

CS 536 Announcements for Thursday, February 20, 2025

Last Time

- syntax-directed translation
- abstract syntax trees
- implementing ASTs

Today

- Java CUP
- approaches to parsing
- bottom-up parsing
- CFG transformations
 - removing useless non-terminals
 - Chomsky normal form (CNF)
- CYK algorithm

Next Time

- Midterm 1, 6:30 – 8 pm (in-class)

Parser generators

Tools that take an SDT spec and build an AST

- **YACC**
- **Java CUP**

Conceptually similar to JLex:

- Input: language rules + actions
- Output: Java code

parser specification → **Java CUP** → **parser source symbols**

Java CUP

parser.java

- constructor takes argument of type `Yylex`
- `parse` method
 - if input correct, returns `Symbol` whose `value` field contains translation of root nonterm
 - if input incorrect, quits on first syntax error
- uses output of `JLex`
 - depends on scanner and `TokenVal` classes
 - `sym.java` defines the communication language
- uses definitions of AST classes (in `ast.java`)

Parts of Java CUP specification

Grammar rules with actions:

```
expr ::= INTLIT
      | ID
      | expr PLUS expr
      | expr TIMES expr
      | LPAREN expr RPAREN
      ;
```

Terminal and nonterminal declarations:

```
terminal INTLIT;
terminal ID;
terminal PLUS;
terminal TIMES;
terminal LPAREN;
terminal RPAREN;
```

```
non terminal expr;
```

Precedence and associativity declarations:

```
precedence left PLUS;
precedence left TIMES;
```

Java CUP Example

Assume:

- Java class `ExpNode` with subclasses `IntLitNode`, `IdNode`, `PlusNode`, `TimesNode`
- `PlusNode` and `TimesNode` each have two children
- `IdNode` has a `String` field (for the identifier)
- `IntLitNode` has an `int` field (for the integer value)
- `INTLIT` token is represented by `IntLitTokenVal` class and has field `intVal`
- `ID` token is represented by `IdTokenVal` class and has field `idVal`

Step 1: add types to terminals and nonterminals

```
/*
 * Terminal declarations
 */
terminal INTLIT;
terminal ID;
terminal PLUS;
terminal TIMES;
terminal LPAREN;
terminal RPAREN;

/*
 * Nonterminal declarations
 */
non terminal expr;
```

Step 2: add precedences and associativities

```
/*
 * Precedence and associativity declarations
 */
precedence left PLUS;
precedence left TIMES;
```

Java CUP Example (cont.)

Step 3: add actions to CFG rules

```
/*
 * Grammar rules with actions
 */
expr ::= INTLIT
      { :
      : }
      | ID
      { :
      : }
      | expr PLUS expr
      { :
      : }
      | expr TIMES expr
      { :
      : }
      | LPAREN expr RPAREN
      { :
      : }
      ;
```

Java CUP Example (cont.)

Input: 2 + 3

Translating lists

Example

$\text{idList} \rightarrow \text{idList COMMA ID} \mid \text{ID}$

Left-recursion or right-recursion?

- for top-down parsers
- for Java CUP

Example

CFG: $\text{idList} \rightarrow \text{idList COMMA ID} \mid \text{ID}$

Goal: the translation of an `idList` is a `LinkedList` of `Strings`

Example

Input: `x , y , z`

Output:

Example (cont.)

Java CUP specification for this syntax-directed translation

Terminal and nonterminal declarations:

Grammar rules and actions:

```
idList ::= idList      COMMA      ID
        {:
```

```
        :}
| ID
  {:
```

```
        :}
;
```

Handling unary minus

```
/*
 * precedences and associativities of operators
 */
precedence left PLUS, MINUS;
precedence left TIMES, DIVIDE;
```

```
/*
 * grammar rules
 */
exp ::= . . .
    | MINUS exp:e
      { : RESULT = new UnaryMinusNode(e);
        : }
    | exp:e1 PLUS exp:e2
      { : RESULT = new PlusNode(e1, e2);
        : }
    | exp:e1 MINUS exp:e2
      { : RESULT = new MinusNode(e1, e2);
        : }
    . . .
;
```


Parsing: two approaches

Top-down / "goal driven"

- start at start nonterminal
- grow parse tree downward until entire sequence is matched

Bottom-up / "data driven"

- start with terminals (sequence)
- generate ever larger subtrees until get to single tree whose root is the start nonterminal

Example:

CFG: $\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$

$\text{term} \rightarrow \text{term} * \text{ID} \mid \text{ID}$

Derive: ID + ID

Cocke – Younger – Kasami (CYK) algorithm

- Works bottom-up
- Time complexity : $O(n^3)$
- Requires grammar to be in Chomsky Normal Form

Chomsky Normal Form (CNF)

- all rules must be in one of two forms
 - $x \rightarrow T$
 - $x \rightarrow a b$
- only rule allowed to derive epsilon is the start symbol s

Why CNF is helpful?

- nonterminals in pairs
- nonterminals (except start) can't derive epsilon

CYK : Dynamic Programming

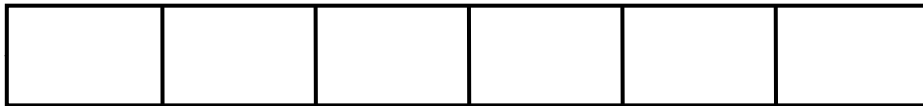
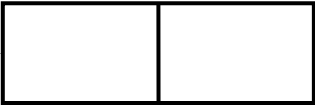
$x \rightarrow T$

$x \rightarrow a b$

Running CYK

Track every viable subtree from leaf to root.

All subspans for a sequence (string) with 6 terminals



CYK Example



$f \rightarrow iw$

$f \rightarrow iy$

$w \rightarrow lx$

$x \rightarrow nr$

$y \rightarrow lr$

$n \rightarrow \mathbf{ID}$

$n \rightarrow iz$

$z \rightarrow cn$

$i \rightarrow \mathbf{ID}$

$l \rightarrow ($

$r \rightarrow)$

$c \rightarrow ,$

Eliminating useless nonterminals

Avoid unnecessary work – remove *useless* rules

1. If a nonterminal cannot derive a sequence of terminal symbols, then it is *useless*
2. If a nonterminal cannot be derived from the start symbol, then it is *useless*

Nonterminals that cannot derive a sequence of terminal symbols

mark all terminal symbols

repeat

 if all symbols on the RHS of a production are marked

 mark the LHS nonterminal

until no more nonterminals can be marked

Example

$s \rightarrow x \mid y$

$x \rightarrow ()$

$y \rightarrow (yy)$

Nonterminals that cannot be derived from the start symbol

mark the start symbol

repeat

 if the LHS of a production is marked

 mark all RHS nonterminals

until no more nonterminals can be marked

Example

$s \rightarrow ab$

$a \rightarrow + \mid - \mid \epsilon$

$b \rightarrow \text{digit} \mid b \text{ digit}$

$c \rightarrow .b$

Chomsky Normal Form

Four steps

- eliminate epsilon productions
- eliminate unit productions
- fix productions with terminal on RHS (along with other symbols)
- fix productions with > 2 nonterminals on RHS

Eliminate (most) epsilon productions

If nonterminal a immediately derives epsilon

- make copies of all rules with a on RHS
- delete all combinations of a in the copies

Example 1

$$f \rightarrow \mathbf{ID} (a)$$

$$a \rightarrow \varepsilon$$

$$a \rightarrow n$$

$$n \rightarrow \mathbf{ID}$$

$$n \rightarrow \mathbf{ID}, n$$

Example 2

$$x \rightarrow a \mathbf{X} a \mathbf{Y} a$$

$$a \rightarrow \varepsilon$$

$$a \rightarrow \mathbf{Z}$$

Chomsky Normal Form (cont.)

Eliminate unit productions

Productions of the form $a \rightarrow b$ are called *unit productions*

If this is the only rule with a on the LHS

- place b anywhere a could have appeared
- remove the unit production $a \rightarrow b$

Example

$f \rightarrow \text{ID} (a)$

$f \rightarrow \text{ID} ()$

$a \rightarrow n$

$n \rightarrow \text{ID}$

$n \rightarrow \text{ID} , n$

If there are multiple rules with a on the LHS,

- for each rule of the form $b \rightarrow \delta$, add $a \rightarrow \delta$
- remove $a \rightarrow b$

Example

$a \rightarrow b \mathbf{X}$

| $c b$

| b

$b \rightarrow \mathbf{Z} \mathbf{Y}$

| $\mathbf{Y} c$

$c \rightarrow \mathbf{Z} a$

Chomsky Normal Form (cont.)

Fix RHS nonterminals

For productions with terminals and something else on the RHS

- for terminal **T**, add rule $x \rightarrow T$
- replace **T** with x in those productions

Example

$f \rightarrow \mathbf{ID} (n)$

$f \rightarrow \mathbf{ID} ()$

$n \rightarrow \mathbf{ID}$

$n \rightarrow \mathbf{ID} , n$

For productions with > 2 nonterminals on the RHS

- replace all but the 1st nonterminal with a new nonterminal
- add rule with new nonterminal on LHS and replaced nonterminal sequence on RHS
- repeat (as necessary)

Example