

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 = <u>x</u> (3); ...	Type checker
Correctness	your favorite program	<u>Tester/User</u>

Error handler should:

- ① Report errors accurately and clearly.
- ② Recover from an error quickly.
- ③ Not slow down compilation of valid code

Three methods:

- ① Panic mode → present
- ② Error productions
- ③ Automatic local or global correction → past

Panic mode is simplest, most popular method

When an error is detected:

- Discard tokens until one with a clear role is found
- Continue from there

Looking for synchronizing tokens

- Typically the statement or expression terminators

- 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



Error productions specify known common mistakes in the grammar.

Example:

Write S^x instead of S^*x

Add the production $E \rightarrow \dots | EEE$

Disadvantage

Complicates the grammar.

(warning but not error)

Error Correction

Idea: find a correct "nearby" program

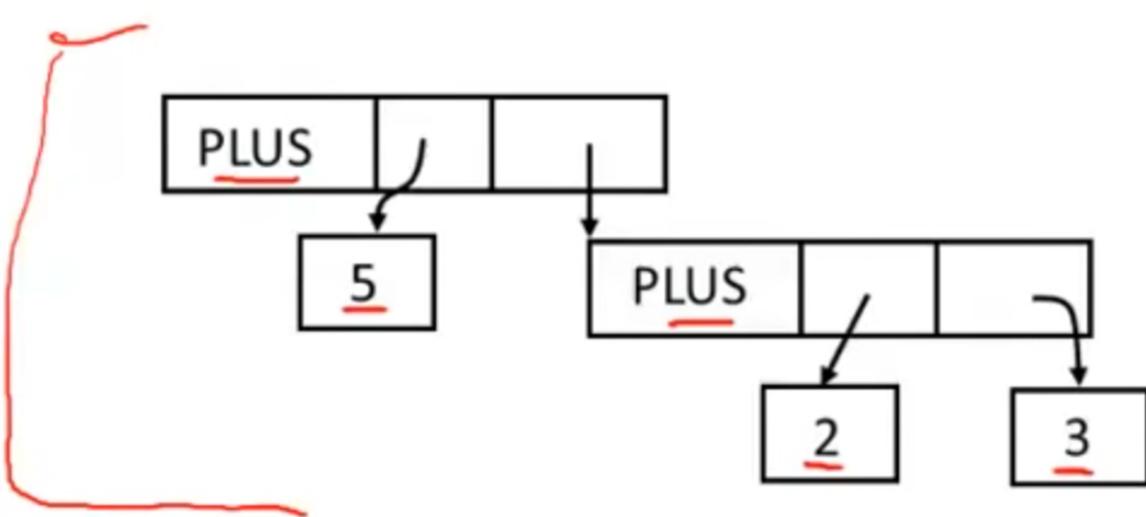
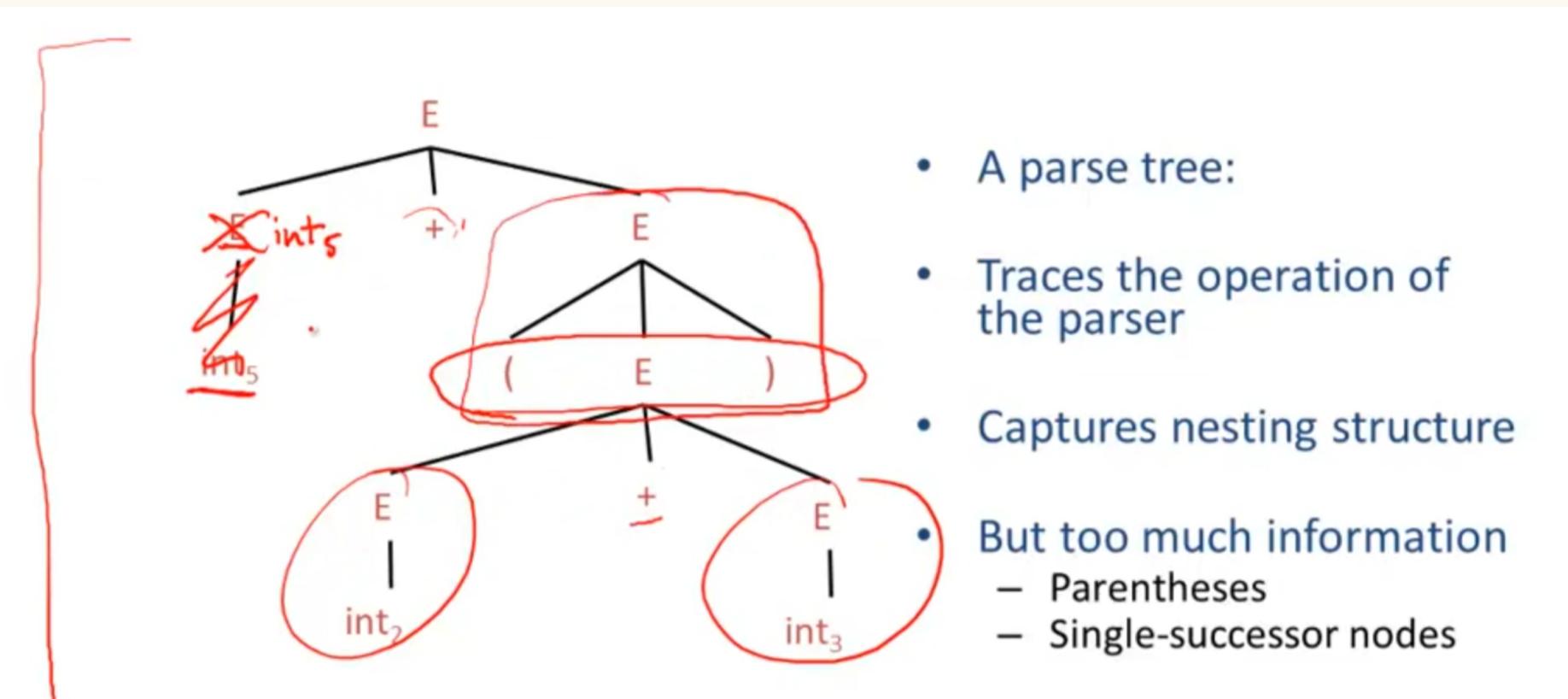
- Try token insertions and deletions
- Exhaustive search

Disadvantages:

- Hard to implement
- Slows down parsing of correct programs
- "Nearby" is not necessarily "the intended" program

Abstract Syntax Trees

AST like parse trees but ignore some details.



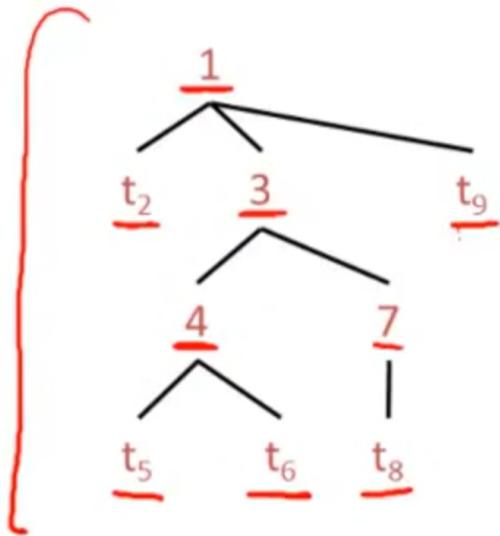
- Also captures the nesting structure
- But abstracts from the concrete syntax
=> more compact and easier to use
- An important data structure in a compiler

Recursive Descent Parsing

Top-down

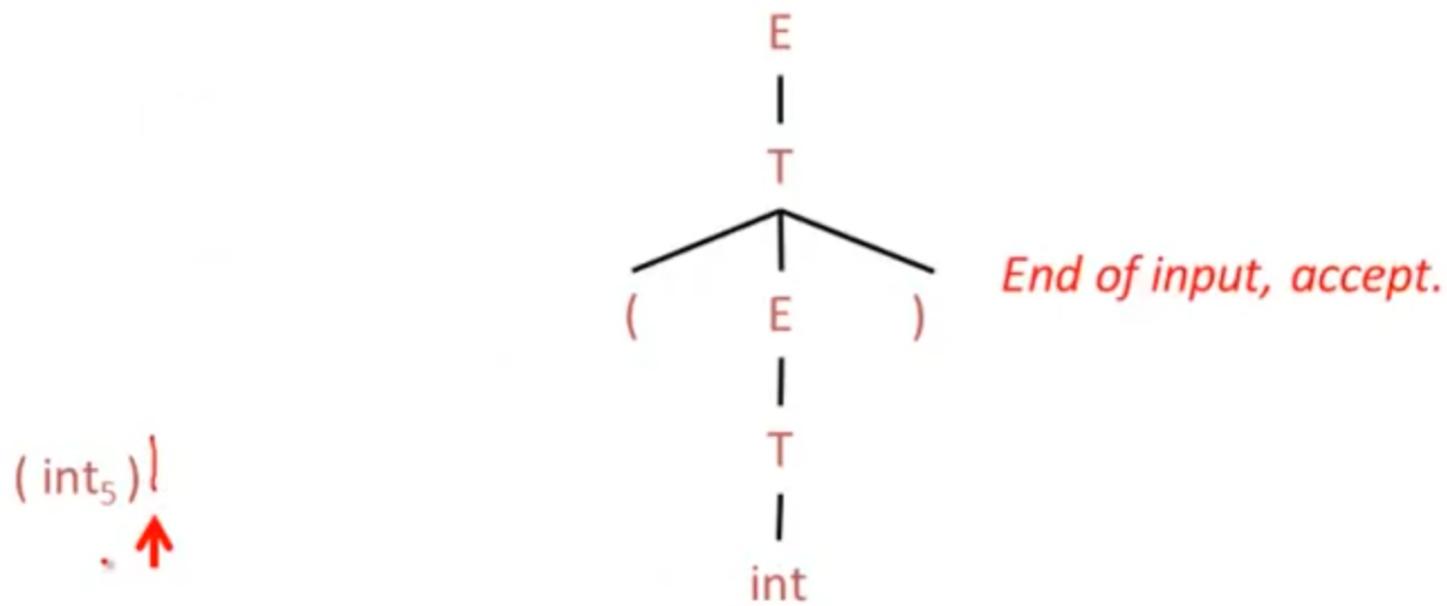
Recursive Descent

- The parse tree is constructed
 - From the top
 - From left to right
- Terminals are seen in order of appearance in the token stream:
t₂ t₅ t₆ t₈ t₉



$E \rightarrow T \mid T + E$

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



Recursive Descent Algorithm

Let TOKEN be the type of tokens

- Special tokens INT, OPEN, CLOSE, PLUS, TIMES.

Let the global next point to the next input token

- Define boolean functions that check for a match of:

- A given token terminal

bool term(TOKEN tok) { return *next == tok; } true/false
↓

- The nth production of S:

→ bool S_n() { ... }

- Try all productions of S:

→ bool S() { ... }

- For production E → T

bool E₁() { return T(); }

- For production E → 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())

(II) (next = save, E₂()) }

- Functions for non-terminal T

```
bool T1() { return term(INT); }  $T \rightarrow \text{int}^+$ 
bool T2() { return term(INT) && term(TIMES) && T(); }  $T \rightarrow \text{int} * T$ 
bool T3() { return term(OPEN) && E() && term(CLOSE); }  $T \rightarrow ( E )$ 
```

5 seconds

```
bool T() {
    TOKEN *save = next;
    return (next = save, T1())  

        || (next = save, T2())  

        || (next = save, T3()); }
```

To start the parser:

- Initialize next to point to first token
- Invoke E()

Easy to implement by hand.

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



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

bool E1() { return T(); }
bool E2() { return T() && term(PLUS) && E(); }

bool E() { TOKEN *save = next; return (next = save, E1())  

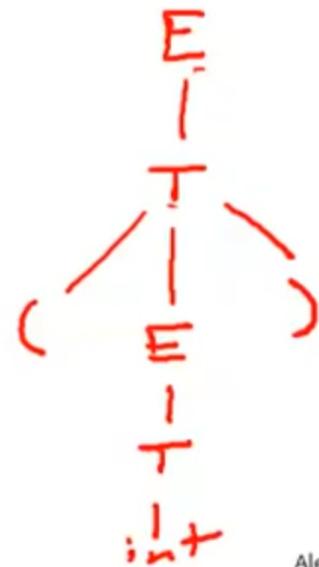
           || (next = save, E2()); }

bool T1() { return term(INT); }
bool T2() { return term(INT) && term(TIMES) && T(); }
bool T3() { return term(OPEN) && E() && term(CLOSE); }

bool T() { TOKEN *save = next; return (next = save, T1())  

           || (next = save, T2())  

           || (next = save, T3()); }
```



$$\begin{cases} E \rightarrow T \mid T + E \\ T \rightarrow \text{int} \mid \underline{\text{int}} * T \mid (E) \end{cases}$$

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

bool E₁() { return T(); }

bool E₂() { return T() && term(PLUS) && E(); }

int * int rejected

What happened?!

bool E() { TOKEN *save = next; return (next = save, E₁()) || (next = save, E₂()); }

bool T₁() { return term(INT); }

bool T₂() { return term(INT) && term(TIMES) && T(); }

bool T₃() { return term(OPEN) && E() && term(CLOSE); }

bool T() { TOKEN *save = next; return (next = save, T₁()) || (next = save, T₂()) || (next = save, T₃()); }

If a production for non-terminal X succeeds

- Cannot backtrack to try a different production for X later.

-

General recursive-descent algorithms support such "full" backtracking

- Can implement any grammar.

Presented recursive-descent algorithm is not general

- But is easy to implement by hand

Sufficient for grammars where for any non-terminal at most one production can succeed

The example grammar can be written to work with the presented algorithm: By left factoring

Left Recursion

- Consider a production $S \rightarrow S a$

~~bool S₁() { return S() && term(a); }~~
~~bool S(){ return S₁(); }~~

S → S_a → S_aa →
S_aa_a → ... → S_a...a →

- `S()` goes into an infinite loop

- A left-recursive grammar has a non-terminal S

$S \rightarrow^+ Sa$ for some α

- Recursive descent does not work in such cases

- Consider the left-recursive grammar

$$S \rightarrow S \underline{\alpha} \mid \underline{\beta}$$

$$S \xrightarrow{} S\underline{\alpha} | \underline{\beta}$$

$$S \xrightarrow{} S\alpha \rightarrow S\alpha\alpha \rightarrow S\alpha\alpha\alpha \rightarrow \dots \rightarrow S\alpha\dots\alpha \rightarrow \underline{\beta\alpha\dots\alpha}$$

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

- Can rewrite using right-recursion

$$S \rightarrow \beta S'$$

$$S' \rightarrow \underline{\alpha S' \mid \varepsilon}$$

$S \rightarrow BS' \rightarrow B\alpha S' \rightarrow B\alpha\alpha S' \rightarrow \dots$
 $\qquad\qquad\qquad \rightarrow B\alpha\dots\alpha S' \rightarrow \underline{B\alpha\dots\alpha}$

- In general

$$S \rightarrow \underline{S} \alpha_1 | \dots | \underline{S} \alpha_n | \underline{\beta_1} | \dots | \underline{\beta_m}$$

- All strings derived from \underline{S} start with one of $\underline{\beta_1}, \dots, \underline{\beta_m}$ and continue with several instances of $\underline{\alpha_1}, \dots, \underline{\alpha_n}$

- Rewrite as

$$\left[\begin{array}{l} S \rightarrow \underline{\beta_1} S' | \dots | \underline{\beta_m} S' \\ S' \rightarrow \underline{\alpha_1} S' | \dots | \underline{\alpha_n} S' | \varepsilon \end{array} \right]$$

- The grammar

$$S \rightarrow \underline{A} \alpha | \delta$$

$$\underline{A} \rightarrow \underline{S} \beta$$

is also left-recursive because

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

$$\underline{S} \rightarrow A \alpha \rightarrow \underline{S} \beta \alpha$$

- This left-recursion can also be eliminated
- See Dragon Book for general algorithm
- Recursive descent
 - Simple and general parsing strategy
 - Left-recursion must be eliminated first
 - ... but that can be done automatically
- Used in production compilers
 - E.g., gcc