

Top-Down Parsing

Lecture 4

Announcements

- PA 1
 - Due today (17 June) at 11:59pm
- WA 2
 - Released today
- PA 2
 - Released today

Problems with Top Down Parsing

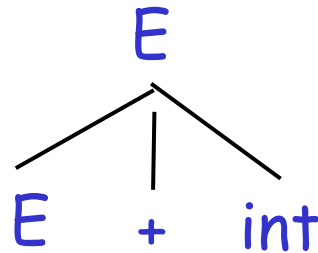
- Left Recursion in CFG may cause parser to loop forever!
- In there is a production of form $A \rightarrow A\alpha$, we say the grammar has left recursion

$$E \rightarrow E + \text{int} \mid \text{int}$$

- Solution: Remove Left Recursion...
 - without changing the Language defined by the Grammar.

Problems with Top Down Parsing (Example)

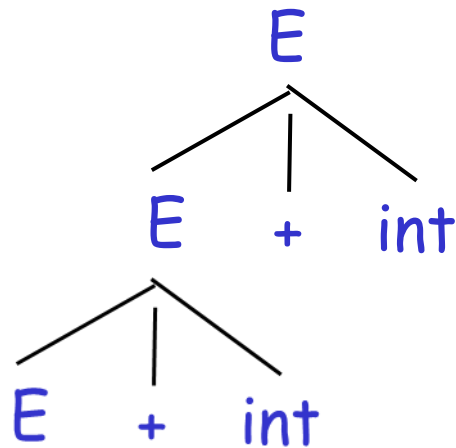
$E \rightarrow E + \text{int} \mid \text{int}$



$\text{int} + \text{int}$
↑

Problems with Top Down Parsing (Example)

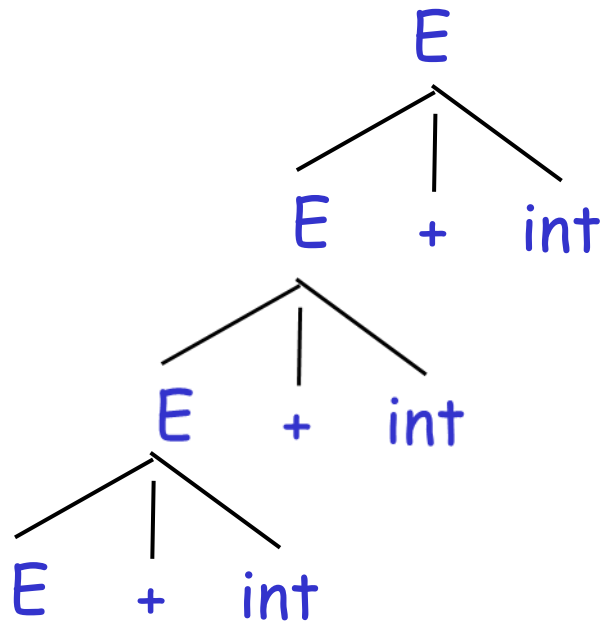
$E \rightarrow E + \text{int} \mid \text{int}$



int + int
↑

Problems with Top Down Parsing (Example)

$E \rightarrow E + \text{int} \mid \text{int}$



int + int
↑

Elimination of Left Recursion

- Consider the left-recursive grammar

$$S \rightarrow S \alpha \mid \beta$$

- S generates all strings starting with a β and followed by a number of α

- Can rewrite using right-recursion

$$S \rightarrow \beta S'$$

$$S' \rightarrow \alpha S' \mid \varepsilon$$

More Elimination of Left-Recursion

- In general

$$S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$

- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$

- Rewrite as

$$S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$$

$$S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \varepsilon$$

General Left Recursion

- The grammar

$$S \rightarrow A \alpha \mid \delta$$

$$A \rightarrow S \beta$$

is also left-recursive because

$$S \rightarrow^+ S \beta \alpha$$

- This left-recursion can also be eliminated
- See Dragon Book for general algorithm
 - Section 4.3.3

Predictive Parsing and Left Factoring

- Consider the grammar

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$

- Hard to predict because
 - For T two productions start with int
 - For E it is not clear how to predict
- We need to left-factor the grammar

Left-Factoring Example

- Recall the grammar

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$

- Factor out common prefixes of productions

$$E \rightarrow T X$$

$$X \rightarrow + E \mid \varepsilon$$

$$T \rightarrow (E) \mid \text{int} Y$$

$$Y \rightarrow * T \mid \varepsilon$$

Predictive Parsers

- Parser can “predict” which production to use
 - By looking at the next few tokens
 - No backtracking
- Predictive parsers accept LL(k) grammars
 - L means “left-to-right” scan of input
 - L means “leftmost derivation”
 - k means “predict based on k tokens of lookahead”
 - In practice, LL(1) is used

LL(1) Parsing Table Example

- Left-factored grammar

$$E \rightarrow TX$$

$$X \rightarrow + E \mid \varepsilon$$

$$T \rightarrow (E) \mid \text{int } Y$$

$$Y \rightarrow * T \mid \varepsilon$$

- The LL(1) parsing table:

next input token

	int	*	+	()	\$
E	TX			TX		
X			+E		ε	ε
T	int Y			(E)		
Y		*T	ε		ε	ε

leftmost non-terminal

rhs of production to use

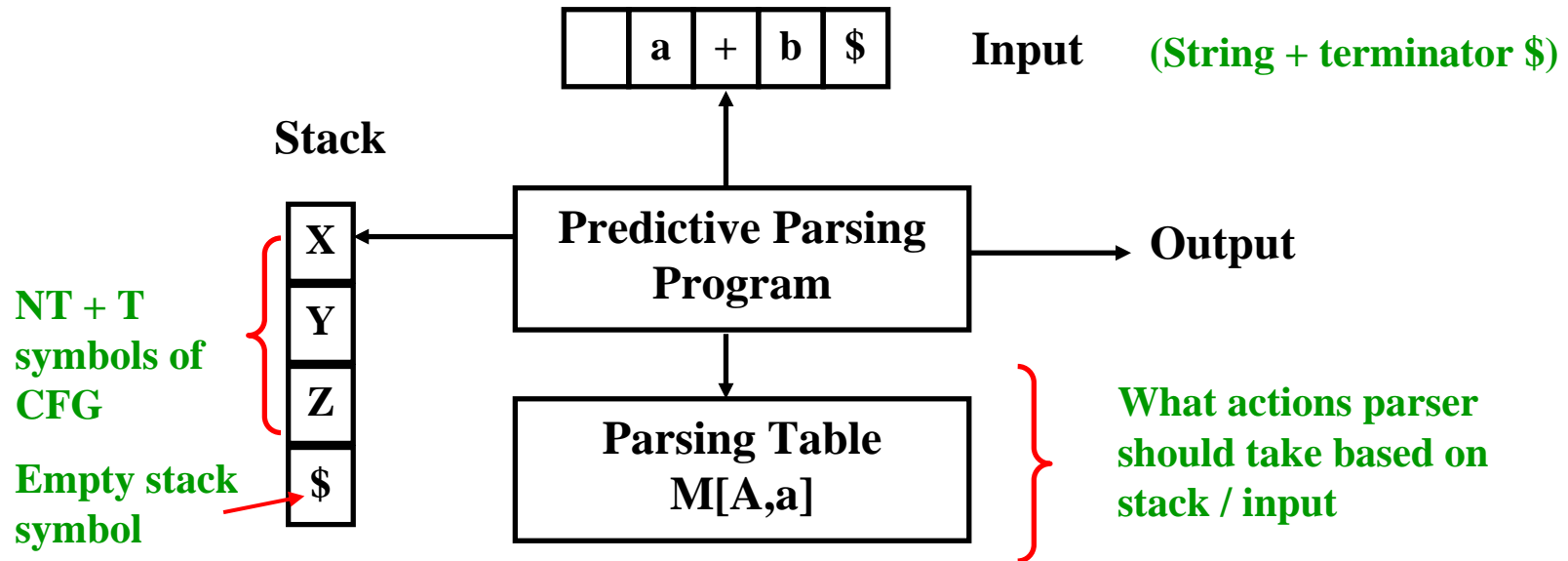
LL(1) Parsing Table Example (Cont.)

- Consider the $[E, \text{int}]$ entry
 - "When current non-terminal is E and next input is int , use production $E \rightarrow TX$ "
 - This can generate an int in the first position
- Consider the $[Y, +]$ entry
 - "When current non-terminal is Y and current token is $+$, get rid of Y "
 - Y can be followed by $+$ only if $Y \rightarrow \varepsilon$

LL(1) Parsing Tables. Errors

- Blank entries indicate error situations
- Consider the $[E, *]$ entry
 - "There is no way to derive a string starting with $*$ from non-terminal E "

LL(1) Parsing Algorithm



General parser behavior: X : top of stack a : current token

1. When $X=a = \$$ halt, accept, success
2. When $X=a \neq \$$, POP X off stack, advance input, go to 1.
3. When X is a non-terminal, examine $M[X, a]$, if it is an error, call recovery routine if $M[X, a] = \{UVW\}$, POP X , PUSH U,V,W , and **DO NOT** advance input

LL(1) Parsing Example

Stack	Input	Action
E \$	int * int \$	T X
T X \$	int * int \$	int Y
int Y X \$	int * int \$	terminal
Y X \$	* int \$	* T
* T X \$	* int \$	terminal
T X \$	int \$	int Y
int Y X \$	int \$	terminal
Y X \$	\$	ϵ
X \$	\$	ϵ
\$	\$	ACCEPT

Constructing Parsing Tables: The Intuition

- Consider non-terminal A , production $A \rightarrow \alpha$, & token t
- $T[A,t] = \alpha$ in two cases:
 - If $\alpha \rightarrow^* t \beta$
 - α can derive a t in the first position
 - We say that $t \in \text{First}(\alpha)$
 - If $A \rightarrow \alpha$ and $\alpha \rightarrow^* \varepsilon$ and $S \rightarrow^* \beta A t \delta$
 - Useful if stack has A , input is t , and A cannot derive t
 - In this case only option is to get rid of A (by deriving ε)
 - Can work only if t can follow A in at least one derivation
 - We say $t \in \text{Follow}(A)$

Constructing LL(1) Parsing Tables

- Construct a parsing table T for CFG G
- For each production $A \rightarrow \alpha$ in G do:
 - For each terminal $t \in \text{First}(\alpha)$ do
 - $T[A, t] = \alpha$
 - If $\varepsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do
 - $T[A, t] = \alpha$
 - If $\varepsilon \in \text{First}(\alpha)$ and $\$ \in \text{Follow}(A)$ do
 - $T[A, \$] = \alpha$

Example 1

$$\begin{array}{l} E \rightarrow TX \\ T \rightarrow (E) \mid \text{int } Y \end{array} \quad \begin{array}{l} X \rightarrow +E \mid \varepsilon \\ Y \rightarrow *T \mid \varepsilon \end{array}$$

	int	*	+	()	\$
E	TX			TX		
X			+E		ε	ε
T	int Y			(E)		
Y		*T	ε		ε	ε

Example 2

$S \rightarrow Sa \mid b$
 $\text{First}(S) = \{b\}$
 $\text{Follow}(S) = \{\$, a\}$

	a	b	\$
S		b, Sa	

Notes on LL(1) Parsing Tables

- If any entry is multiply defined then G is not LL(1)
 - If G is ambiguous
 - If G is left recursive
 - If G is not left-factored
 - And in other cases as well
- Most programming language CFGs are not LL(1)

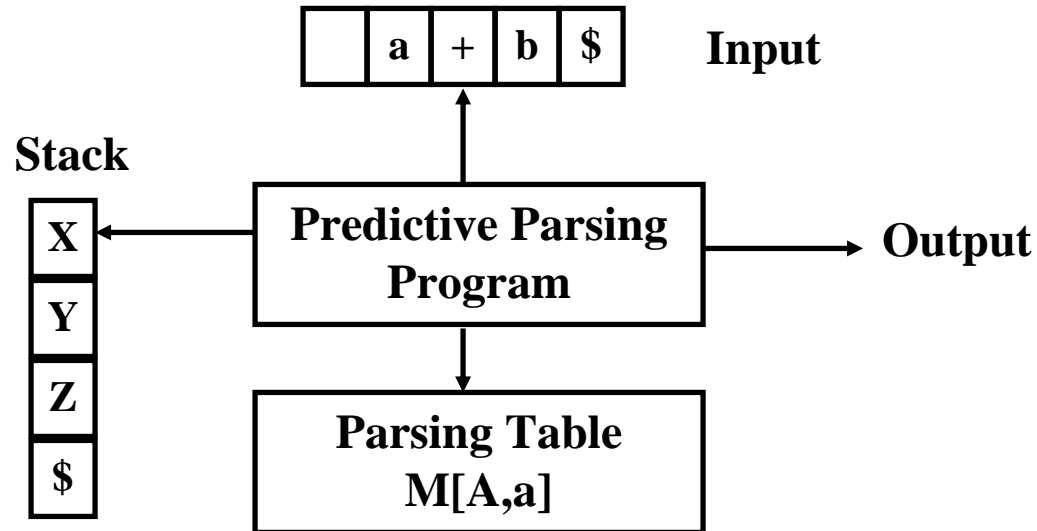
Notes on LL(1) Grammars

Grammar is LL(1) \Leftrightarrow when for all $A \rightarrow \alpha \mid \beta$

1. $\text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$; besides, only one of α or β can derive ϵ
2. if α derives ϵ , then $\text{Follow}(A) \cap \text{First}(\beta) = \emptyset$

It may not be possible for a grammar to be manipulated into an LL(1) grammar

Implementing Panic Mode in LL(1)



Error situations include:

- 1.If X is a terminal and it doesn't match current token.
- 2.If $M[X, \text{Input}]$ is empty - No allowable actions

Panic-Mode Recovery

- Assume in a syntax error, non-terminal A is on the top of the stack.
- The choice for a synchronizing set is important.
 - define the synchronizing set of A to be $\text{Follow}(A)$. Then skip input until a token in $\text{Follow}(A)$ appears and then pop A from the stack. Resume parsing...
 - add symbols of $\text{FIRST}(A)$ to the synchronizing set. In this case, we skip input and once we find a token in $\text{FIRST}(A)$, we resume parsing from A .

Panic-Mode Recovery (Cont.)

Modify the empty cells of the Parsing Table.

1. if $M[A, a] = \{\text{empty}\}$ and a belongs to $\text{Follow}(A)$ then we set $M[A, a] = \text{"synch"}$

Error-recovery Strategy :

If $A = \text{top-of-the-stack}$ and $a = \text{current-token}$,

1. If A is NT and $M[A, a] = \{\text{empty}\}$ then skip a from the input.
2. If A is NT and $M[A, a] = \{\text{synch}\}$ then pop A .
3. If A is a terminal and $A \neq a$ then pop A (This is essentially inserting A before a).

Parse Table / Example

	id	+	*	()	\$
E	TE'			TE'	synch	synch
E'		+TE'			ε	ε
T	FT'	synch		FT'	synch	synch
T'		ε	*FT'		ε	ε
F	id	synch	synch	(E)	synch	synch

Pop top of stack NT
for "synch" cells

Skip current-token
for empty cells

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid id$$

Parsing Example

	id	+	*	()	\$
E	TE'			TE'	synch	synch
E'		+TE'			∈	∈
T	FT'	synch		FT'	synch	synch
T'		∈	*FT'		∈	∈
F	id	synch	synch	(E)	synch	synch

STACK	INPUT	Remark
E \$	+ id * + id \$	error, skip +
E \$	id * + id \$	
TE' \$	id * + id \$	
FT' E' \$	id * + id \$	
id T' E' \$	id * + id \$	
T' E' \$	* + id \$	
* FT' E' \$	* + id \$	
FT' E' \$	+ id \$	

Possible Error Msg:
 "Misplaced +
 I am skipping it"

$E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

Parsing Example (Cont.)

	id	+	*	()	\$
E	TE'			TE'	synch	synch
E'		+TE'			ε	ε
T	FT'	synch		FT'	synch	synch
T'		ε	*FT'		ε	ε
F	id	synch	synch	(E)	synch	synch

STACK	INPUT	Remark
FT'E'\$	+ id \$	error, M[F,+] = synch , F is popped
T'E'\$	+ id \$	
E'\$	+ id \$	
+TE'\$	+ id \$	
TE'\$	id \$	
FT'E'\$	id \$	
idT'E'\$	id \$	
T'E'\$	\$	
E'\$	\$	
\$	\$	

Possible Error Msg:
"Missing Term"

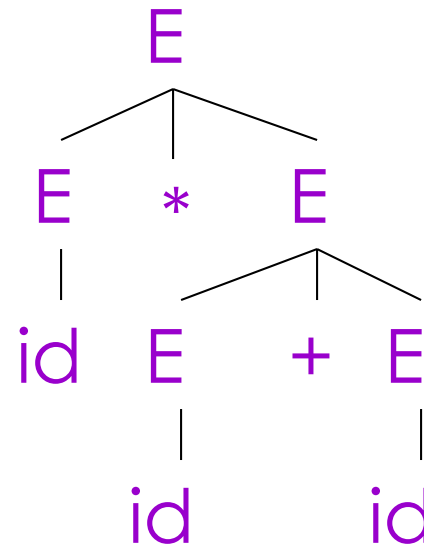
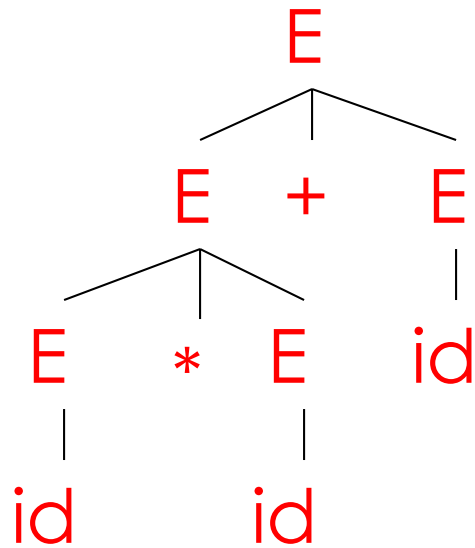
$E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

Ambiguity

- Grammar $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- String $id * id + id$

Ambiguity (Cont.)

This string has two parse trees



Ambiguity (Cont.)

- A grammar is *ambiguous* if it has more than one parse tree for some string
 - Equivalently, there is more than one right-most or left-most derivation for some string
- Ambiguity is **BAD**
 - Leaves meaning of some programs ill-defined

Dealing with Ambiguity

- There are several ways to handle ambiguity
- Most direct method is to rewrite grammar unambiguously

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

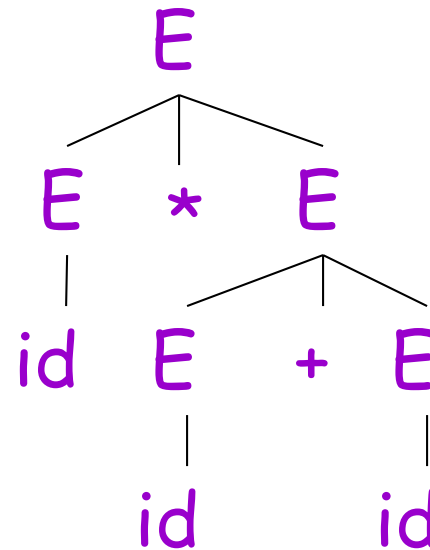
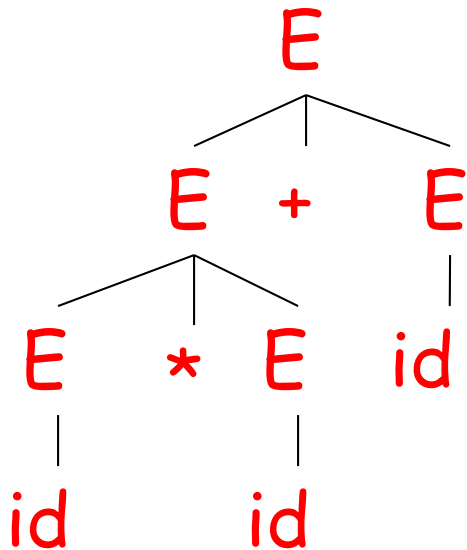
- Enforces precedence of $*$ over $+$

Ambiguity in Arithmetic Expressions

- Recall the grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

- The string $id * id + id$ has two parse trees:

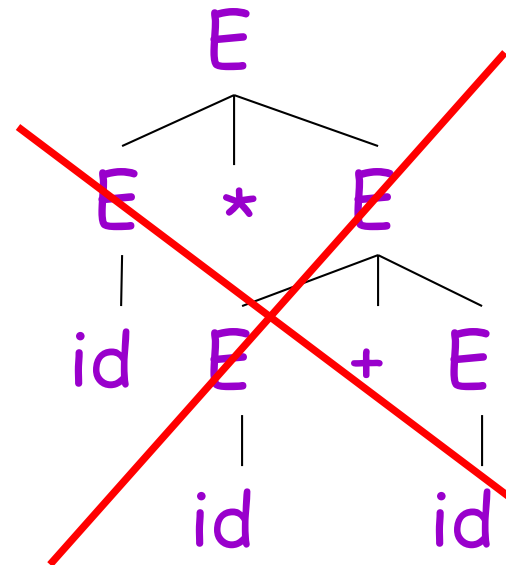
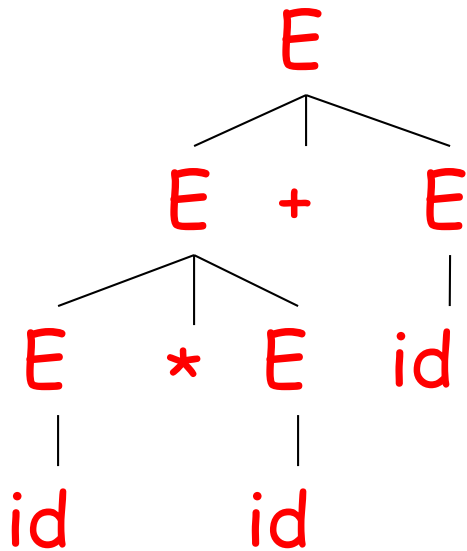


Ambiguity in Arithmetic Expressions

- Recall the grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

- The string $id * id + id$ has two parse trees:



Ambiguity: The Dangling Else

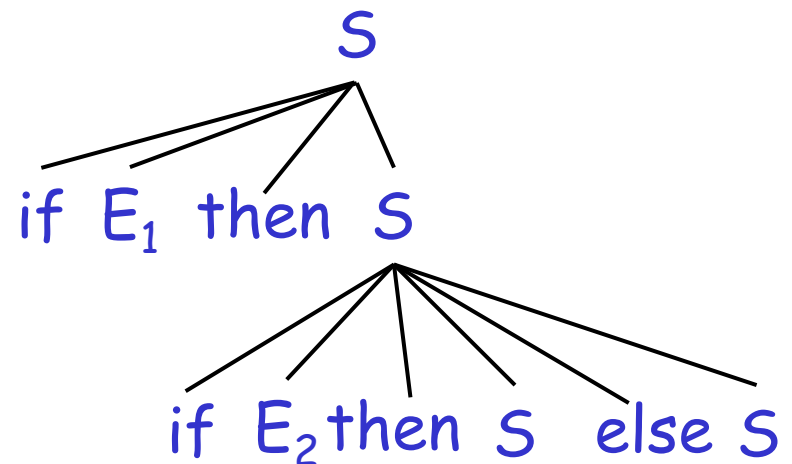
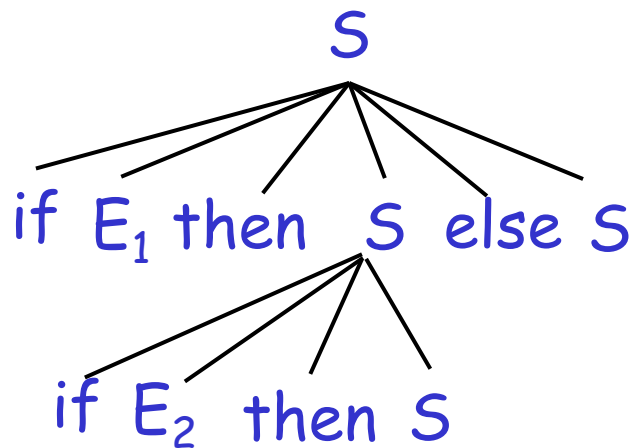
- Consider the grammar
$$S \rightarrow \text{if } E \text{ then } S$$
$$| \text{if } E \text{ then } S \text{ else } S$$
$$| \text{OTHER}$$
- This grammar is also ambiguous

The Dangling Else: Example

- The expression

if E_1 then if E_2 then S else S

has two parse trees



- Typically we want the second form

The Dangling Else: A Fix

- `else` matches the closest unmatched `then`
- We can describe this in the grammar

$S \rightarrow$ MIF /* all `then` are matched */
 | UIF /* some `then` is unmatched */

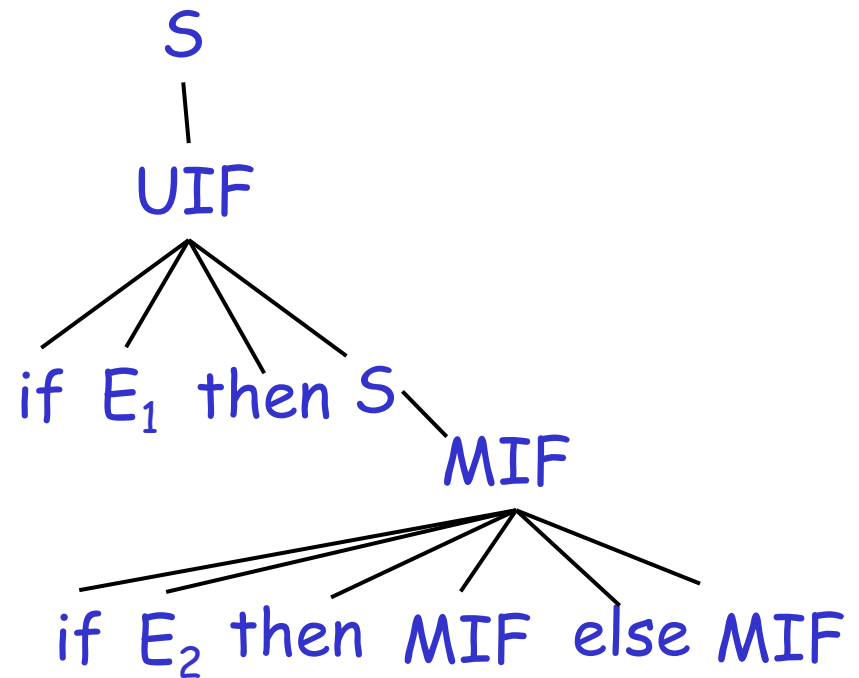
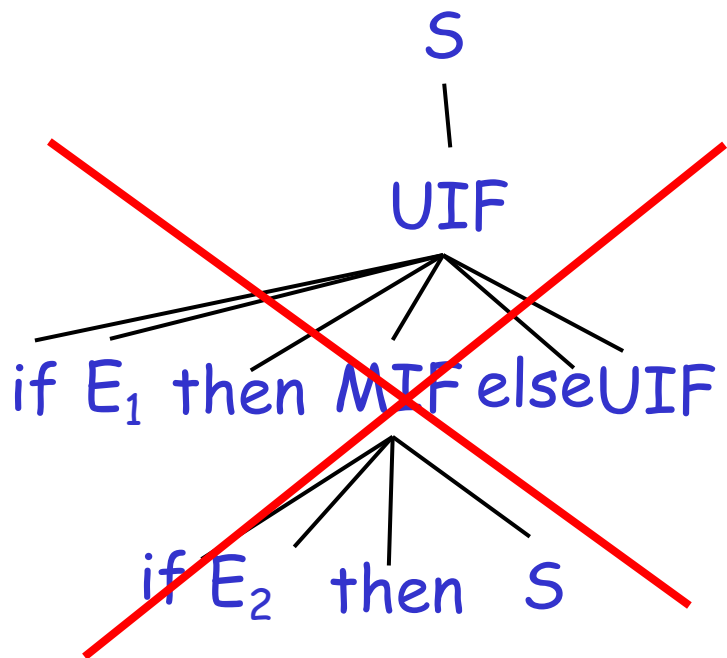
MIF \rightarrow if E then MIF else MIF
 | OTHER

UIF \rightarrow if E then S
 | if E then MIF else UIF

- Describes the same set of strings

The Dangling Else: Example Revisited

- The expression `if E1 then if E2 then S else S`



- Not valid because the `then` expression is not a `MIF`

Prof. Aiken [slightly modified]

- A valid parse tree (for a `UIF`)

Ambiguity

- No general techniques for handling ambiguity
- Impossible to convert automatically an ambiguous grammar to an unambiguous one
- Used with care, ambiguity can simplify the grammar
 - Sometimes allows more natural definitions
 - We need disambiguation mechanisms

Question?

Choose the alternative that correctly left factors “if” statements in the given grammar

$EXPR \rightarrow$ if BOOL then { EXPR }
 | if BOOL then { EXPR } else { EXPR }
 | ...
 $BOOL \rightarrow$ true | false

$EXPR \rightarrow$ if true then { EXPR }
 | if false then { EXPR }
 | if true then { EXPR } else { EXPR }
 | if false then { EXPR } else { EXPR }
 | ...

$EXPR \rightarrow$ EXPR' | EXPR' else { EXPR }
 $EXPR' \rightarrow$ if BOOL then { EXPR }
 | ...
 $BOOL \rightarrow$ true | false

$EXPR \rightarrow$ if BOOL EXPR'
 | ...
 $EXPR' \rightarrow$ then { EXPR }
 | then { EXPR } else { EXPR }
 $BOOL \rightarrow$ true | false

$EXPR \rightarrow$ if BOOL then { EXPR } EXPR'
 | ...
 $EXPR' \rightarrow$ else { EXPR } | ϵ
 $BOOL \rightarrow$ true | false

Question?

Choose the next parse state given the grammar, parse table, and current state below. The initial string is:

if true then { true } else { if false then { false } } \$

	if	then	else	{	}	true	false	\$
E	if Bthen { E } E'				ϵ	B	B	ϵ
E'			else { E }		ϵ			ϵ
B						true	false	

- | | Stack | Input |
|-----------------------|----------------------------|-------------------------------------|
| Current | E' \$ | else { if false then { false } } \$ |
| <input type="radio"/> | \$ | \$ |
| <input type="radio"/> | else { E } \$ | else { if false then { false } } \$ |
| <input type="radio"/> | E } \$ | if false then { false } } \$ |
| <input type="radio"/> | else {if Bthen {E} E' } \$ | else { if false then { false } } \$ |

$E \rightarrow \text{if B then } \{ E \} E' \mid B \mid \epsilon$
 $E' \rightarrow \text{else } \{ E \} \mid \epsilon$
 $B \rightarrow \text{true} \mid \text{false}$

Question?

For the given grammar, find the First and Follow of Non-terminals and the Parse table

$S \rightarrow i E t S S' \mid a$	First(S) =	Follow(S) =
$S' \rightarrow e S \mid \epsilon$	First(S') =	Follow(S') =
$E \rightarrow b$	First(E) =	Follow(E) =

	a	b	e	i	t	\$
S						
S'						
E						

Question?

For the given grammar,
find the First and Follow
of Non-terminals and
the Parse table

$E \rightarrow T E'$
 $E' \rightarrow + T E' \mid \epsilon$
 $T \rightarrow F T'$
 $T' \rightarrow * F T' \mid \epsilon$
 $F \rightarrow (E) \mid id$

First(E,T,F) =

First(E') =

First(T') =

Follow(E) =

Follow(T) =

Follow(F) =

Follow(E') =

Follow(T') =

	id	+	*	()	\$
E						
E'						
T						
T'						
F						

Question?

Which of the following statements are true about the given grammar?

$$S \rightarrow a T U b \mid \varepsilon$$

$$T \rightarrow c U c \mid b U b \mid a U a$$

$$U \rightarrow S b \mid c c$$

Choose all that are correct.

- The follow set of S is $\{ \$, b \}$
- The first set of U is $\{ a, b, c \}$
- The first set of S is $\{ \varepsilon, a, b \}$
- The follow set of T is $\{ a, b, c \}$

Question?

Consider the following grammar:

$$\begin{aligned} S &\rightarrow A (S) B \mid \varepsilon \\ A &\rightarrow S \mid S B x \mid \varepsilon \\ B &\rightarrow S B \mid y \end{aligned}$$

What are the first and follow sets of S

- First: $\{x, y, (, \varepsilon\}$ Follow: $\{y, x, (,)\}$
- First: $\{x, \varepsilon\}$ Follow: $\{\$, y, x, (,)\}$
- First: $\{y, (, \varepsilon\}$ Follow: $\{\$, y, (,)\}$
- First: $\{x, y, (, \varepsilon\}$ Follow: $\{\$, y, x, (,)\}$
- First: $\{x, y, (\}$ Follow: $\{\$, y, x, (,)\}$
- First: $\{x, (\}$ Follow: $\{\$, y, x\}$

Question?

Choose the grammar that correctly

eliminates left recursion from the given grammar: $E \rightarrow E + T \mid T$

$T \rightarrow id \mid (E)$

$E \rightarrow E + id \mid E + (E)$
 $\quad \mid id \mid (E)$

$E \rightarrow TE'$
 $E' \rightarrow + TE' \mid \varepsilon$
 $T \rightarrow id \mid (E)$

$E \rightarrow E' + T \mid T$
 $E' \rightarrow id \mid (E)$
 $T \rightarrow id \mid (E)$

$E \rightarrow id + E \mid E + T \mid T$
 $T \rightarrow id \mid (E)$

Question?

Consider the following grammar. Adding which one of the following rules will cause the grammar to be left-recursive?
[Choose all that apply]

$$S \rightarrow A$$

$$A \rightarrow B \mid C$$

$$B \rightarrow (C)$$

$$C \rightarrow B + C \mid D$$

$$D \rightarrow 1 \mid 0$$

- $D \rightarrow A$
- $A \rightarrow D$
- $B \rightarrow C$
- $D \rightarrow B$
- $C \rightarrow 1 C$

Question?

Which of the following grammars are ambiguous?

- $S \rightarrow SS \mid a \mid b$
- $E \rightarrow E + E \mid id$
- $S \rightarrow Sa \mid Sb$
- $E \rightarrow E' \mid E' + E$
 $E' \rightarrow -E' \mid id \mid (E)$

Question?

Choose the unambiguous version
of the given ambiguous grammar: $S \rightarrow SS | a | b | \varepsilon$

$S \rightarrow Sa | Sb | \varepsilon$

$S \rightarrow SS'$
 $S' \rightarrow a | b$

$S \rightarrow S | S'$
 $S' \rightarrow a | b$

$S \rightarrow Sa | Sb$

Question?

Consider the following grammar. How many unique parse trees are there for the string $5 * 3 + (2 * 7) + 4$?

- 2
- 1
- 7
- 8
- 5
- 4

$$E \rightarrow E * E \mid E + E \mid (E) \mid \text{int}$$