

Parsing

Stefan D. Bruda

CS 310, Winter 2025



PARSING

- Interface to lexical analysis:

```
typename vocab;      /* alphabet + end-of-string */  
const vocab EOS;    /* end-of-string pseudo-token */  
vocab gettoken(void); /* returns next token */
```

- Parsing = determining whether the current input belongs to the given language

- In practice a parse tree is constructed in the process as well

- General method: Not as efficient as for finite automata

- Various paths with different states and also **different stack content** along them
 - Careful housekeeping (**dynamic programming**) reduces the otherwise exponential complexity to $O(n^3)$
 - We want linear time instead = **deterministic PDA**
 - General idea: **try to determine what to do next based on the next token in the input**



RECURSIVE DESCENT PARSING

- Basic idea: simulate a PDA for the given language
 - Stack manipulation = ...



RECURSIVE DESCENT PARSING

- Basic idea: simulate a PDA for the given language
 - Stack manipulation = function calls and returns
- Construct a function for each nonterminal
- Decide which function to call based on the next input token = linear complexity

```
vocab t;
void MustBe (vocab ThisToken) {
    if (t != ThisToken) { printf("reject"); exit(0); }
    t = gettoken();
}

void Balanced (void) {
    switch (t) {
        case EOS:
        case ONE: /* <empty> */
            break;
        default: /* 0 <balanced> 1 */
            MustBe(ZERO);
            Balanced();
            MustBe(ONE);
    } /* switch */
} /* Balanced */

int main (void) {
    t = gettoken();
    Balanced();
    /* accept iff
       t == EOS */
}
```

RECURSIVE DESCENT PARSING: LEFT FACTORING



- Not all grammars are suitable for recursive descent:

$$\begin{aligned}\langle \text{stmt} \rangle &::= \langle \text{empty} \rangle \\&\quad | \quad \text{VAR} := \langle \text{exp} \rangle \\&\quad | \quad \text{IF } \langle \text{exp} \rangle \text{ THEN } \langle \text{stmt} \rangle \text{ ELSE } \langle \text{stmt} \rangle \\&\quad | \quad \text{WHILE } \langle \text{exp} \rangle \text{ DO } \langle \text{stmt} \rangle \\&\quad | \quad \text{BEGIN } \langle \text{seq} \rangle \text{ END} \\ \langle \text{seq} \rangle &::= \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{seq} \rangle\end{aligned}$$

- Both rules for $\langle \text{seq} \rangle$ begin with the same nonterminal
- Impossible to decide which one to apply based only on the next token
- Fortunately concatenation is distributive over union so we can fix the grammar (**left factoring**):

$$\begin{aligned}\langle \text{seq} \rangle &::= \langle \text{stmt} \rangle \langle \text{seqTail} \rangle \\ \langle \text{seqTail} \rangle &::= \langle \text{empty} \rangle \mid ; \langle \text{seq} \rangle\end{aligned}$$

RECURSIVE DESCENT PARSING: AMBIGUITY



- Some programming constructs are inherently ambiguous

```
<stmt> ::= if ( <exp> ) <stmt>
          | if ( <exp> ) <stmt> else <stmt>
```

- Solution: choose one path and stick to it (e.g., match the else-statement with the nearest else-less if statement)

```
case IF:
    t = gettoken();
    MustBe(OPEN_PAREN);
    Expression();
    MustBe(CLS_PAREN);
    Statement();
    if (t == ELSE) {
        t = gettoken();
        Statement();
    }
```



- Consider the following code:

```
int y;
template <class T> void g(T& v) {
    T::x(y);
}
```

- The statement `T::x(y)` can be

AMBIGUITY (C++): IT'S A TYPE, IT'S A PLANE



- Consider the following code:

```
int y;
template <class T> void g(T& v) {
    T::x(y);
}
```

- The statement `T::x(y)` can be
 - A function call (member function `x` of `T` applied to `y`), or
 - A declaration of `y` as a variable of type `T::x`

AMBIGUITY (C++): IT'S A TYPE, IT'S A PLANE



- Consider the following code:

```
int y;
template <class T> void g(T& v) {
    T::x(y);
}
```

- The statement `T::x(y)` can be
 - A function call (member function `x` of `T` applied to `y`), or
 - A declaration of `y` as a variable of type `T::x`
- Resolution: unless otherwise specified, an identifier is assumed to refer to something that is **not** a type or template.
 - If we want something else, we use the keyword `typename`:

```
T::x(y);           // function x of T applied to y
typename T::x(y); // y is a variable of type T::x
```
 - The sole purpose of `typename` is to resolve the ambiguity in the desired direction



RECURSIVE DESCENT PARSING: CLOSURE, ETC.

- Any left recursion in the grammar will cause the parser to go into an infinite loop:

$$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle \langle \text{addop} \rangle \langle \text{term} \rangle \mid \langle \text{term} \rangle$$

- Solution: eliminate left recursion using a closure

$$\begin{aligned}\langle \text{exp} \rangle &::= \langle \text{term} \rangle \langle \text{closure} \rangle \\ \langle \text{closure} \rangle &::= \langle \text{empty} \rangle \\ &\quad \mid \langle \text{addop} \rangle \langle \text{term} \rangle \langle \text{closure} \rangle\end{aligned}$$

- Not the same language theoretically, but differences not relevant in practice
- This being said, some languages are simply not parseable using recursive descent

$$\langle \text{palindrome} \rangle ::= \langle \text{empty} \rangle \mid 0 \mid 1 \mid 0 \langle \text{palindrome} \rangle 0 \mid 1 \langle \text{palindrome} \rangle 1$$

- No way to know when to choose the $\langle \text{empty} \rangle$ rule
- No way to choose between the second and the fourth rule
- No way to choose between the third and the fifth rule



RECURSIVE DESCENT PARSING: SUFFICIENT CONDITIONS

- $\text{first}(\alpha)$ = set of all initial tokens in the strings derivable from α
- $\text{follow}(\langle N \rangle)$ = set of all initial tokens in nonempty strings that may follow $\langle N \rangle$ (possibly including EOS)
- Sufficient conditions for a grammar to allow recursive descent parsing:
 - For $\langle N \rangle ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ must have $\text{first}(\alpha_i) \cap \text{first}(\alpha_j) = \emptyset$,
 $1 \leq i < j \leq n$
 - Whenever $\langle N \rangle \Rightarrow^* \varepsilon$ must have $\text{follow}(\langle N \rangle) \cap \text{first}(\langle N \rangle) = \emptyset$
- Grammars that do not have these properties may be fixable using left factoring, closure, etc.

RECURSIVE DESCENT PARSING: SUFFICIENT CONDITIONS



- $\text{first}(\alpha)$ = set of all initial tokens in the strings derivable from α
- $\text{follow}(\langle N \rangle)$ = set of all initial tokens in nonempty strings that may follow $\langle N \rangle$ (possibly including EOS)
- Sufficient conditions for a grammar to allow recursive descent parsing:
 - For $\langle N \rangle ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ must have $\text{first}(\alpha_i) \cap \text{first}(\alpha_j) = \emptyset$, $1 \leq i < j \leq n$
 - Whenever $\langle N \rangle \Rightarrow^* \varepsilon$ must have $\text{follow}(\langle N \rangle) \cap \text{first}(\langle N \rangle) = \emptyset$
- Grammars that do not have these properties may be fixable using left factoring, closure, etc.
- Method for constructing the recursive descent function $N()$ for the nonterminal $\langle N \rangle$ with rules $\langle N \rangle ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$:
 - ① For $\alpha_i \neq \varepsilon$ apply the rewriting rule $\langle N \rangle ::= \alpha_i$ whenever the next token in the input is in $\text{first}(\alpha_i)$
 - ② For $\alpha_i = \varepsilon$ apply the rewriting rule $\langle N \rangle ::= \alpha_i$ (that is, $\langle N \rangle ::= \varepsilon$) whenever the next token in the input is in $\text{follow}(\langle N \rangle)$
 - ③ Signal a syntax error in all the other cases



RECURSIVE DESCENT EXAMPLE

```
typedef enum { VAR, EQ, IF, ELSE, WHILE, OPN_BRACE, CLS_BRACE,
               OPN_PAREN, CLS_PAREN, SEMICOLON, EOS } vocab;
vocab gettoken() {...}
vocab t;
void MustBe(vocab ThisToken) {...}

void Statement();
void Sequence();

int main() {
    t = gettoken();
    Statement();
    if (t != EOS) printf("String not accepted\n");
    return 0;
}
void Sequence() {
    if (t == CLS_BRACE) /* <empty> */;
    else { /* <statement> <sequence> */
        Statement();
        Sequence();
    }
}
```



RECURSIVE DESCENT EXAMPLE (CONT'D)

```
void Statement() {
    switch(t) {
        case SEMICOLON: /* ; */
            t = gettoken();
            break;
        case VAR: /* VAR = <exp> */
            t = gettoken();
            MustBe(EQ);
            Expression();
            MustBe(SEMICOLON);
            break;
        case IF: /* if (<expr>) <statement> else <statement> */
            t = gettoken();
            MustBe(OPEN_PAREN);
            Expression();
            MustBe(CLS_PAREN);
            Statement();
            MustBe(ELSE);
            Statement();
            break;
    }
}
```



RECURSIVE DESCENT EXAMPLE (CONT'D)

```
case WHILE: /* while (exp) <statement> */
    t = gettoken();
    MustBe(OPEN_PAREN);
    Expression();
    MustBe(CLS_PAREN);
    Statement();
    break;
default: /* { <sequence > } */
    MustBe(OPN_BRACE);
    Sequence();
    MustBe(CLS_BRACE);
} /* switch */
} /* Statement () */
```



CONSTRUCTING THE PARSE TREE

- The parse tree can be constructed through the recursive calls:
 - Each function creates a current node
 - The children are populated through recursive calls
 - The current node is then returned

```
class Node {...};

Node* Sequence() {
    Node* current = new Node(SEQ, ...);
    if (t == CLS_BRACE) /* <empty> */ ;
    else { /* <statement> <sequence> */
        current.addChild(Statement());
        current.addChild(Sequence());
    }
    return current;
}
```



CONSTRUCTING THE PARSE TREE (CONT'D)

```
Node* Statement() {
    Node* current;
    switch(t) {
        case SEMICOLON: /* ; */
            t = gettoken();
            return new Node(EMPTY);
            break;
        case VAR: /* VAR = <exp> */
            current = new Node(ASSIGN, ...);
            current.addChild(VAR, ...);
            t = gettoken();
            MustBe(EQ);
            current.addChild(Expression());
            MustBe(SEMICOLON);
            break;
        case IF: /* if (<expr>) <statement> else <statement> */
            current = new Node(COND, ...);
            /* ... */
    }
    return current;
}
```



PARSING IN PRACTICE: BOTTOM-UP PARSERS

- Substantial effort may be needed to make grammars suitable for recursive descent
- Recursive descent code relatively difficult to generate
- Recursive descent code uses processor stack frames (expensive space-wise), it is more efficient to simulate a pushdown automaton
- Bottom-up parsing** better in practice: Given $G = (N, \Sigma, S, R)$ construct the automaton $M = (\{p, q\}, \Sigma, N + \Sigma, \Delta, s, \{q\})$ with Δ containing exactly all the transitions:

$$\begin{array}{ll} \text{shift} & \forall a \in \Sigma : ((p, a, \varepsilon), (p, a)) \\ \text{reduce} & \forall A \rightarrow \alpha \in R : ((p, \varepsilon, \alpha^R), (p, A)) \\ \text{done} & ((p, \varepsilon, S), (q, \varepsilon)) \end{array}$$

- Result highly nondeterministic, but nondeterminism fixable for **weak-precedence grammars** = most programming languages:
 - When to shift and when to reduce?
 - Establish a **precedence relation** $P \subseteq (N \cup \Sigma) \times \Sigma$
 - If $(\text{stack-top}, \text{input}) \in P$ then we reduce, else we shift
 - When we reduce, with what rule we reduce?
 - We use the **longest rule** = **greedy** (eat up the longest stack top)



EXAMPLE OF BOTTOM-UP PARSING

$E \rightarrow E + T$	P	()	y	+	*	EOS		
$E \rightarrow T$	(Input	Stack
$T \rightarrow T * F$)		✓		✓	✓	✓	$y + y^*y$	ϵ
$T \rightarrow F$	y		✓		✓	✓	✓	shift	$+ y^*y$
$F \rightarrow (E)$	+							red	$+ y^*y$
$F \rightarrow y$	*							red	F
	E							red	T
	T		✓		✓		✓	red	E
	F		✓		✓	✓	✓	shift	$+ y^*y$
$((p, a, \epsilon), (p, a)), a \in \{+, *, ., (,), y\}$								y^*y	$+ E$
$((p, \epsilon, T + E), (p, E))$								shift	$y + E$
$((p, \epsilon, T), (p, E))$								red	$F + E$
$((p, \epsilon, F * T), (p, T))$								red	$T + E$
$((p, \epsilon, F), (p, T))$								shift	$* T + E$
$((p, \epsilon,)E(), (p, F))$								shift	$y * T + E$
$((p, \epsilon, y), (p, F))$								red	$F * T + E$
$((p, \epsilon, E) (q, \epsilon))$								g-red	$T + E$
$((p, \epsilon, E) (q, \epsilon))$								g-red	E
$((p, \epsilon, E) (q, \epsilon))$								done	ϵ
									ϵ

red = reduce (unambiguous)

g-red = greedy reduce (longest rule)