- The parse tree is constructed
  - From the top
  - From left to right

- Terminals are seen in order of appearance in the token stream:

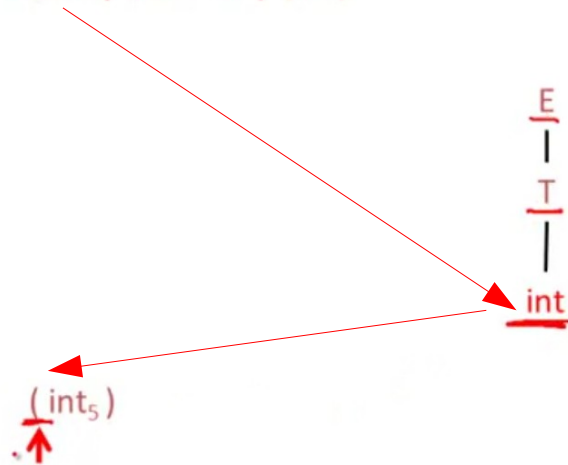  $t_2 \ t_5 \ t_6 \ t_8 \ t_9$

- Consider the grammar

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

- Token stream is:   ( $\text{int}_5$ )

- Start with top-level non-terminal $E$
  - Try the rules for $E$ in order

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid ( E )$

$E$
$|$
$T$
$|$
$int$

*Mismatch: int does not match (*
*Backtrack …*

$( int_5 )$

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid ( E )$

$E$
$|$
$T$
$int$     $*$     $T$

*Mismatch: int does not match (*
*Backtrack …*

$( int_5 )$

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid (E)$

```
        E
        |
        T
       /|\
      ( E )        Match!  Advance input.
```

( int₅ )

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid (E)$

```
        E
        |
        T
       /|\
      ( E )        Match!  Advance input.
        |
        T
        |
       int
```

( int₅ )

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid (E)$

E
|
T
/ | \
(   E   )
    |
    T
    |
    int

*Match! Advance input.*

( int$_5$ )

---

$E \rightarrow T \mid T + E$

$T \rightarrow int \mid int * T \mid (E)$

E
|
T
/ | \
(   E   )
    |
    T
    |
    int

*End of input, accept.*

( int$_5$ )

```
/* recognize tokens for the calculator and print them out */

%{
    enum yytokentype {
        NUMBER = 258,
        ADD = 259,
        SUB = 260,
        MUL = 261,
        DIV = 262,
        ABS = 263,
        EOL = 264 /* end of line */
    };
```

Enumeration (or enum) is a user defined data type in C. It is mainly used to assign names to integral constants, the names make a program easy to read and maintain.

- 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; }
```

advances next, returns boolean

$$E \rightarrow T \mid T + E$$
$$T \rightarrow int \mid int * T \mid ( E )$$

– The nth production of S:  bool $S_n()$ { ... }

– Try all productions of S:  bool $S()$ { ... }

- For production $E \rightarrow T$  bool $E_1()$ { return T(); }

these advance next

- For production $E \rightarrow T + E$

bool $E_2()$ { return T() && term(PLUS) && E(); }

&& evaluates arguments in left to right order

$$E \rightarrow T \mid T + E$$
$$T \rightarrow int \mid int * T \mid ( E )$$

- For all productions of E (with backtracking)

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

|| if first branch succeeds, do not bother with second branch

backtracking

if they all fail, the higher level will do the backtracking

$$E \rightarrow T \mid T + E$$
$$T \rightarrow int \mid int * T \mid ( E )$$

- Functions for non-terminal T

```
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()); }
```

# A limitation of recursive descent

$$E \rightarrow T \mid T + E$$
$$T \rightarrow int \mid int * T \mid ( E )$$

```
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, ˆT₃()); }
```

int * int will be rejected

once a non terminal succeeds
no way to try another production

- 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 rewritten to work with the presented algorithm
  - By *left factoring*

# Left Recursion

In the formal language theory of computer science, left recursion is a special case of recursion where a string is recognized as part of a language by the fact that it decomposes into a string from that same language (on the left) and a suffix (on the right).

- Consider a production $S \rightarrow S\ a$

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

    ```
    bool S() { return  S₁(); }
    ```

non empty sequence of rewrites

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

    $$S \rightarrow^+ S\alpha \quad \text{for some } \alpha$$

- Recursive descent does not work in such cases

- Consider the left-recursive grammar
$$S \rightarrow S\alpha \mid \beta$$

$$S \rightarrow S\alpha \rightarrow S\alpha\alpha \rightarrow S\alpha\alpha\alpha \rightarrow \cdots \rightarrow S\alpha \cdots \alpha \rightarrow \beta\alpha \cdots \alpha$$

- S generates all strings starting with a $\beta$ and followed by any number of $\alpha$'s

zero or more

- Can rewrite using right-recursion
$$S \rightarrow \beta S'$$

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

$$S \rightarrow \beta S' \rightarrow \beta \alpha S' \rightarrow \beta \alpha\alpha S' \rightarrow \cdots$$
$$\rightarrow \beta\alpha \cdots \alpha S' \rightarrow \beta\alpha \cdots \alpha$$

- 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 $\alpha_1,\dots,\alpha_n$

zero or more

- Rewrite as

$$S \rightarrow \beta_1\,S' \mid \dots \mid \beta_m\,S'$$
$$S' \rightarrow \alpha_1\,S' \mid \dots \mid \alpha_n\,S' \mid \varepsilon$$

- 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

- 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

# Predictive Parsers

- Like recursive-descent but parser can "predict" which production to use
  - By looking at the next few tokens
  - No backtracking

  *lookahead restricted grammars*

- Predictive parsers accept LL(k) grammars

  *left-to-right* *k tokens lookahead .* *left-most derivation*

  always k=1

A deterministic model of computation is a model of computation such that the successive states of the machine and the operations to be performed are completely determined by the preceding state.

- In recursive descent,
  - At each step, many choices of production to use
  - Backtracking used to undo bad choices

- In LL(1),
  - At each step, only one choice of production

LL(1) grammars cannot be left recursive since the leftmost nonterminal is the same as the LHS. This would result in infinite recursion.

Left factoring is removing the common left factor that appears in two productions of the same non-terminal.

It is done to avoid back-tracing by the parser.

Suppose the parser has a look-ahead consider this example

A -> qB | qC
where A,B,C are non-terminals and q is a sentence. In this case, the parser will be confused as to which of the two productions to choose and it might have to back-trace.
After left factoring, the grammar is converted to

A -> qD

D -> B | C

In this case, a parser with a look-ahead will always choose the right production.

- Recall the grammar

$$E \rightarrow T + E \mid T$$
$$T \rightarrow int \mid int * T \mid (E)$$

- 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

commom prefix

$$E \rightarrow T X$$

$$E \rightarrow T X$$
$$X \rightarrow + E \mid \epsilon$$

suffixes

$$T \rightarrow int Y \mid (E)$$
$$Y \rightarrow * T \mid \epsilon$$

- Left-factored grammar

$$E \rightarrow T\ X \qquad\qquad X \rightarrow +\ E \mid \varepsilon$$
$$T \rightarrow (\ E\ ) \mid int\ Y \qquad\qquad Y \rightarrow *\ T \mid \varepsilon$$

- The LL(1) parsing table:

*next input token*

|   | int | * | + | ( | ) | $ |
|---|-----|---|---|---|---|---|
| E | T X |   |   | T X |   |   |
| X |     |   | + E |   | ε | ε |
| T | int Y |   |   | ( E ) |   |   |
| Y |     | * T | ε |   | ε | ε |

*leftmost non-terminal*

*rhs of production to use*

- Consider the [E, int] entry
  - "When current non-terminal is E and next input is int, use production $E \rightarrow T\ X$"

|   | int | * | + | ( | ) | $ |
|---|-----|---|---|---|---|---|
| E | T X |   |   | T X |   |   |
| X |     |   | + E |   | ε | ε |
| T | int Y |   |   | ( E ) |   |   |
| Y |     | * T | ε |   | ε | ε |

- 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 → ε

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | ε | ε |
| T | int Y | | | ( E ) | | |
| Y | | * T | ε | | ε | ε |

- Consider the [E,*] entry
  - "There is no way to derive a string starting with * from non-terminal E"

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | ε | ε |
| T | int Y | | | ( E ) | | |
| Y | | * T | ε | | ε | ε |

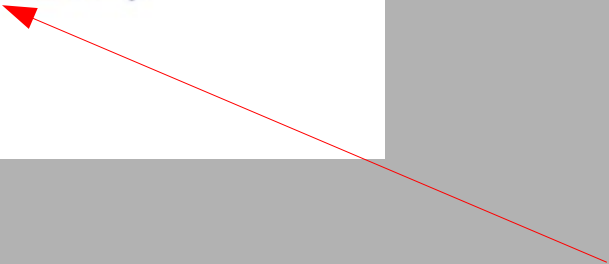- Method similar to recursive descent, except
  - For the leftmost non-terminal S
  - We look at the next input token a
  - And choose the production shown at [S,a]

- A stack records frontier of parse tree
  - Non-terminals that have yet to be expanded
  - Terminals that have yet to matched against the input
  - Top of stack = leftmost pending terminal or non-terminal

- Reject on reaching error state
- Accept on end of input & empty stack

```
initialize stack = <Ṡ $> and next
repeat
  case stack of
    <X, rest>  : if T[X,*next] = Y₁...Yₙ
                   then stack ← <Y₁... Yₙ rest>;
                   else  error ();
    <t, rest>   : if t == *next ++
                    then  stack ← <rest>;
                    else error ();
until stack == < >
```

$$\text{initialize stack} = \langle \dot{S}\ \$ \rangle \text{ and next}$$

terminal on top of stack
matches input, pop and advance
input

- Left-factored grammar

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

$$T \rightarrow ( E ) \mid int\ Y \qquad Y \rightarrow * T \mid \varepsilon$$

- The LL(1) parsing table:

*next input token*

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | $\varepsilon$ | $\varepsilon$ |
| T | int Y | | | ( E ) | | |
| Y | | * T | $\varepsilon$ | | $\varepsilon$ | $\varepsilon$ |

*leftmost non-terminal*

*rhs of production to use*

| Stack | Input | Action |
|---|---|---|
| E $ | int * int $ | T X |
| T X $ | int * int $ | int Y |

- ## The LL(1) parsing table:

*next input token*

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | ε | ε |
| T | int Y | | | ( E ) | | |
| Y | | * T | ε | | ε | ε |

*leftmost non-terminal*

*rhs of production to use*

| | | |
|---|---|---|
| int Y X $ | int * int $ | terminal |
| | | |
| Y X $ | * int $ | * T |



| | | |
|---|---|---|
| * T X $ | * int $ | terminal |
| | | |
| T X $ | int $ | int Y |

- Left-factored grammar

$$E \rightarrow T\,X \qquad\qquad X \rightarrow +\,E \mid \varepsilon$$
$$T \rightarrow (\,E\,) \mid int\ Y \qquad\qquad Y \rightarrow *\,T \mid \varepsilon$$

- The LL(1) parsing table:   *next input token*

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | ε | ε |
| T | int Y | | | ( E ) | | |
| Y | | * T | ε | | ε | ε |

*leftmost non-terminal*   *rhs of production to use*

| int Y X $ | int $ | terminal |
|---|---|---|
| Y X $ | $ | ε |
| X $ | $ | ε |

– "When current non-terminal is Y and current token is +, get rid of Y"

- Left-factored grammar

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

$T \rightarrow ( E ) \mid int\ Y$   $Y \rightarrow * T \mid \varepsilon$

- The LL(1) parsing table:   *next input token*

| | int | * | + | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | T X | | | T X | | |
| X | | | + E | | $\varepsilon$ | $\varepsilon$ |
| T | int Y | | | ( E ) | | |
| Y | | * T | $\varepsilon$ | | $\varepsilon$ | $\varepsilon$ |

*leftmost non-terminal*   *rhs of production to use*

X $   $   $\varepsilon$

$   $   ACCEPT