

# Error Handling

## Syntax-Directed Translation

### Recursive Descent Parsing

#### Lecture 6

Prof. Aiken CS 143 Lecture 6

1

## Announcements

- PA1
  - Due today at midnight
  - README, test case
  - Your name(s)!
- WA1
  - Due today at 5pm
- PA2
  - Assigned today
- WA2
  - Assigned Tuesday

Prof. Aiken CS 143 Lecture 6

2

## Outline

- Extensions of CFG for parsing
  - Precedence declarations
  - Error handling
  - Semantic actions
- Constructing a parse tree
- Recursive descent

Prof. Aiken CS 143 Lecture 6

3

## Error Handling

- Purpose of the compiler is
  - To detect non-valid programs
  - To translate the valid ones
- Many kinds of possible errors (e.g. in C)

Error kind	Example	Detected by ...
Lexical	... \$ ...	Lexer
Syntax	... x *% ...	Parser
Semantic	... int x; y = x(3); ...	Type checker
Correctness	your favorite program	Tester/User

Prof. Aiken CS 143 Lecture 6

4

## Syntax Error Handling

- Error handler should
  - Report errors accurately and clearly
  - Recover from an error quickly
  - Not slow down compilation of valid code
- Good error handling is not easy to achieve

Prof. Aiken CS 143 Lecture 6

5

## Approaches to Syntax Error Recovery

- From simple to complex
  - Panic mode
  - Error productions
  - Automatic local or global correction
- Not all are supported by all parser generators

Prof. Aiken CS 143 Lecture 6

6

### Error Recovery: Panic Mode

- Simplest, most popular method
- When an error is detected:
  - Discard tokens until one with a clear role is found
  - Continue from there
- Such tokens are called synchronizing tokens
  - Typically the statement or expression terminators

Prof. Aiken CS 143 Lecture 6

7

### Syntax Error Recovery: Panic Mode (Cont.)

- Consider the erroneous expression  
 $(1 + + 2) + 3$
- Panic-mode recovery:
  - Skip ahead to next integer and then continue
- Bison: use the special terminal `error` to describe how much input to skip  
 $E \rightarrow \text{int} \mid E + E \mid ( E ) \mid \text{error int} \mid ( \text{error} )$

Prof. Aiken CS 143 Lecture 6

8

### Syntax Error Recovery: Error Productions

- Idea: specify in the grammar known common mistakes
- Essentially promotes common errors to alternative syntax
- Example:
  - Write `5 x` instead of `5 * x`
  - Add the production  $E \rightarrow \dots \mid E E$
- Disadvantage
  - Complicates the grammar

Prof. Aiken CS 143 Lecture 6

9

### Error Recovery: Local and Global Correction

- Idea: find a correct "nearby" program
  - Try token insertions and deletions
  - Exhaustive search
- Disadvantages:
  - Hard to implement
  - Slows down parsing of correct programs
  - "Nearby" is not necessarily "the intended" program
  - Not all tools support it

Prof. Aiken CS 143 Lecture 6

10

### Syntax Error Recovery: Past and Present

- Past
  - Slow recompilation cycle (even once a day)
  - Find as many errors in one cycle as possible
  - Researchers could not let go of the topic
- Present
  - Quick recompilation cycle
  - Users tend to correct one error/cycle
  - Complex error recovery is less compelling
  - Panic-mode seems enough

Prof. Aiken CS 143 Lecture 6

11

### Abstract Syntax Trees

- So far a parser traces the derivation of a sequence of tokens
- The rest of the compiler needs a structural representation of the program
- Abstract syntax trees
  - Like parse trees but ignore some details
  - Abbreviated as AST

Prof. Aiken CS 143 Lecture 6

12

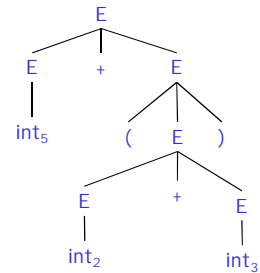
## Abstract Syntax Tree. (Cont.)

- Consider the grammar  
 $E \rightarrow \text{int} \mid (E) \mid E + E$
- And the string  
 $5 + (2 + 3)$
- After lexical analysis (a list of tokens)  
 $\text{int}_5 \text{'+' ' (' int}_2 \text{'+' int}_3 \text{' )'}$
- During parsing we build a parse tree ...

Prof. Aiken CS 143 Lecture 6

13

## Example of Parse Tree

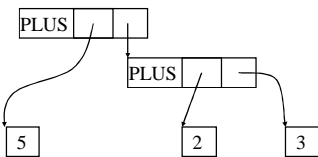


- Traces the operation of the parser
- Does capture the nesting structure
- But too much info
  - Parentheses
  - Single-successor nodes

Prof. Aiken CS 143 Lecture 6

14

## Example of Abstract Syntax Tree



- Also captures the nesting structure
- But abstracts from the concrete syntax  
 $\Rightarrow$  more compact and easier to use
- An important data structure in a compiler

Prof. Aiken CS 143 Lecture 6

15

## Semantic Actions

- This is what we'll use to construct ASTs
- Each grammar symbol may have attributes
  - For terminal symbols (lexical tokens) attributes can be calculated by the lexer
- Each production may have an action
  - Written as:  $X \rightarrow Y_1 \dots Y_n \quad \{ \text{action} \}$
  - That can refer to or compute symbol attributes

Prof. Aiken CS 143 Lecture 6

16

## Semantic Actions: An Example

- Consider the grammar  
 $E \rightarrow \text{int} \mid E + E \mid (E)$
- For each symbol  $X$  define an attribute  $X.\text{val}$ 
  - For terminals,  $\text{val}$  is the associated lexeme
  - For non-terminals,  $\text{val}$  is the expression's value (and is computed from values of subexpressions)
- We annotate the grammar with actions:
 

$E \rightarrow \text{int}$	$\{ E.\text{val} = \text{int.val} \}$
$\mid E_1 + E_2$	$\{ E.\text{val} = E_1.\text{val} + E_2.\text{val} \}$
$\mid (E_1)$	$\{ E.\text{val} = E_1.\text{val} \}$

Prof. Aiken CS 143 Lecture 6

17

## Semantic Actions: An Example (Cont.)

- String:  $5 + (2 + 3)$
- Tokens:  $\text{int}_5 \text{'+' ' (' int}_2 \text{'+' int}_3 \text{' )'}$

### Productions

$E \rightarrow E_1 + E_2$   
 $E_1 \rightarrow \text{int}_5$   
 $E_2 \rightarrow (E_3)$   
 $E_3 \rightarrow E_4 + E_5$   
 $E_4 \rightarrow \text{int}_2$   
 $E_5 \rightarrow \text{int}_3$

### Equations

$E.\text{val} = E_1.\text{val} + E_2.\text{val}$   
 $E_1.\text{val} = \text{int}_5.\text{val} = 5$   
 $E_2.\text{val} = E_3.\text{val}$   
 $E_3.\text{val} = E_4.\text{val} + E_5.\text{val}$   
 $E_4.\text{val} = \text{int}_2.\text{val} = 2$   
 $E_5.\text{val} = \text{int}_3.\text{val} = 3$

Prof. Aiken CS 143 Lecture 6

18

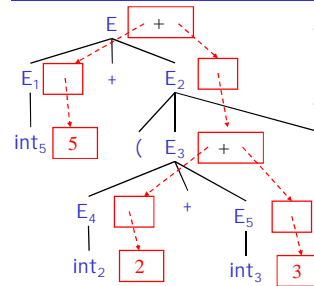
## Semantic Actions: Notes

- Semantic actions specify a system of equations
  - Order of resolution is not specified
- Example:
  - $E_3.val = E_4.val + E_5.val$
  - Must compute  $E_4.val$  and  $E_5.val$  before  $E_3.val$
  - We say that  $E_3.val$  depends on  $E_4.val$  and  $E_5.val$
- The parser must find the order of evaluation

Prof. Aiken CS 143 Lecture 6

19

## Dependency Graph



- Each node labeled  $E$  has one slot for the  $val$  attribute
- Note the dependencies

Prof. Aiken CS 143 Lecture 6

20

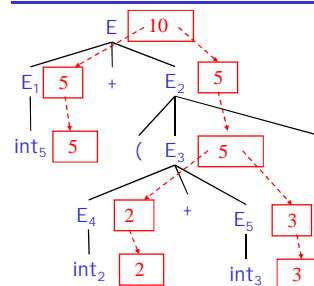
## Evaluating Attributes

- An attribute must be computed after all its successors in the dependency graph have been computed
  - In previous example attributes can be computed bottom-up
- Such an order exists when there are no cycles
  - Cyclically defined attributes are not legal

Prof. Aiken CS 143 Lecture 6

21

## Dependency Graph



Prof. Aiken CS 143 Lecture 6

22

## Semantic Actions: Notes (Cont.)

- Synthesized attributes
  - Calculated from attributes of descendants in the parse tree
  - $E.val$  is a synthesized attribute
  - Can always be calculated in a bottom-up order
- Grammars with only synthesized attributes are called S-attributed grammars
  - Most common case

Prof. Aiken CS 143 Lecture 6

23

## Inherited Attributes

- Another kind of attribute
- Calculated from attributes of parent and/or siblings in the parse tree
- Example: a line calculator

Prof. Aiken CS 143 Lecture 6

24

## A Line Calculator

- Each line contains an expression  
 $E \rightarrow \text{int} \mid E + E$
- Each line is terminated with the = sign  
 $L \rightarrow E = \mid + E =$
- In second form the value of previous line is used as starting value
- A program is a sequence of lines  
 $P \rightarrow \varepsilon \mid P L$

Prof. Aiken CS 143 Lecture 6

25

## Attributes for the Line Calculator

- Each  $E$  has a synthesized attribute  $val$ 
  - Calculated as before
- Each  $L$  has an attribute  $val$ 
  - $L \rightarrow E = \{ L.val = E.val \}$
  - $\mid + E = \{ L.val = E.val + L.prev \}$
- We need the value of the previous line
- We use an inherited attribute  $L.prev$

Prof. Aiken CS 143 Lecture 6

26

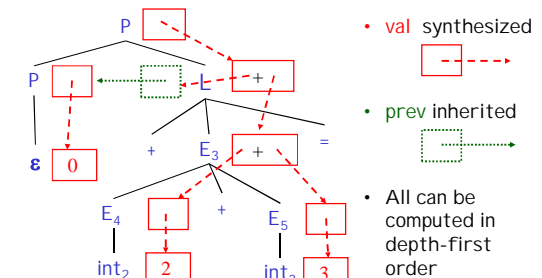
## Attributes for the Line Calculator (Cont.)

- Each  $P$  has a synthesized attribute  $val$ 
  - The value of its last line
    - $P \rightarrow \varepsilon \{ P.val = 0 \}$
    - $\mid P_1 L \{ P.val = L.val; L.prev = P_1.val \}$
  - Each  $L$  has an inherited attribute  $prev$
  - $L.prev$  is inherited from sibling  $P_1.val$
- Example ...

Prof. Aiken CS 143 Lecture 6

27

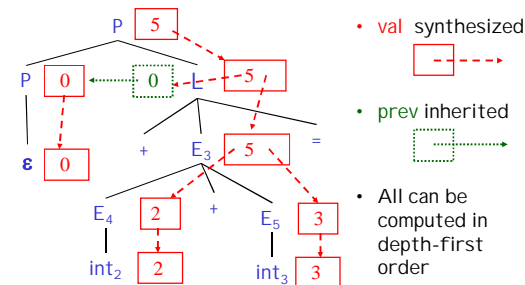
## Example of Inherited Attributes



Prof. Aiken CS 143 Lecture 6

28

## Example of Inherited Attributes



Prof. Aiken CS 143 Lecture 6

29

## Semantic Actions: Notes (Cont.)

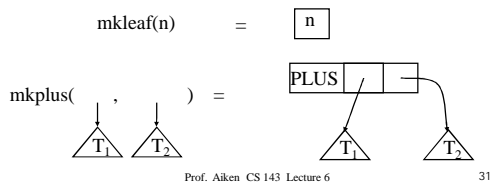
- Semantic actions can be used to build ASTs
- And many other things as well
  - Also used for type checking, code generation, ...
- Process is called syntax-directed translation
  - Substantial generalization over CFGs

Prof. Aiken CS 143 Lecture 6

30

## Constructing An AST

- We first define the AST data type
  - Supplied by us for the project
- Consider an abstract tree type with two constructors:



## Constructing a Parse Tree

- We define a synthesized attribute `ast`
  - Values of `ast` values are ASTs
  - We assume that `int.lexval` is the value of the integer lexeme
  - Computed using semantic actions

$E \rightarrow \text{int} \quad E.\text{ast} = \text{mkleaf}(\text{int.lexval})$   
 $E \rightarrow E_1 + E_2 \quad E.\text{ast} = \text{mkplus}(E_1.\text{ast}, E_2.\text{ast})$   
 $E \rightarrow (E_1) \quad E.\text{ast} = E_1.\text{ast}$

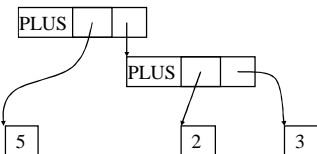
Prof. Aiken CS 143 Lecture 6

32

## Parse Tree Example

- Consider the string `int5 '+' '(' int2 '+' int3 ')'`
- A bottom-up evaluation of the `ast` attribute:

$E.\text{ast} = \text{mkplus}(\text{mkleaf}(5),$   
 $\quad \text{mkplus}(\text{mkleaf}(2), \text{mkleaf}(3)))$



Prof. Aiken CS 143 Lecture 6

33

## Summary

- We can specify language syntax using CFG
- A parser will answer whether  $s \in L(G)$ 
  - ... and will build a parse tree
  - ... which we convert to an AST
  - ... and pass on to the rest of the compiler

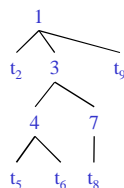
Prof. Aiken CS 143 Lecture 6

34

## Intro to Top-Down Parsing: The Idea

- The parse tree is constructed
  - From the top
  - From left to right
- Terminals are seen in order of appearance in the token stream:

$t_2 \ t_5 \ t_6 \ t_8 \ t_9$



Prof. Aiken CS 143 Lecture 6

35

## Recursive Descent Parsing

- Consider the grammar
 
$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Token stream is: `int5 * int2`
- Start with top-level non-terminal `E`
  - Try the rules for `E` in order

Prof. Aiken CS 143 Lecture 6

36

### Recursive Descent Parsing. Example (Cont.)

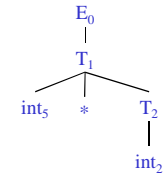
- Try  $E_0 \rightarrow T_1 + E_2$ 
  - Then try a rule for  $T_1 \rightarrow ( E_3 )$ 
    - But  $($  does not match input token  $int_5$
  - Try  $T_1 \rightarrow int$ . Token matches.
    - But  $+$  after  $T_1$  does not match input token  $*$
  - Try  $T_1 \rightarrow int * T_2$ 
    - This will match but  $+$  after  $T_1$  will be unmatched
- We've exhausted the choices for  $T_1$ 
  - **Failure!**
  - Backtrack, try another choice for  $E_0$

Prof. Aiken CS 143 Lecture 6

37

### Recursive Descent Parsing. Example (Cont.)

- Try  $E_0 \rightarrow T_1$
- Follow same steps as before for  $T_1$ 
  - And succeed with  $T_1 \rightarrow int * T_2$  and  $T_2 \rightarrow int$
  - With the following parse tree



Prof. Aiken CS 143 Lecture 6

38

### A Recursive Descent Parser. Preliminaries

- Let TOKEN be the type of tokens
  - Special tokens INT, OPEN, CLOSE, PLUS, TIMES
- Let the global `next` point to the next token

Prof. Aiken CS 143 Lecture 6

39

### A Recursive Descent Parser (2)

- Define boolean functions that check the token string for a match of
  - A given token terminal  
`bool term(TOKEN tok) { return *next++ == tok; }`
  - A given production of S (the  $n^{\text{th}}$ )  
`bool Sn() { ... }`
  - Any production of S:  
`bool S() { ... }`

Prof. Aiken CS 143 Lecture 6

40

### A Recursive Descent Parser (3)

- For production  $E \rightarrow T$   
`bool E1() { return T(); }`
- For production  $E \rightarrow T + E$   
`bool E2() { return T() && term(PLUS) && E(); }`
- For all productions of E (with backtracking)  
`bool E() {  
 TOKEN *save = next;  
 return (next = save, E1())  
 || (next = save, E2());  
}`

Prof. Aiken CS 143 Lecture 6

41

### A Recursive Descent Parser (4)

- Functions for non-terminal T  
`bool T1() { return term(OPEN) && E() && term(CLOSE); }`  
`bool T2() { return term(INT) && term(TIMES) && T(); }`  
`bool T3() { return term(INT); }`  
  
`bool T() {  
 TOKEN *save = next;  
 return (next = save, T1())  
 || (next = save, T2())  
 || (next = save, T3());  
}`

Prof. Aiken CS 143 Lecture 6

42

## Recursive Descent Parsing. Notes.

- To start the parser
  - Initialize `next` to point to first token
  - Invoke `E()`
- Notice how this simulates the example parse
- Easy to implement by hand
- But does not always work ...

Prof. Aiken CS 143 Lecture 6

43

## Example

```
E → T + E | T
T → int | int * T | ( E )

int * int

bool term(TOKEN tok) { return *next++ == tok; }

bool E_0() { return T_0(); }
bool E_1() { return T_0() && term(PLUS) && E_0(); }

bool E() { TOKEN *save = next; return (next = save, E_0())
|| (next = save, E_1()); }

bool T_0() { return term(OPEN) && E() && term(CLOSE); }
bool T_1() { return term(INT) && term(TIMES) && T_0(); }
bool T_2() { return term(INT); }

bool T() { TOKEN *save = next; return (next = save, T_0())
|| (next = save, T_1())
|| (next = save, T_2()); }
```

Prof. Aiken CS 143 Lecture 6

44

## When Recursive Descent Does Not Work

- Consider a production  $S \rightarrow S a$   
`bool S_1() { return S() && term(a); }`  
`bool S() { return S_1(); }`
- `S()` goes into an infinite loop
- A left-recursive grammar has a non-terminal  $S$   
 $S \rightarrow^+ S \alpha$  for some  $\alpha$
- Recursive descent does not work in such cases

Prof. Aiken CS 143 Lecture 6

45

## Elimination of Left Recursion

- Consider the left-recursive grammar  
 $S \rightarrow S \alpha \mid \beta$
- $S$  generates all strings starting with a  $\beta$  and followed by a number of  $\alpha$
- Can rewrite using right-recursion  
 $S \rightarrow \beta S'$   
 $S' \rightarrow \alpha S' \mid \epsilon$

Prof. Aiken CS 143 Lecture 6

46

## 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  $\beta_1, \dots, \beta_m$  and continue with several instances of  $\alpha_1, \dots, \alpha_n$
- Rewrite as  
 $S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$   
 $S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \epsilon$

Prof. Aiken CS 143 Lecture 6

47

## 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

Prof. Aiken CS 143 Lecture 6

48



## Summary of Recursive Descent

---

- Simple and general parsing strategy
  - Left-recursion must be eliminated first
  - ... but that can be done automatically
- Unpopular because of backtracking
  - Thought to be too inefficient
- In practice, backtracking is eliminated by restricting the grammar