

# CS 536 Announcements for Wednesday, February 28, 2024

Midterm 1, Thursday, February 29, 7:30 – 9 pm, S429 Chemistry

## Last Time

- review for Midterm 1

## Today

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

## Next Time

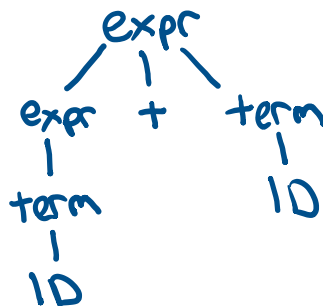
- wrap up CYK
- classes of grammars
- top-down parsing

In addition to printed copies of the overheads, there are full-sized versions of the diagrams for running the CYK algorithm available

## Parsing: two approaches

### Top-down / "goal driven"

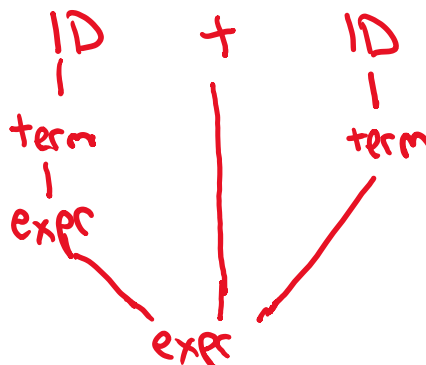
- start at start nonterminal
- grow parse tree downward until entire sequence is matched



### 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

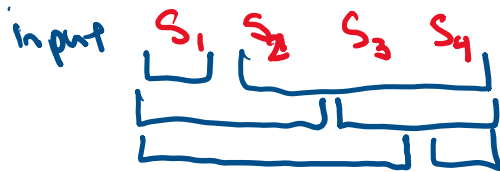
- Works bottom-up
- Time complexity:  $O(n^3)$   $n = \text{length of input (}\# \text{ of tokens in sequence)}$
- Requires grammar to be in Chomsky Normal Form

### 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?

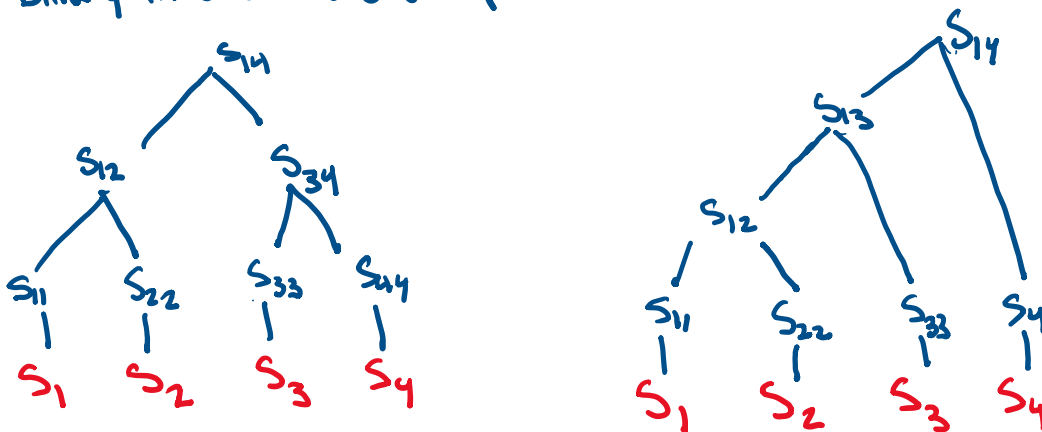
- nonterminals in pairs  $\rightarrow$  can think of a subtree as a subspan
- nonterminals (except start) can't derive epsilon  $\rightarrow$  each subspan has at least 1 token



### CYK : Dynamic Programming

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

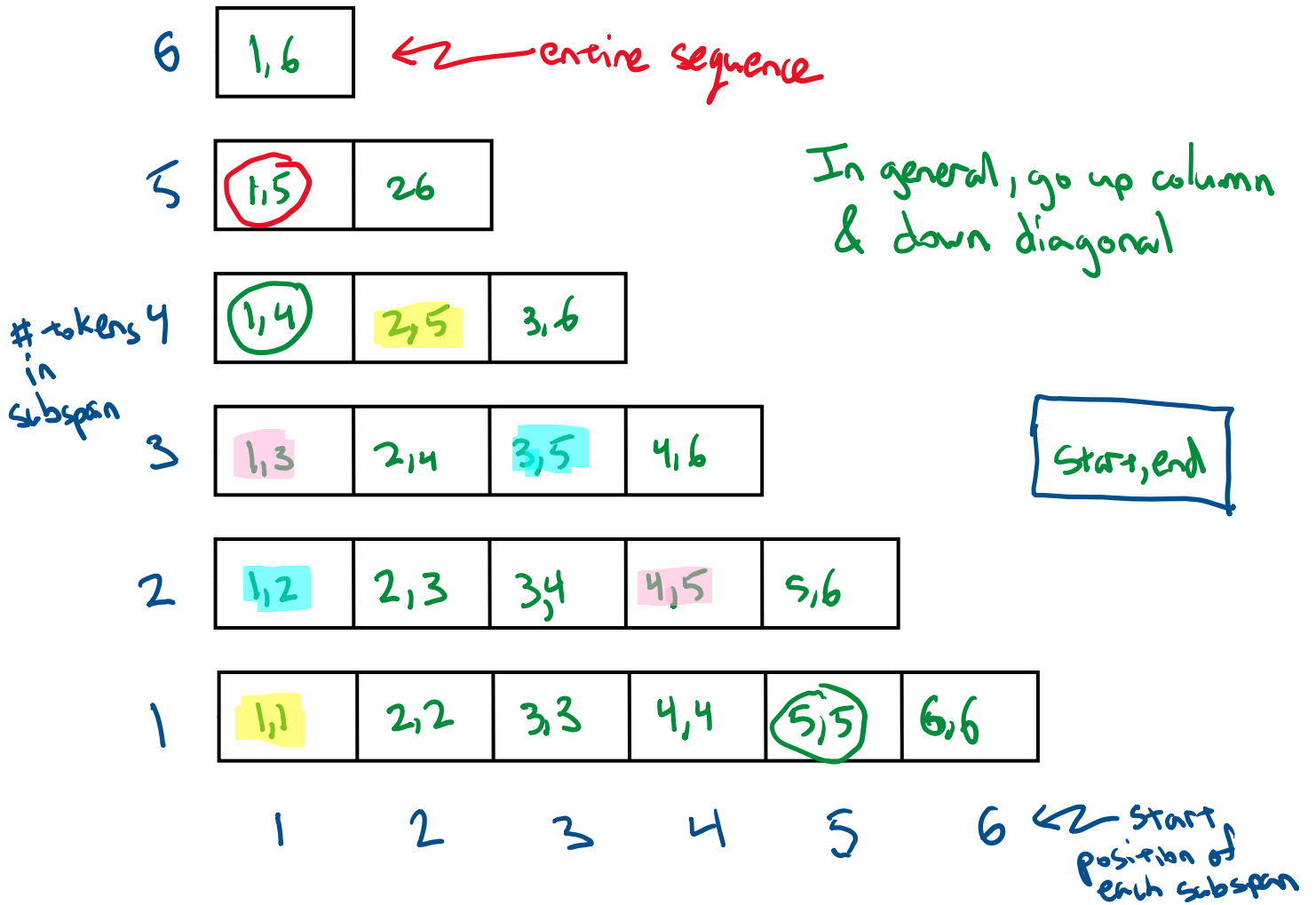
$x \rightarrow a b \rightarrow$  binary interior node of parse tree



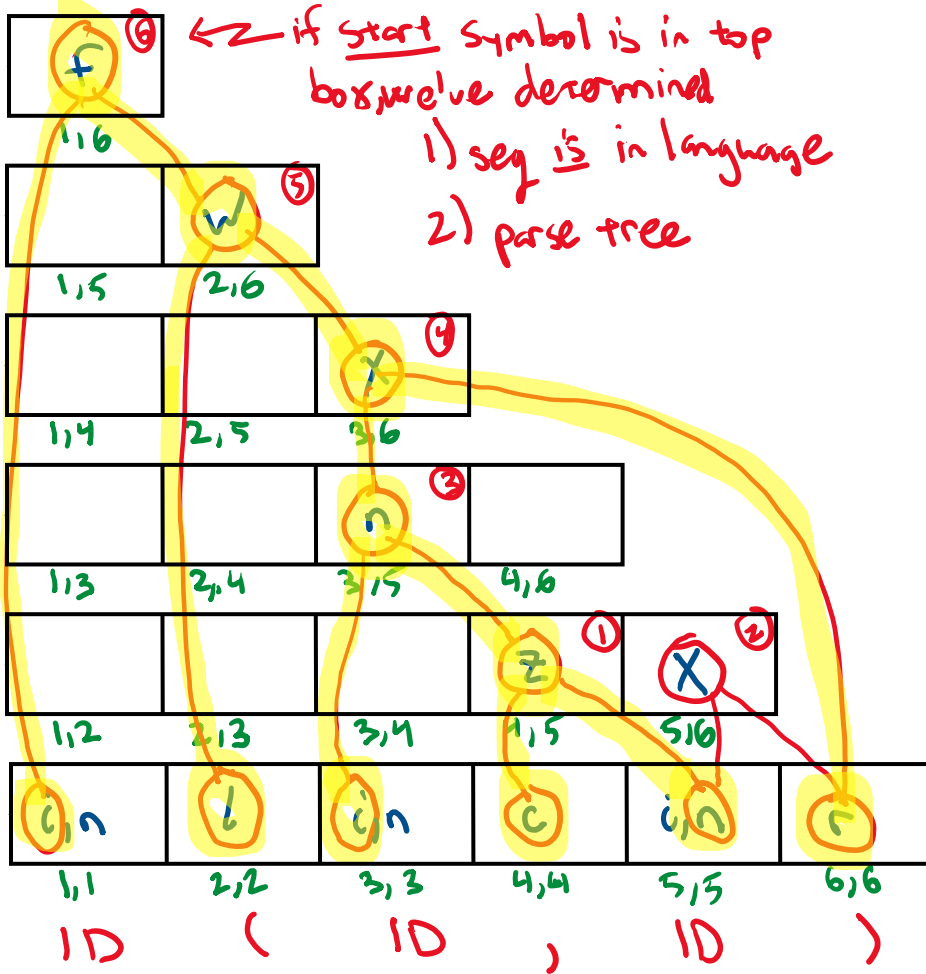
## Running CYK

Track every viable subtree from leaf to root.

All subspans for a sequence (string) with 6 terminals



# CYK Example



- f → iw (6)
- f → iy
- w → lx (5)
- x → nr (2) (1)
- 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

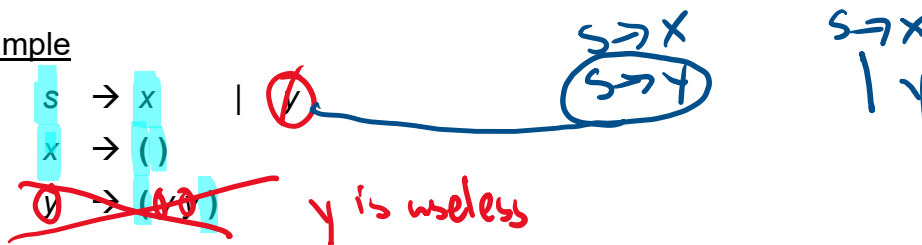
repeat

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



### Nonterminals that cannot be derived from the start symbol

mark the **start** symbol

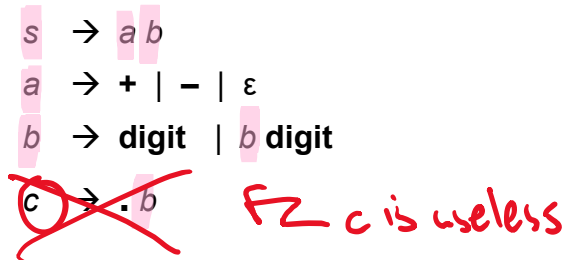
repeat

if the LHS of a production is marked

mark all **RHS nonterminals** (wherever they show up)

until no more nonterminals can be marked

#### Example

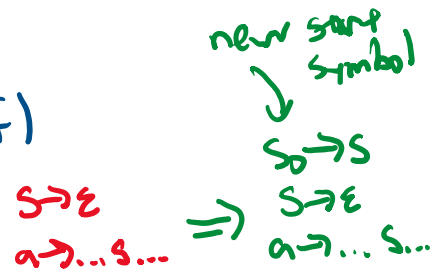


# Chomsky Normal Form

## Four steps

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

ok to have start  $S \rightarrow \epsilon$   
but if so can't have start 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 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 LHS,  
place  $b$  anywhere  $a$  could have appeared  
& remove unit production

#### Example

$$f \rightarrow ID(a)$$

$$f \rightarrow ID()$$

$$a \rightarrow n$$

$$n \rightarrow ID$$

$$n \rightarrow ID, n$$

$$f \rightarrow ID(n)$$

$$| \quad ID()$$

$$n \rightarrow ID$$

$$| \quad ID, n$$

If there are multiple rules with  $a$  on LHS,  
for each rule of the form  $b \rightarrow \delta$ ,  
add  $a \rightarrow \delta$  & remove  $a \rightarrow b$

$$a \rightarrow bX$$

$$| \quad cb$$

$$| \quad b$$

$$b \rightarrow ZY$$

$$| \quad Yc$$

$$c \rightarrow Za$$



$$a \rightarrow bX$$

$$| \quad cb$$

$$| \quad ZY$$

$$| \quad Yc$$

$$b \rightarrow ZY$$

$$| \quad Yc$$

$$c \rightarrow Za$$

## Chomsky Normal Form (cont.)

### Fix RHS nonterminals

For productions with terminals and something else on the RHS

- for terminal  $T$ , add rule  $x \rightarrow T$   
where  $x$  is a new non-terminal
- replace  $T$  with  $x$  in those productions

### Example

$f \rightarrow ID(n)$

$f \rightarrow ID()$

$n \rightarrow ID$

$n \rightarrow ID, n$

$i \rightarrow ID$

$l \rightarrow ($

$r \rightarrow )$

$c \rightarrow ,$

$f \rightarrow i l n r$

$f \rightarrow i l r$

$n \rightarrow ID$

$n \rightarrow i c n$

For productions with  $> 2$  nonterminals on the RHS

- replace all but 1 nonterm with new nonterm
- add rule from new nonterm to replaced nonterm sequence
- repeat

### Example

$f \rightarrow i l n r$

$\Rightarrow$

$f \rightarrow i w$

$\Rightarrow$

$f \rightarrow i w$

$w \rightarrow l \underline{x}$

$\underline{x} \rightarrow \underline{n r}$

Left for you:

$f \rightarrow i l r$

$n \rightarrow i c n$