

## CS 536 Announcements for Monday, March 4, 2024

### Last Time

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

can't derive a seq of tokens  
can't be derived from start non-term  
 $x \rightarrow T$      $x \rightarrow ab$     only start can derive  $\epsilon$

### Today

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

### Next Time

- building a predictive parser
- FIRST and FOLLOW sets

## Parsing (big picture)

### Context-free grammars (CFGs) Given CFG $G$

- language generation:  $G \rightarrow w \in L(G)$
- language recognition: given  $w$ , is  $w \in L(G)$ ?

### Translation

- given  $w \in L(G)$ , create a parse tree for  $w$
- given  $w \in L(G)$ , create an AST for  $w$

↳ passed on to next phase of our compiler

## CYK algorithm

Step 1: get grammar in Chomsky Normal Form (CNF)

Step 2: build all possible parse trees **bottom-up**

- start with runs of 1 terminal
- connect 1-terminal runs into 2-terminal runs
- connect 1- and 2-terminal runs into 3-terminal runs
- connect 1- and 3- or 2- and 2-terminal runs into 4-runs
- ...
- if we can connect entire tree, rooted at start symbol, we've found a valid parse

Pros: able to parse an arbitrary CFG *including ambiguous grammars*

Cons:  $O(n^3)$  time complexity *← too slow!*

For special classes of grammars, we can parse in  $O(n)$  time

*↳ eg LL(1) & LALR(1)*

## Classes of grammars

LL(1)  
 scan from L to R  
 ↳ 1 token lookahead  
 leftmost derivation

LALR(1)  
 look ahead (technical detail of LR parsers)  
 ↳ 1 token lookahead  
 rightmost derivation (in reverse)  
 scan L to R

CFG  
 ∪  
 LR(1)  
 ∪  
 LALR(1)  
 ∪  
 LL(1)  
 ∪  
 RG

Both are accepted by parser generators

LALR(1)

- parsed by bottom-up parsers
- harder to understand

*Java CUP generates a LALR(1) parser*

LL(1)

- parsed by top-down parsers

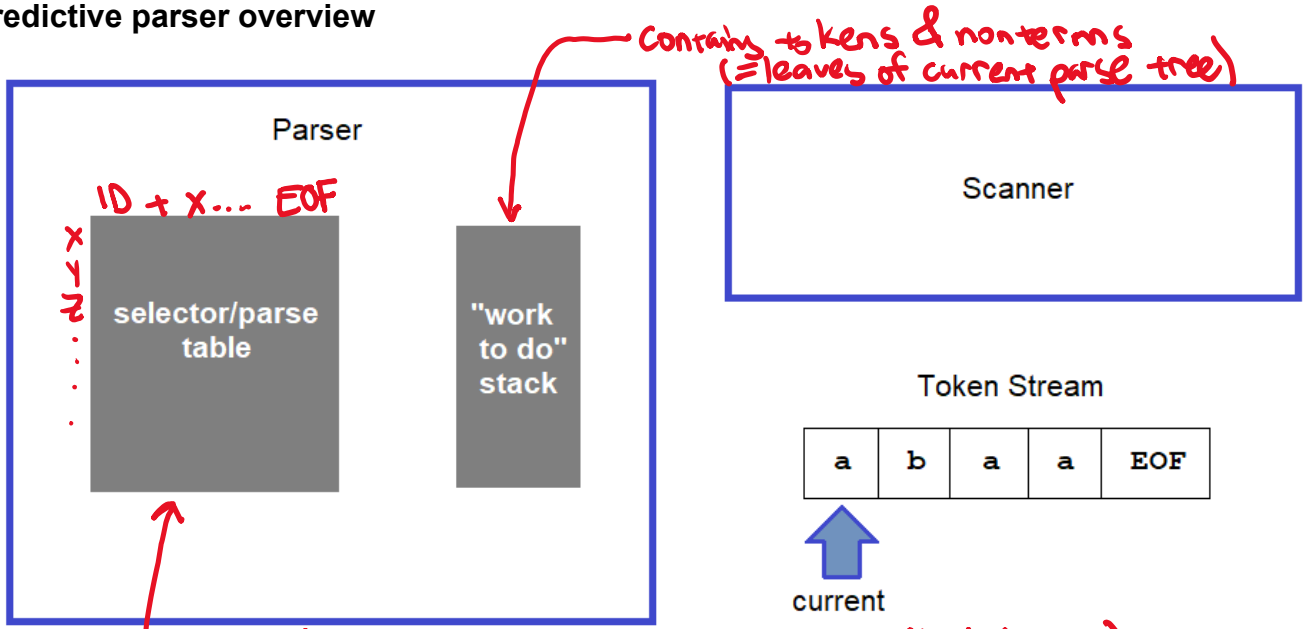
*↳ or predictive parsers or recursive descent parsers*

*= regular grammar - langs that can be recognized by DFAs*

# Top-down parsers

- Start at start symbol
- Repeatedly "predict" what production to use

## Predictive parser overview



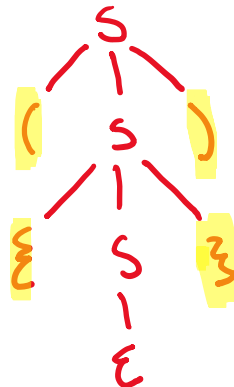
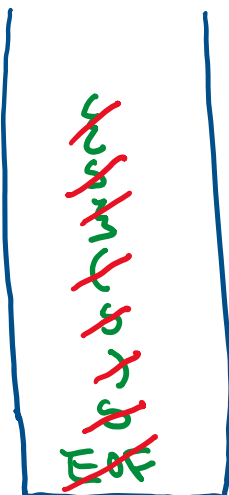
## Example

CFG:  $s \rightarrow (s) | \{s\} | \epsilon$

Parse table:

	(	)	{	}	EOF
s	(s)	$\epsilon$	{s}	$\epsilon$	$\epsilon$

Input: ( { } ) EOF  
 ↑ ↑ ↑ ↑ ↑



## Predictive parser algorithm

stack.push(**EOF**)

stack.push(start nonterm)

**T** = scanner.getToken( )

initial  
work  
Stack

Start  
EOF

repeat

if stack.top is terminal **Y**  
 match **Y** with **T**  
 pop **Y** from stack  
**T** = scanner.getToken()

if stack.top is nonterminal **x**  
 get table[x, current token **T**]  
 pop **x** from stack  
 push production's **RHS** (each symbol from R to L)

note: don't push  $\epsilon$

until one of the following:

stack is empty — accept input

stack.top is a terminal that does not match **T**

stack.top is a nonterm and parse-table entry is empty

reject input

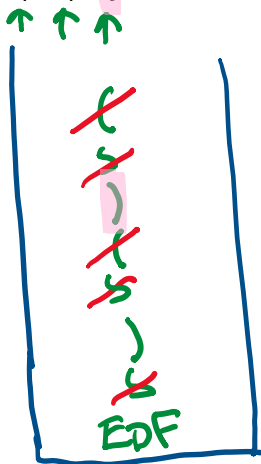
### Example

CFG:  $s \rightarrow (s) | \{s\} | \epsilon$

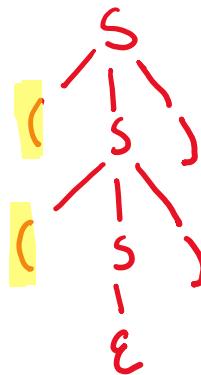
Parse table:

	(	)	{	}	EOF
s	(s)	$\epsilon$	{s}	$\epsilon$	$\epsilon$

Input: ( ( } EOF



no match  
↓  
reject input



## Consider

CFG:  $s \rightarrow (s) | \{s\} | ( | \{ | \epsilon$

Parse table:

	(	)	{	}	EOF
S					

(s) or ( ) ?

If could look ahead  $\geq$  tokens,  
we could make a good choice

This grammar is not LL(1), but it is LL(2)

Some grammars are not LL(k) for any k

### Two issues

- 1) How do we know if the language is LL(1)?
- 2) How do we build the selector table?

Answer: If we can build  
a parse table (selector table)  
& each entry has (at most) 1  
production in it, then the  
grammar is LL(1)

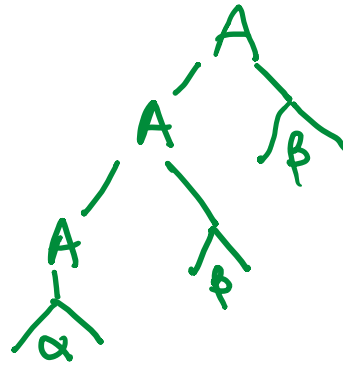


## Removing left-recursion

We can remove immediate left recursion without "changing" the grammar:

Consider:  $A \rightarrow A\beta$

$\alpha$   
 $\uparrow$   
 doesn't start with  
 non-term A



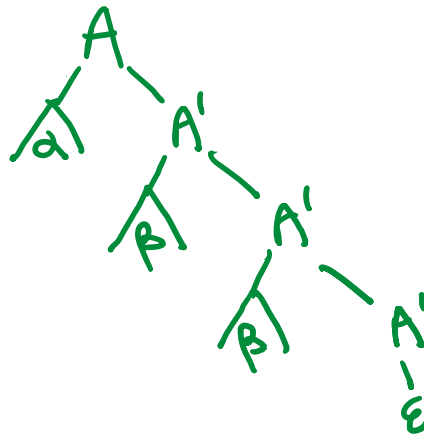
Solution: introduce new nonterminal  $A'$  and new productions:

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta A' \mid \epsilon$$

Messes up associativity  
of parse tree

(we'll fix this when  
we build the AST)



More generally,

$$A \rightarrow \underline{\alpha_1} \mid \underline{\alpha_2} \mid \dots \mid \underline{\alpha_n} \mid A\underline{\beta_1} \mid A\underline{\beta_2} \mid \dots \mid A\underline{\beta_p}$$

transforms to

$$A \rightarrow \underline{\alpha_1} A' \mid \underline{\alpha_2} A' \mid \dots \mid \underline{\alpha_n} A'$$

$$A' \rightarrow \underline{\beta_1} A' \mid \underline{\beta_2} A' \mid \dots \mid \underline{\beta_p} A' \mid \boxed{\epsilon}$$

## Grammars that are not left-factored

If a nonterminal has two productions whose right-hand sides have a common prefix, the grammar is not left-factored.

Example:  $s \rightarrow (s) \mid ()$

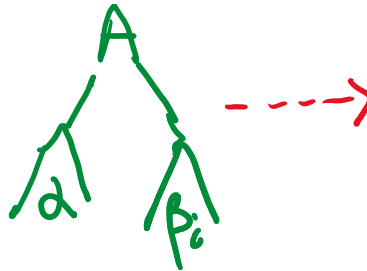
$S \rightarrow (s'$   
 $s' \rightarrow s) \mid )$

gets collapsed  
when AST is  
built

Given:  $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$

transform it to

$A \rightarrow \alpha A'$   
 $A' \rightarrow \beta_1 \mid \beta_2$



More generally,

$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_p$

transforms to

$A \rightarrow \alpha A' \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_p$   
 $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

### Combined example

$\text{exp} \rightarrow (\text{exp})$   
 $\mid \text{exp exp}$   
 $\mid ()$



$\text{exp} \rightarrow (\text{exp})\text{exp}' \mid ()\text{exp}'$   
 $\text{exp}' \rightarrow \text{exp exp}' \mid \epsilon$



$\text{exp} \rightarrow (\text{exp}''$   
 $\text{exp}'' \rightarrow \text{exp})\text{exp}' \mid )\text{exp}'$   
 $\text{exp}' \rightarrow \text{exp exp}' \mid \epsilon$