

Top-Down Parsing

CIS*4650 (Winter 2020)

Parsing Algorithms

● Top-down parsing:

- Parse input tokens by tracing out steps in a leftmost derivation
- Predictive parsers: e.g., recursive descent parsing, LL(k)
- Backtracking parsers: more powerful but much slower and may run into infinite loops
- Relatively easy to write by hand

● Bottom-up parsing:

- Parse input tokens by performing reductions to uncover the parsing steps in a rightmost derivation
- E.g., LR(k), LALR(k)
- Most parser generators use bottom-up parsing

Recursive Descent Parsing

○ Basic idea:

- Construct a parse tree from root and create the nodes of the parse tree in the depth-first order
- Non-terminals become recursive functions
- Productions become clauses in the non-terminal function
- Clauses implement explicit calls to consume the expected tokens
- Parse errors are generated if the expected tokens can't be consumed

e.g., `STMT -> if EXPR then STMT else STMT`
 `| while EXPR do STMT`
 `| begin STMT_LIST end`

Recursive Descent Parsing

```
enum token_t { IF, THEN, ELSE, WHILE, DO, BEGIN, END };
```

```
extern enum token_t getNextToken();
```

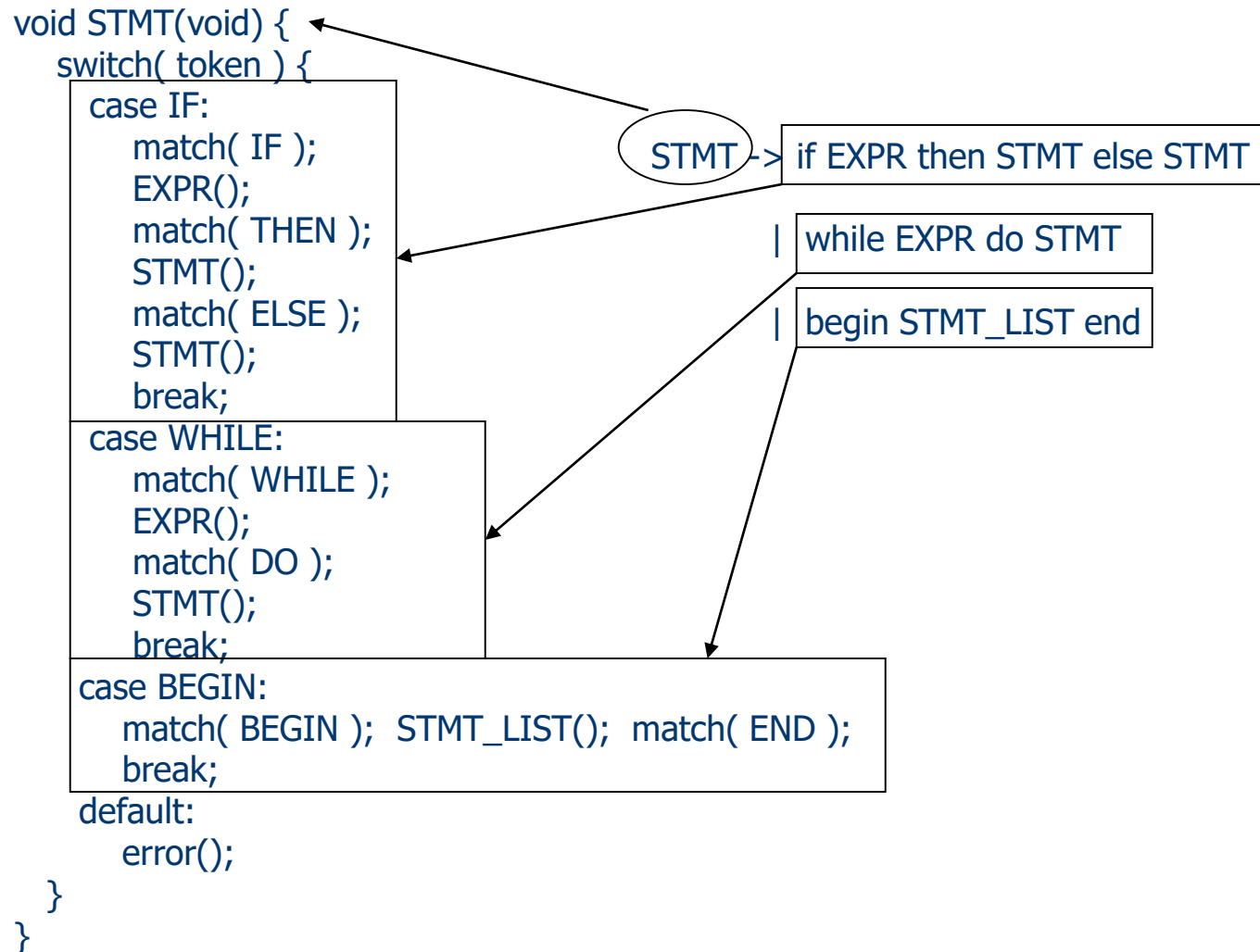
```
enum token_t token;
```

```
void advance() { token = getNextToken(); }
```

```
void match(enum token_t candidate) {  
    if( candidate == token )  
        advance();  
    else  
        error();  
}
```

```
void error(void) { fprintf(stderr, "Error\n"); }
```

Recursive Descent Parsing



Recursive Descent Parsing

● Handling choices:

$\langle \text{if-stmt} \rangle \rightarrow \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle$
| $\text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

```
void if-stmt( void ) {  
    match( IF );  
    match( LPAREN );  
    exp();  
    match( RPAREN );  
    stmt();  
    if( token == ELSE ) {  
        match( ELSE );  
        stmt();  
    }  
}
```

In EBNF:

$\langle \text{if-stmt} \rangle \rightarrow \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle [\text{ else } \langle \text{stmt} \rangle]$

Recursive Descent Parsing

● Handling repetitions:

$\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle \langle \text{addop} \rangle \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\langle \text{addop} \rangle \rightarrow + \mid -$

```
void exp( void ) {  
    term();  
    while( token == PLUS ||  
           token == MINUS ) {  
        match( token );  
        term();  
    }  
}
```

In EBNF:

$\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle \{ \langle \text{addop} \rangle \langle \text{term} \rangle \}$
 $\langle \text{addop} \rangle \rightarrow + \mid -$

Recursive Descent Parsing

- Not always easy to convert grammar in BNF to that in EBNF
- May be difficult to decide which case to use if they all start with non-terminals
- Need to handle $A \rightarrow \varepsilon$ productions
- Early detect errors if possible: e.g., “)3 - 2)”

LL(1) Parsing

- LL(1): Left-to-right parse, Leftmost derivation, and 1 symbol lookahead
- Use an explicit stack rather than recursive calls

e.g.: $S \rightarrow (S) S \mid \epsilon$

| Steps | Parsing stack | Input | Action |
|-------|---------------|--------|--------------------------|
| 1 | \$ S | () \$ | $S \rightarrow (S) S$ |
| 2 | \$ S) S (| () \$ | match |
| 3 | \$ S) S |) \$ | $S \rightarrow \epsilon$ |
| 4 | \$ S) |) \$ | match |
| 5 | \$ S | \$ | $S \rightarrow \epsilon$ |
| 6 | \$ | \$ | accept |

LL(1) Parsing

● Generation action:

- Replace a non-terminal A at the top of the stack by a string α using the production $A \rightarrow \alpha$
- α is pushed in reverse order: e.g., $S \rightarrow (S) S$ becomes $\$ S) S ($

● Match action:

- Match a token on top of the stack with the next input token

● Generations correspond to the leftmost derivation:

$S \Rightarrow (S) S$
 $\Rightarrow () S$
 $\Rightarrow ()$

$S \rightarrow (S) S$
 $S \rightarrow \epsilon$
 $S \rightarrow \epsilon$

LL(1) Parsing

- LL(1) Parsing table: $M[N,T]$ where N is a non-terminal and T is a token.

e.g.: $S \rightarrow (S)S \mid \epsilon$

| $M[N,T]$ | (|) | \$ |
|----------|----------------------|--------------------------|--------------------------|
| S | $S \rightarrow (S)S$ | $S \rightarrow \epsilon$ | $S \rightarrow \epsilon$ |

- LL(1) grammar: if the associated LL(1) parsing table has at most one production in each table entry.

LL(1) Parsing

● A non-LL(1) example:

$\langle \text{stmt} \rangle \rightarrow \langle \text{if-stmt} \rangle \mid \text{other}$

$\langle \text{if-stmt} \rangle \rightarrow \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \langle \text{else-part} \rangle$

$\langle \text{else-part} \rangle \rightarrow \text{else} \langle \text{stmt} \rangle \mid \epsilon$

$\langle \text{exp} \rangle \rightarrow 0 \mid 1$

| M[N,T] | if | other | else | 0 | 1 | \$ |
|------------------------------------|--|--|---|--|--|---|
| $\langle \text{stmt} \rangle$ | $\langle \text{stmt} \rangle \rightarrow$ $\langle \text{if-stmt} \rangle$ | $\langle \text{stmt} \rangle \rightarrow$ other | | | | |
| $\langle \text{if-stmt} \rangle$ | $\langle \text{if-stmt} \rangle \rightarrow$ if ($\langle \text{exp} \rangle$) $\langle \text{stmt} \rangle$ $\langle \text{else-part} \rangle$ | | | | | |
| $\langle \text{else-part} \rangle$ | | | $\langle \text{else-part} \rangle \rightarrow$ else $\langle \text{stmt} \rangle$ $\langle \text{else-part} \rangle \rightarrow \epsilon$ | | | $\langle \text{else-part} \rangle \rightarrow \epsilon$ |
| $\langle \text{exp} \rangle$ | | | | $\langle \text{exp} \rangle \rightarrow 0$ | $\langle \text{exp} \rangle \rightarrow 1$ | |

Left Recursion Removal

What's the problem with left recursion?

$$A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid A \alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

e.g., $\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle + \langle \text{term} \rangle \mid \langle \text{exp} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$

Removing left recursion:

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon \end{aligned}$$

e.g., $\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle \langle \text{exp}' \rangle$
 $\langle \text{exp}' \rangle \rightarrow + \langle \text{term} \rangle \langle \text{exp}' \rangle \mid - \langle \text{term} \rangle \langle \text{exp}' \rangle \mid \varepsilon$

Left Recursion Removal

How about left associativity?

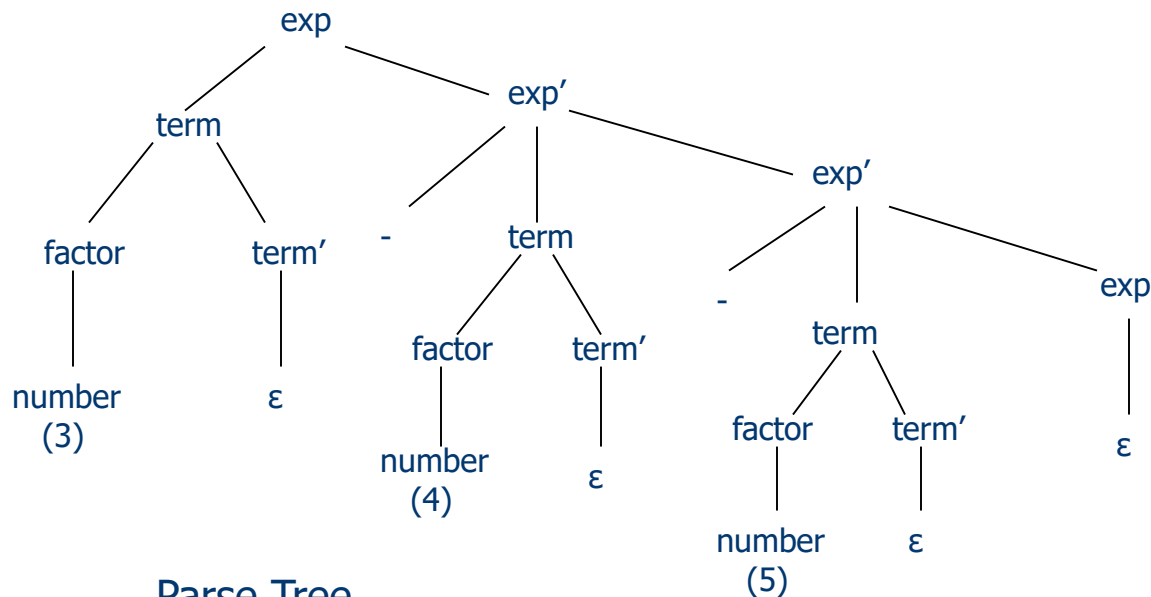
$\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle \langle \text{exp}' \rangle$

$\langle \text{exp}' \rangle \rightarrow + \langle \text{term} \rangle \langle \text{exp}' \rangle \mid - \langle \text{term} \rangle \langle \text{exp}' \rangle \mid \epsilon$

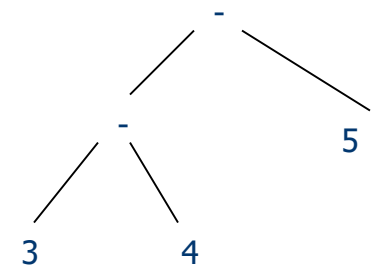
$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \langle \text{term}' \rangle$

$\langle \text{term}' \rangle \rightarrow * \langle \text{factor} \rangle \langle \text{term}' \rangle \mid / \langle \text{factor} \rangle \langle \text{term}' \rangle \mid \epsilon$

$\langle \text{factor} \rangle \rightarrow (\langle \text{exp} \rangle) \mid \text{number} \mid \text{id}$



Parse Tree



Syntax Tree

Left Factoring

- What's the problem with a common prefix of symbols?

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

e.g., $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle \mid \langle \text{stmt} \rangle$

e.g., $\langle \text{if-stmt} \rangle \rightarrow \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle$
 $\mid \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

- Left factoring the common prefix:

$$A \rightarrow \alpha A' \qquad A' \rightarrow \beta \mid \gamma$$

e.g., $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \langle \text{stmts}' \rangle$

$\langle \text{stmts}' \rangle \rightarrow ; \langle \text{stmts} \rangle \mid \epsilon$

e.g., $\langle \text{if-stmt} \rangle \rightarrow \text{if} (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \langle \text{else-part} \rangle$

$\langle \text{else-part} \rangle \rightarrow \text{else } \langle \text{stmt} \rangle \mid \epsilon$

Left Factoring

● Left factoring vs. left recursion removal:

e.g., $\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{exp} \rangle \mid \langle \text{term} \rangle$

Left factoring:

$\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle \langle \text{exp}' \rangle$

$\langle \text{exp}' \rangle \rightarrow + \langle \text{exp} \rangle \mid \varepsilon$

Substituting $\langle \text{term} \rangle \langle \text{exp}' \rangle$ for $\langle \text{exp} \rangle$:

$\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle \langle \text{exp}' \rangle$

$\langle \text{exp}' \rangle \rightarrow + \langle \text{term} \rangle \langle \text{exp}' \rangle \mid \varepsilon$

● Both left factoring and left recursion removal may obscure the associativity

FIRST and FOLLOW Sets

● Notation

- α --- string of terminals or non-terminals
- A --- single terminal or non-terminal

● $\text{FIRST}(\alpha)$ --- set of terminals that can start any string derived from α

- If $\alpha \Rightarrow^* \varepsilon$, α is nullable and $\varepsilon \in \text{FIRST}(\alpha)$

● $\text{FOLLOW}(A)$ --- set of terminals that can immediately follow A in some phrase form

Computing FIRST Sets

for all non-terminals A **do** $\text{FIRST}(A) = \{\}$

for all terminals A **do** $\text{FIRST}(A) = \{A\}$

repeat

for each production $A \rightarrow X_1 X_2 \dots X_n$ **do**

$k = 1$

while $k \leq n$ **do**

 add $\text{FIRST}(X_k) - \{\epsilon\}$ to $\text{FIRST}(A)$

if $\epsilon \notin \text{FIRST}(X_k)$ **then** break;

$k = k + 1$

if $k > n$ **then** add ϵ to $\text{FIRST}(A)$

until no changes to any $\text{FIRST}(A)$

Computing FIRST Sets

Given the following grammar:

$\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle \langle \text{addop} \rangle \langle \text{term} \rangle \mid \langle \text{term} \rangle$

$\langle \text{addop} \rangle \rightarrow + \mid -$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle \langle \text{mulop} \rangle \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{mulop} \rangle \rightarrow * \mid /$

$\langle \text{factor} \rangle \rightarrow (\langle \text{exp} \rangle) \mid \text{number}$

We get the following FIRST sets:

$\text{FIRST}(\text{exp}) = \{ (, \text{number} \}$

$\text{FIRST}(\text{term}) = \{ (, \text{number} \}$

$\text{FIRST}(\text{factor}) = \{ (, \text{number} \}$

$\text{FIRST}(\text{addop}) = \{ +, - \}$

$\text{FIRST}(\text{mulop}) = \{ *, / \}$

Computing FIRST Sets

| Production | Pass 1 | Pass 2 | Pass 3 |
|---|--|--|---|
| $\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle$ $\quad \langle \text{addop} \rangle \langle \text{term} \rangle$ | | | |
| $\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle$ | | | $\text{FIRST}(\text{exp}) = \{ (, \text{number} \}$ |
| $\langle \text{addop} \rangle \rightarrow +$ | $\text{FIRST}(\text{addop}) = \{ + \}$ | | |
| $\langle \text{addop} \rangle \rightarrow -$ | $\text{FIRST}(\text{addop}) = \{ +, - \}$ | | |
| $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle$ $\quad \langle \text{mulop} \rangle \langle \text{factor} \rangle$ | | | |
| $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle$ | | $\text{FIRST}(\text{term}) = \{ (, \text{number} \}$ | |
| $\langle \text{mulop} \rangle \rightarrow *$ | $\text{FIRST}(\text{mulop}) = \{ * \}$ | | |
| $\langle \text{mulop} \rangle \rightarrow /$ | $\text{FIRST}(\text{mulop}) = \{ *, / \}$ | | |
| $\langle \text{factor} \rangle \rightarrow (\langle \text{exp} \rangle)$ | $\text{FIRST}(\text{factor}) = \{ (\}$ | | |
| $\langle \text{factor} \rangle \rightarrow \text{number}$ | $\text{FIRST}(\text{factor}) = \{ (, \text{number} \}$ | | |

Computing FOLLOW Sets

$\text{FOLLOW}(\text{start-symbol}) = \{\$ \}$

for all non-terminals $A \neq \text{start-symbol}$ **do** $\text{FOLLOW}(A) = \{ \}$

repeat

for each production $A \rightarrow X_1 X_2 \dots X_n$ **do**

for each X_i that is a non-terminal **do**

 add $\text{FIRST}(X_{i+1} X_{i+2} \dots X_n) - \{\epsilon\}$ to $\text{FOLLOW}(X_i)$

 /* note that if $i = n$, then $X_{i+1} X_{i+2} \dots X_n = \epsilon$ */

if $\epsilon \in \text{FIRST}(X_{i+1} X_{i+2} \dots X_n)$ **then** add $\text{FOLLOW}(A)$ to $\text{FOLLOW}(X_i)$

until no changes to any $\text{FOLLOW}(A)$

Computing FOLLOW Sets

| Production | Pass 1 | Pass 2 |
|---|--|--|
| $\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle$ $\quad \langle \text{addop} \rangle \langle \text{term} \rangle$ | $\text{FOLLOW}(\text{exp}) = \{ \$, +, - \}$ $\text{FOLLOW}(\text{addop}) = \{ (, \text{number} \}$ $\text{FOLLOW}(\text{term}) = \{ \$, +, - \}$ | $\text{FOLLOW}(\text{term}) = \{ \$, +, -, *, /,) \}$ |
| $\langle \text{exp} \rangle \rightarrow \langle \text{term} \rangle$ | | |
| $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle$ $\quad \langle \text{mulop} \rangle \langle \text{factor} \rangle$ | $\text{FOLLOW}(\text{term}) = \{ \$, +, -, *, / \}$ $\text{FOLLOW}(\text{mulop}) = \{ (, \text{number} \}$ $\text{FOLLOW}(\text{factor}) = \{ \$, +, -, *, / \}$ | $\text{FOLLOW}(\text{factor}) = \{ \$, +, -, *, /,) \}$ |
| $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle$ | | |
| $\langle \text{factor} \rangle \rightarrow (\langle \text{exp} \rangle)$ | $\text{FOLLOW}(\text{exp}) = \{ \$, +, -,) \}$ | |

Another Example

| Production | Pass 1 | Pass 2 |
|--------------------------|--|--|
| $X \rightarrow a$ | $\text{FIRST}(X) = \{a\}$ | |
| $X \rightarrow Y$ | | $\text{FIRST}(X) = \{\epsilon, a, c\}$ |
| $Y \rightarrow \epsilon$ | $\text{FIRST}(Y) = \{\epsilon\}$ | |
| $Y \rightarrow c$ | $\text{FIRST}(Y) = \{\epsilon, c\}$ | |
| $Z \rightarrow d$ | $\text{FIRST}(Z) = \{d\}$ | |
| $Z \rightarrow XYZ$ | $\text{FIRST}(Z) = \{a, d\}$ $\text{FOLLOW}(X) = \{a, c, d\}$ $\text{FOLLOW}(Y) = \{a, c, d\}$ | $\text{FIRST}(Z) = \{a, c, d\}$ |

LL(1) Parsing Table

Construction Rule:

- for each a in $\text{FIRST}(\alpha)$, enter $A \rightarrow \alpha$ in $M[A, a]$.
- if ϵ is in $\text{FIRST}(\alpha)$, enter $A \rightarrow \alpha$ in $M[A, f]$ for each f in $\text{FOLLOW}(A)$.

| | a | c | d |
|---|--|---|--|
| X | X \rightarrow a X \rightarrow Y | X \rightarrow Y | X \rightarrow Y |
| Y | Y \rightarrow ϵ | Y \rightarrow ϵ Y \rightarrow c | Y \rightarrow ϵ |
| Z | Z \rightarrow XYZ | Z \rightarrow XYZ | Z \rightarrow d Z \rightarrow XYZ |

Error Recovery

- Report the first error and tell where the error has occurred
- Report error(s) and allow parsing to continue
 - deletion: skip tokens until reaching one in FIRST or FOLLOW sets
 - worst case: consume all input tokens
- Error repair
 - insert/replace/delete/switch tokens based on minimal distance
 - dangerous (possible infinite loop if repair generates another error)
 - expensive (usually limited to simple cases)

Error Recovery

● Panic mode for recursive descent parsers:

- Skip input tokens in order to resume parsing
- may consume all input tokens, but always terminate

```
void scanto( synchset ) {  
    while( token  $\notin$  synchset  $\cup$  { $\$$ } )  
        advance();  
}
```

```
void checkinput( firstset, followset ) {  
    if( token  $\notin$  firstset ) {  
        error();  
        scanto( firstset  $\cup$  followset );  
    }  
}
```

Error Recovery

● Panic mode (continued):

```
void exp( synchset ) {
    checkinput( { (, number }, synchset );
    if( token ∉ synchset ) {
        term( synchset );
        while( token == + || token == - ) {
            match( token );
            term( synchset );
        }
        checkinput( synchset, { (, number } );
    }
}
```

```
void factor( synchset ) {
    checkinput( { (, number }, synchset );
    if( token ∉ synchset ) {
        switch( token ) {
            case (:
                match( ( );
                exp( { ) } );
                match( ) );
                break;
            case number:
                match( number );
                break;
            default: error();
        }
        checkinput( synchset, { (, number } );
    }
}
```

Error Recovery in LL(1)

| M[N,T] | (| number |) | + | - | * | / | \$ |
|--------|----------------------------------|----------------------------------|-------------------------|--|--|---|---|-------------------------|
| exp | <exp> -> <term> <exp'> | <exp> -> <term> <exp'> | pop | scan | scan | scan | scan | pop |
| exp' | scan | scan | <exp'> -> ϵ | <exp'> -> <addop> <term> <exp'> | <exp'> -> <addop> <term> <exp'> | scan | scan | <exp'>-> ϵ |
| addop | pop | pop | scan | <addop> -> + | <addop>-> - | scan | scan | pop |
| term | <term> -> <factor> <term'> | <term> -> <factor> <term'> | pop | pop | pop | scan | scan | pop |
| term' | scan | scan | <term'>-> ϵ | <term'> -> ϵ | <term'> -> ϵ | <term'>-> <mulop> <factor> <term'> | <term'>-> <mulop> <factor> <term'> | <term'>-> ϵ |
| mulop | pop | pop | scan | scan | scan | <mulop>-> * | <mulop>-> / | pop |
| factor | <factor> -> (<exp>) | <factor> -> number | pop | pop | pop | pop | pop | pop |

Error Recovery in LL(1)

- Enter “pop”:
if token == \$ or token in FOLLOW(A)
- Enter “scan”:
if token != \$ and
token not in FIRST(A) or FOLLOW(A)
- Enter normal entries:
if token != \$ and token in FIRST(A)