{*} \hat{\alpha}\beta \right) \vee \left( \hat{\alpha} = \epsilon \wedge \alpha \xRightarrow{*} \epsilon \right) \right\}$$

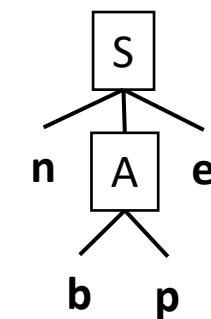


# A Parse Tree Perspective

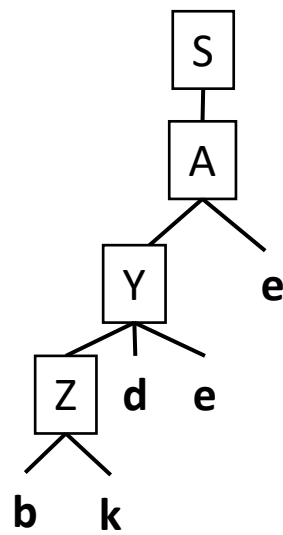
*Building LL(1) Selector Table: FIRST sets, single symbol*

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

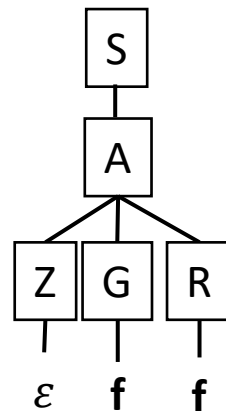
What does the parse tree say about  $\text{FIRST}(A)$ ?



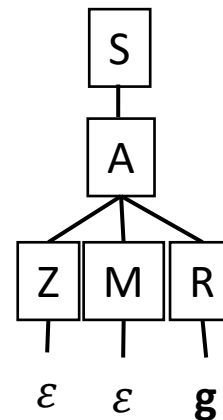
$\text{FIRST}(A)$   
includes  
 $\{b\}$



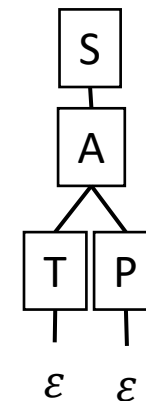
Again,  $\text{FIRST}(A)$   
includes  
 $\{b\}$



$\text{FIRST}(A)$   
Includes  
 $\{f\}$



$\text{FIRST}(A)$   
includes  
 $\{g\}$



$\text{FIRST}(A)$   
includes  
 $\{\epsilon\}$

If these were the only possible parse trees, then  $\text{FIRST}(A) = \{b, f, g, \epsilon\}$



# A Parse Tree Perspective

*Building LL(1) Selector Table: FIRST sets, single symbol*

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

## **This isn't how you build FIRST sets**

- Looking at parse trees is illustrative for concepts only
- We need to derive FIRST sets directly from the grammar



# Building FIRST Sets: Methodology

## *Building Parser Tables*

### **First sets exist for any arbitrary string of symbols $\alpha$**

- Defined in terms of FIRST sets for a single symbol
  - FIRST of an alphabet terminal
  - FIRST for  $\epsilon$
  - FIRST for a nonterminal
- Use single-symbol FIRST to construct symbol-string FIRSTS



# Rules for Single Symbols

## *Building Parser Tables*

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

### Building FIRST for terminals

$\text{FIRST}(\mathbf{t}) = \{ \mathbf{t} \}$  for  $\mathbf{t}$  in  $\Sigma$

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



### Building FIRST( $X$ ) for nonterminal $X$

For each  $X ::= \alpha_1 \alpha_2 \dots \alpha_n$

$C_1$ : add  $\text{FIRST}(\alpha_1) - \varepsilon$

$C_2$ : If  $\varepsilon$  could “prefix”  $\text{FIRST}(\alpha_k)$ , add  $\text{FIRST}(\alpha_k) - \varepsilon$

$C_3$ : If  $\varepsilon$  is in every FIRST set  $\alpha_1 \dots \alpha_n$ , add  $\varepsilon$



# Rules for Single Symbols

## Building LL(1) Parsers

### **Building FIRST( $X$ ) for nonterminal $X$**

For each  $X ::= \alpha_1 \alpha_2 \dots \alpha_n$

$C_1$ : add FIRST( $\alpha_1$ ) -  $\varepsilon$

$C_2$ : If  $\varepsilon$  could “prefix” FIRST( $\alpha_k$ ), add FIRST( $\alpha_k$ ) -  $\varepsilon$

$C_3$ : If  $\varepsilon$  is in every FIRST set  $\alpha_1 \dots \alpha_n$ , add  $\varepsilon$



# Rules for Single Symbols

## Building LL(1) Parsers

### Building FIRST( $X$ ) for nonterminal $X$

For each  $X ::= \alpha_1 \alpha_2 \dots \alpha_n$

$C_1$ : add FIRST( $\alpha_1$ ) -  $\epsilon$

$C_2$ : If  $\epsilon$  could “prefix” FIRST( $\alpha_k$ ), add FIRST( $\alpha_k$ ) -  $\epsilon$

$C_3$ : If  $\epsilon$  is in every FIRST set  $\alpha_1 \dots \alpha_n$ , add  $\epsilon$

Say there's a production

$X ::= Y Z R T$

and we know

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

FIRST( $Z$ ) = {  $\epsilon$ , **b**, **m** }

FIRST( $R$ ) = { **c** }

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

By  $C_2$  clause FIRST( $X$ ) includes **b**, **m** and **c**

**b,m** because FIRST of every symbol before the 2<sup>nd</sup> includes  $\epsilon$ )

*Z in this case* 

**c** because FIRST of every symbol before the 3<sup>rd</sup> includes  $\epsilon$ )

*R in this case* 

FIRST( $X$ ) does not add **d** in this clause  
because not every FIRST set before the T  
includes  $\epsilon$



# Building FIRST Sets for Symbol Strings

## Building LL(1) Parsers

### **Building FIRST( $\alpha$ )**

Let  $\alpha$  be composed of symbols  $\alpha_1 \alpha_2 \dots \alpha_n$

$C_1$ : add FIRST( $\alpha_1$ ) -  $\varepsilon$

$C_2$ : If  $\alpha_1 \dots \alpha_{k-1}$  is nullable, add FIRST( $\alpha_k$ ) -  $\varepsilon$

$C_3$ : If  $\alpha_1 \dots \alpha_n$  is nullable, add  $\varepsilon$

### **Base Cases:**

$\alpha_i$  is a terminal  $t$ . Add  $t$

$\alpha_i$  is a nonterminal  $X$ . Add every leaf symbol that could begin an  $X$  subtree  
(this gets a bit complicated due to dependencies)



# Summary: Explored the LL(1) Mindset

FIRST Sets

## LL(1) “Parseability” Qualification

- Knowing the leftmost terminal of a parse (sub)tree is enough to pick the next derivation step

## Elusive Conditions

- Two different rules could start with the same terminal (not left factored)
- The same rule(s) could be applied repeatedly (left recursive)

## Began choosing matching productions to input

- What terminal could the production be the start of (FIRST)?