

CS 536 Announcements for Thursday, February 20, 2025

Last Time

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

← end of Midterm 1 material

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)

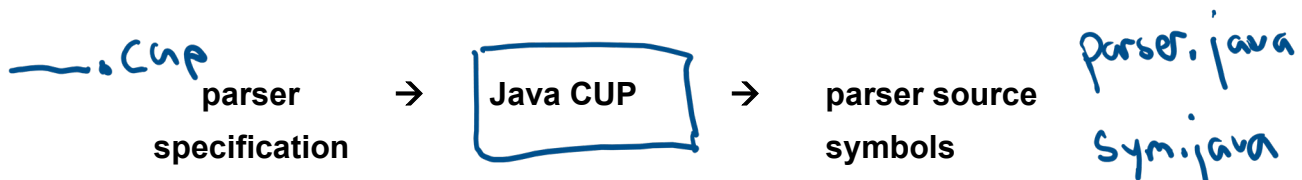
Parser generators

Tools that take an SDT spec and build an AST

- YACC *Yet Another C Compiler* - creates a parser in C
- Java CUP *Constructor of Useful Parsers* - creates a parser in Java

Conceptually similar to JLex:

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



Java CUP

parser.java

- **constructor** takes argument of type `Yylex` ← from scanner
- **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** ← back.jlex
 - depends on scanner and `TokenVal` classes
 - `sym.java` defines the communication language = defines token names used by both JLex & Java CUP
- 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
      ;
  
```

not shown (yet)

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;
  
```

↑ associativity

can do:
precedence nonassoc LESS;

order (in `_.cup`)
indicates precedence
low
↓
high

Java CUP Example

defined in ast.java

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

defined in bach.jlex

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

Need type if we want to use value associated with token

terminal `IntLitTokenVal INTLIT;`
terminal `IdTokenVal ID;`

↳ from scanner (bach.jlex)

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

Type required for all non-terms

↳ non terminal `ExpNode expr;`

↳ from ast.java

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 : i ← type is IntLitTokenVal
      {
        RESULT = new IntLitNode(i.intVal);
      }
      | ID : i
      {
        RESULT = new IdNode(i.idVal);
      }
      | expr : e1 PLUS expr : e2
      {
        RESULT = new PlusNode(e1, e2);
      }
      | expr : e1 TIMES expr : e2
      {
        RESULT = new TimesNode(e1, e2);
      }
      | LPAREN expr : e RPAREN
      {
        RESULT = e;
      }
      ;

```

Subclasses
of
ExpNode

General format:

```

nonterm ::= rule1
        {
          // action for rule1
          RESULT = ... ;
        }
        | rule2
        {
          RESULT = ... ;
        }
        | ...
        ;

```

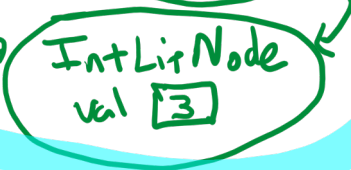
Java CUP Example (cont.)

Input: 2 + 3



Parse tree w/translation

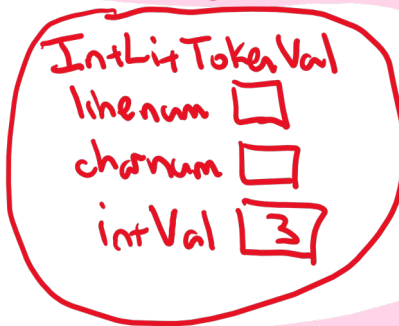
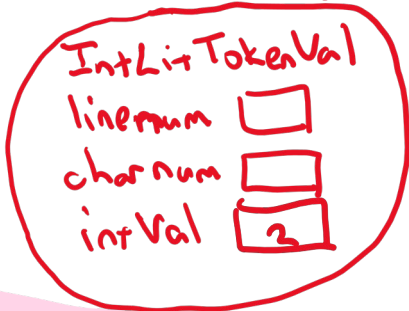
Built by parser



INTLIT

INTLIT

Built by scanner



Translating lists

Example

idList → idList COMMA ID | ID

← left recursive

Left-recursion or right-recursion?

- for **top-down** parsers must use right recursion
left-recursion leads to an infinite loop
- for Java CUP use left recursion
↳ **bottom-up** parser

Example

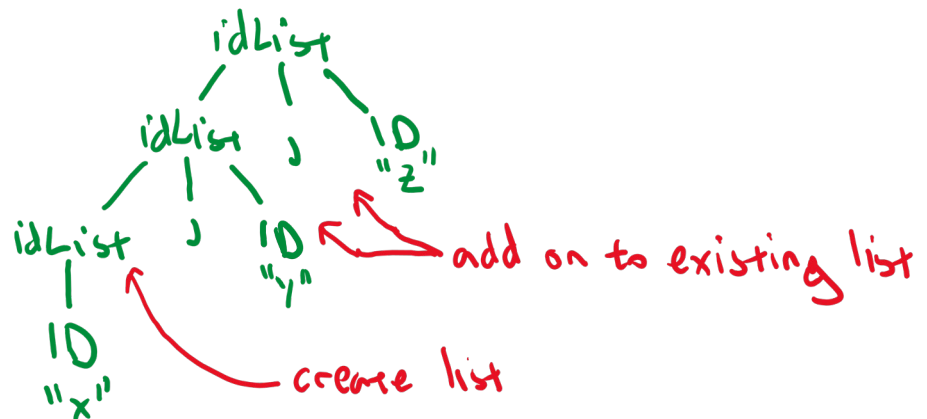
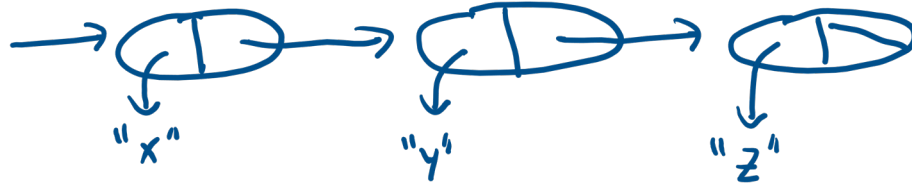
CFG: idList → idList COMMA ID | ID

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

Example

Input: x , y , z

Output:



Example (cont.)

idList → *idList* COMMA *ID*
| *ID*

Java CUP specification for this syntax-directed translation

Terminal and nonterminal declarations:

terminal *IdToken* Val *ID*;
terminal COMMA;
non terminal *LinkedList*<*String*> *idList*;

Grammar rules and actions:

```
idList ::= idList : L COMMA ID : i
        { :
          L.addLast(i.idVal);
          RESULT = L;
        : }
| ID : i
  { :
    LinkedList<String> L = new LinkedList<String>();
    L.add(i.idVal);
    RESULT = L;
  : }
;
```

Handling unary minus

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

binary minus has lowest precedence
"phony" token (never returned by scanner)
unary minus has highest precedence

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

Precedence of a rule is that of the last token in the rule, unless assigned a specific precedence using %prec

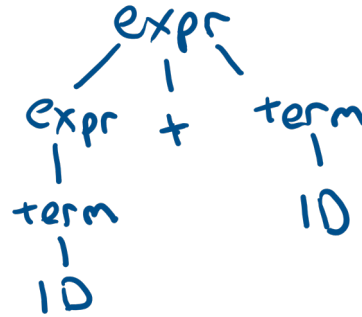
Resume at 7:40 pm

Full-sized versions of diagrams for running the CYK algorithm (pages 11-12) available

Parsing: two approaches

Top-down / "goal driven"

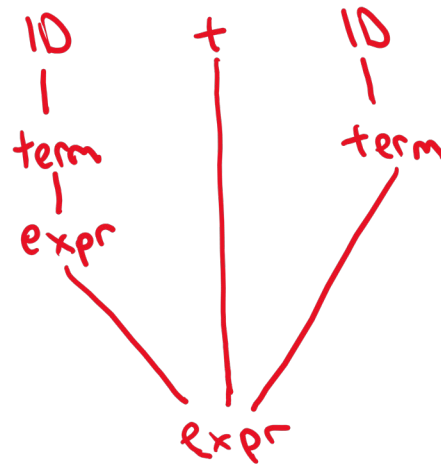
- start at start nonterminal
- grow parse tree downward until entire sequence is matched



Bottom-up / "data driven"

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

(note: parse tree is upside down)



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

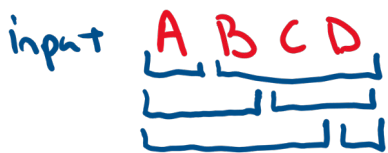
- Works bottom-up
- Time complexity : $O(n^3)$ $n = \text{length of input (\# of tokens in sequence)}$
- Requires grammar to be in Chomsky Normal Form

Chomsky Normal Form (CNF)

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

Why CNF is helpful?

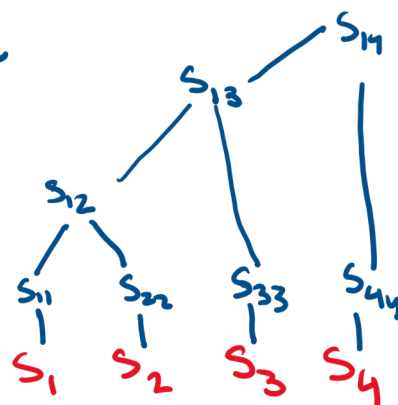
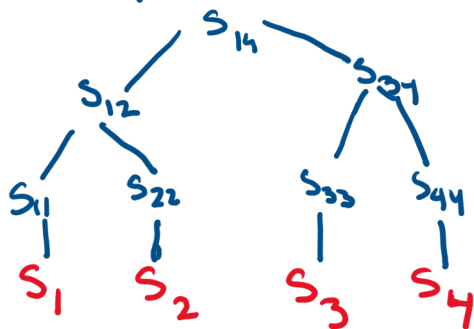
- nonterminals in pairs \rightarrow can think of a subtree as a subspan
- nonterminals (except start) can't derive epsilon \rightarrow each subspan has at least 1 token



CYK : Dynamic Programming

$x \rightarrow T \rightarrow$ forms leaf of parse tree

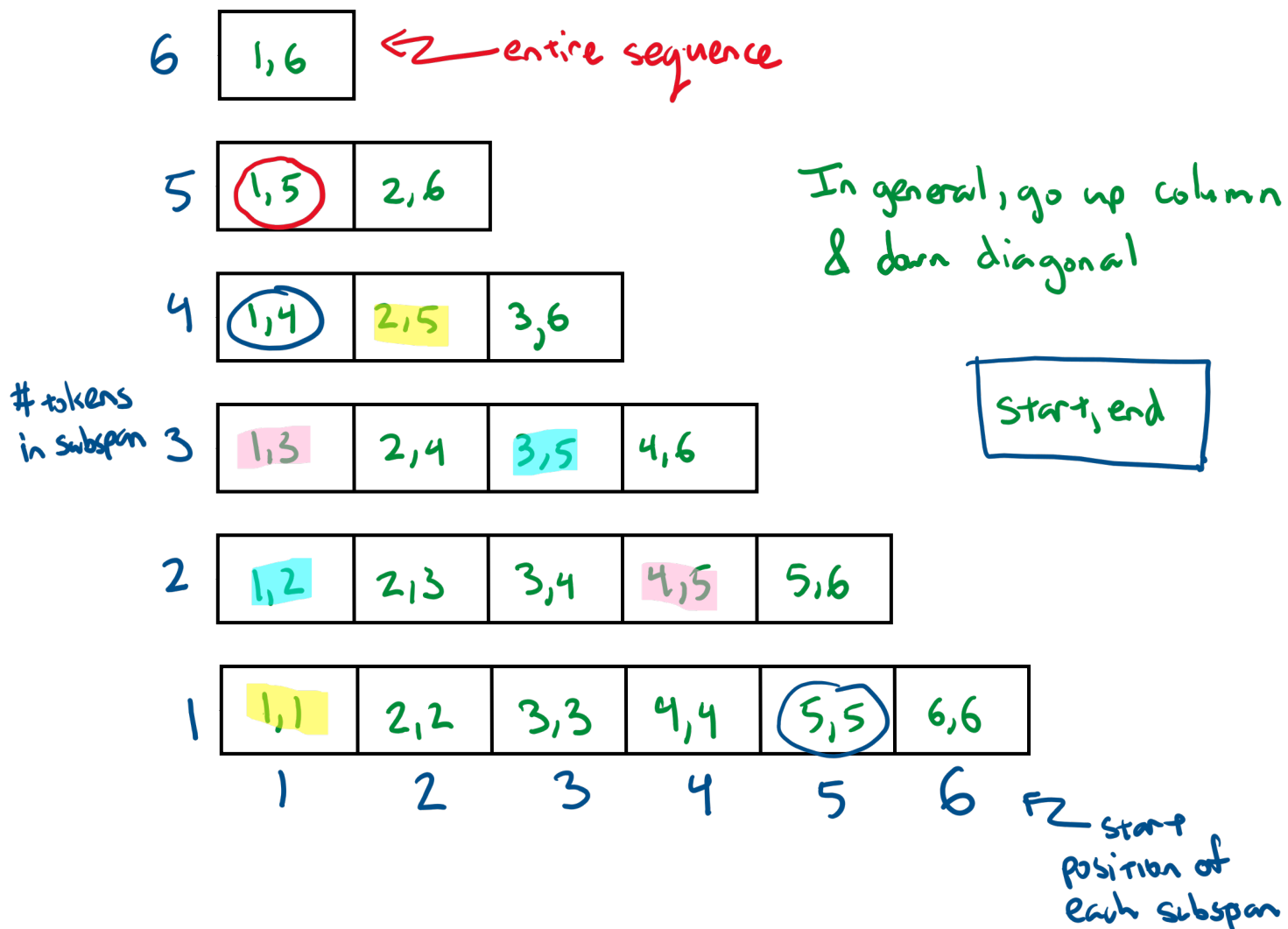
$x \rightarrow a b \rightarrow$ binary interior node of parse tree



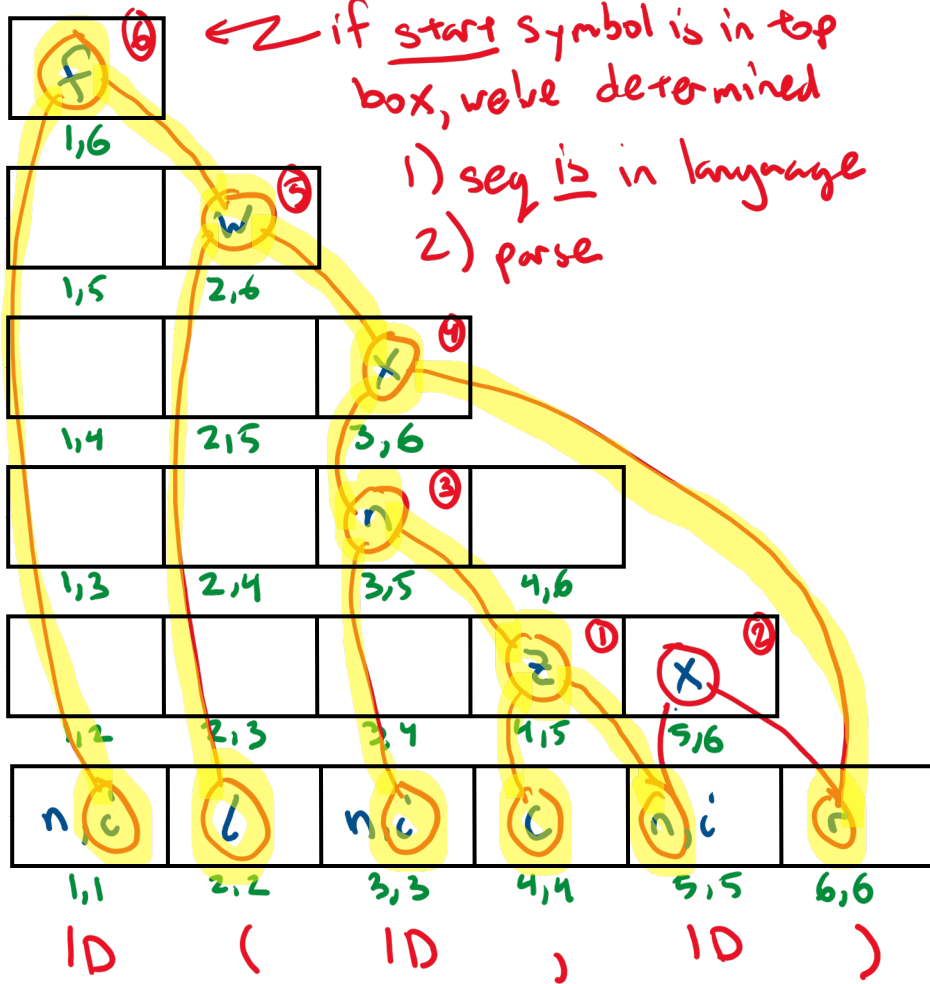
Running CYK

Track every viable subtree from leaf to root.

All subspans for a sequence (string) with 6 terminals



CYK Example



- f → iw ⁶
- f → iy ⁵
- w → lx ⁵
- x → nr ^{2 4}
- y → lr
- n → ID
- n → iz ³
- z → cn ¹
- i → ID
- l → (
- r →)
- c → ,

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** (everywhere it shows up)

until no more nonterminals can be marked

Example

$s \rightarrow x$
 $x \rightarrow ()$
 ~~$y \rightarrow (yy)$~~

~~y~~

$s \rightarrow x$
 $| y$
ie $s \rightarrow x$
 ~~$s \rightarrow y$~~

y is useless

$s \rightarrow x$
 $x \rightarrow ()$

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** (whenever they show up)

until no more nonterminals can be marked

Example

$s \rightarrow ab$
 $a \rightarrow + | - | \epsilon$
 $b \rightarrow \text{digit} | b \text{ digit}$

~~$c \rightarrow .b$~~ $\leftarrow c$ is useless

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

OK to have start $S \rightarrow \epsilon$
but if so, can't have start on RHS

if we have
 $S \rightarrow \epsilon$
 $a \rightarrow \dots S \dots$
 ↓ translate to
 make new start symbol
 S_0 & add production
 $S_0 \rightarrow S$
 $S \rightarrow \epsilon$
 $a \rightarrow \dots S \dots$

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

$f \rightarrow ID(a)$
 $a \rightarrow \epsilon$
 $a \rightarrow n$
 $n \rightarrow ID$
 $n \rightarrow ID, n$

$f \rightarrow ID(a)$
 $f \rightarrow ID(\cancel{a}) \Rightarrow f \rightarrow ID()$
 $a \rightarrow n$
 $n \rightarrow ID$
 $n \rightarrow ID, n$

Example 2

$x \rightarrow aXaYa$
 $a \rightarrow \epsilon$
 $a \rightarrow Z$

$x \rightarrow aXaYa$
 $| aXaY$
 $| aXYa$
 $| XaYa$
 $| aXY$
 $| XaY$
 $| XYa$
 $| XY$
 $a \rightarrow Z$

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

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

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

$a \rightarrow n$

$n \rightarrow \text{ID}$

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

$f \rightarrow \text{ID}(n)$
 $\quad | \text{ID}()$
 $n \rightarrow \text{ID}$
 $\quad | \text{ID}, n$

If there are multiple rules with a on the LHS,

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

← note: on-line reading is incorrect on how to handle this

Example

$a \rightarrow bX$

$\quad | cb$

$\quad | \underline{b}$

$\underline{b} \rightarrow \underline{ZY}$

$\quad | \underline{Yc}$

$c \rightarrow Za$

$a \rightarrow bX$
 $\quad | cb$
 $\quad | \underline{ZY}$
 $\quad | \underline{Yc}$
 $b \rightarrow ZY$
 $\quad | Yc$
 $c \rightarrow Za$

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

$f \rightarrow ID(n)$

$f \rightarrow ID()$

$n \rightarrow ID$

$n \rightarrow ID, n$

$i \rightarrow ID$

$L \rightarrow ($

$r \rightarrow)$

$c \rightarrow ,$

$f \rightarrow iLnr$

$f \rightarrow iLr$

$n \rightarrow ID$

$n \rightarrow icn$

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

$f \rightarrow iLnr \Rightarrow f \rightarrow iw \Rightarrow$

$w \rightarrow Lnr$

$x \rightarrow nr$

Left for you:

$f \rightarrow iLr$

$n \rightarrow icn$