

CS3200 - More on CFGs + parsers

Note Title

2/3/2012

## Announcements

- HW due Saturday

## Other parsing algorithms

CYK is still pretty slow, especially for large programs.

After it was developed, a lot of work was put into figuring out what grammars could have faster algorithms.

Two big (& useful) classes have  $O(n)$  time parsers: LL & LR.

## LL & LR grammars

"LL" is left-to-right, leftmost derivation

"LR" is left-to-right, rightmost derivation

• So parser will scan left to right either way.

• LL will make a leftmost derivation  
(so right-leaning tree)

## LL versus LR

- LL are a bit simpler, so we'll start with them
- Note: LR is a larger class (so more grammars are LR than are LL)
- Both are used in production compilers today

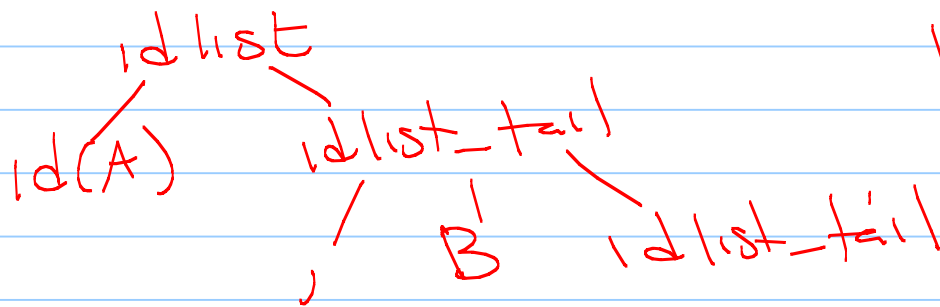
Example: LL parsing

$idlist \rightarrow id \ idlist\_tail$

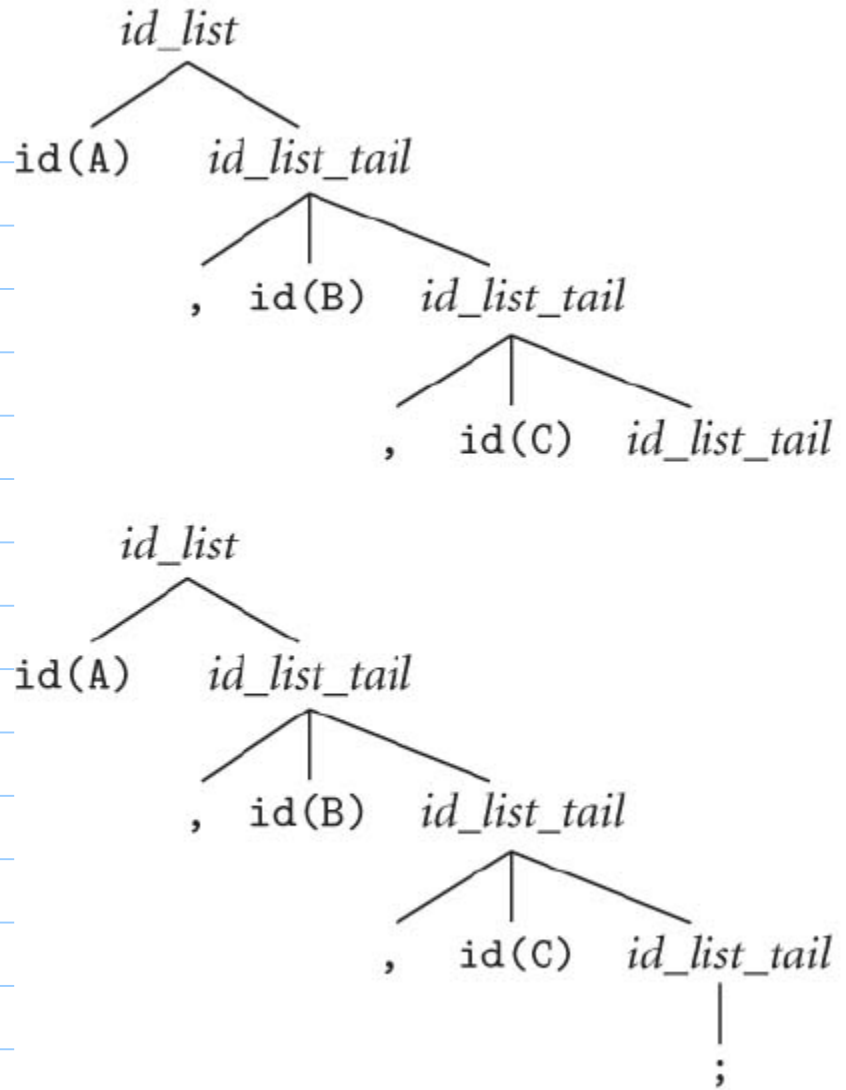
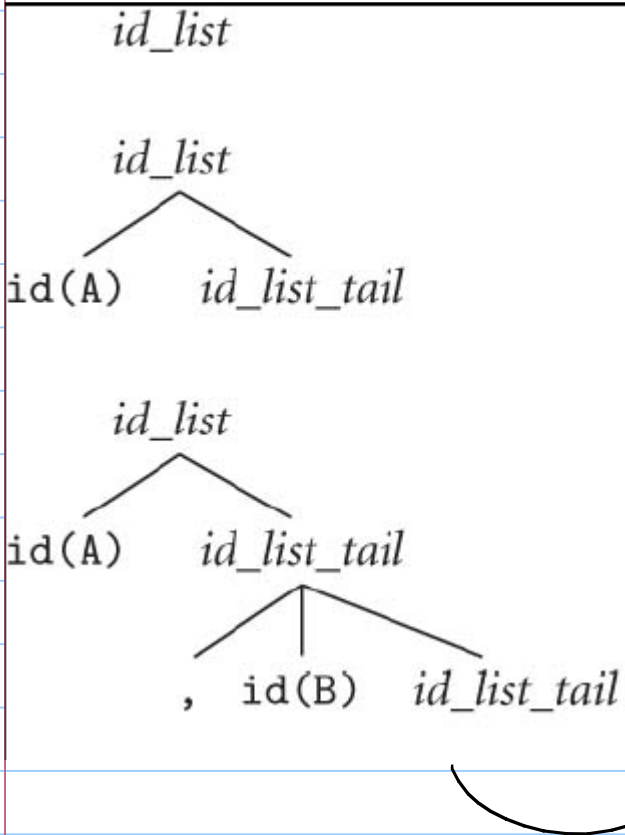
$idlist\_tail \rightarrow , \ id \ idlist\_tail$

$idlist\_tail \rightarrow ;$

Parse tree for "A,B,C;"  $\rightarrow id(A), id(B), id(C);$



$idlist \rightarrow A \ idlist\_tail$   
 $\rightarrow A; B \ idlist\_tail$



## LL(k) + LR(k)

When LL or LR is written with (1), (2), etc, it refers to how much look-ahead is allowed.

LL(1) means we can only look 1 token ahead when making our decision of which rule to match

Most commercial ones are LR(1), but exceptions exist (such as ANTLR).

A non LL(1) example: Left recursion

id\_list  $\rightarrow$  id  
 $\rightarrow$  id\_list, id

Imagine: Scanning left to right, &  
U encounter an id token. U

Which parse tree do we build?

A, B, C  
↑



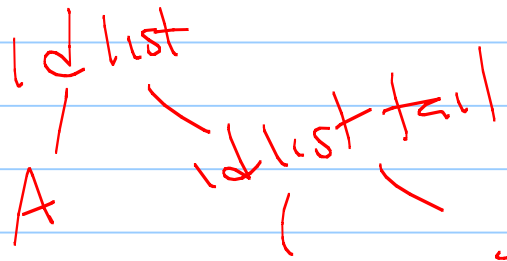
Making the grammar LL(1):

$\text{id\_list} \rightarrow \text{id id\_list\_tail}$

$\text{id\_list\_tail} \rightarrow , \text{id id\_list\_tails}$

$\rightarrow \epsilon$

A, B, C



Another non-LL(0) example: common prefixes  
PASCAL

stmt  $\rightarrow$  id := expr  
stmt  $\rightarrow$  id (argument\_list)

So when next token is an id,  
don't know which rule to use.

Fix?

• stmt  $\rightarrow$  id stmt\_tail

stmt\_tail  $\rightarrow$  := expr  
 $\rightarrow$  (arg\_list)

c = x+y;  
c(x);

Some grammars are non-LL:

- Eliminating left recursion and common prefixes is a very mechanical procedure which can be applied to any grammar.
- However, might not work! There are examples of inherently non-LL grammars.
- In these cases, generally add some heuristic to deal with odd cases

## PASCAL if statement

Example: non-LL language: optional else  
stmt  $\rightarrow$  if condition then\_clause else\_clause

then\_clause  $\rightarrow$  then stmt

else\_clause  $\rightarrow$  else stmt  
 $\rightarrow \epsilon$

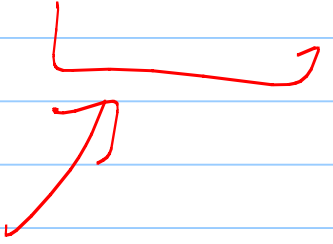
What syntax?

## if statements

Ex: if  $C_1$  then if  $C_2$  then  $S_1$  else  $S_2$

Parse tree:

inherently  
ambiguous



## Back to LL-parsing

We have seen mostly top-down parsing.

Start with  $S_0$ , the start token, & try to construct the tree based on the next input.

Also called predictive parsing - matches the rule based on current token/state plus the next input!  
or next  $k$  inputs

## LR grammars

Bottom-up parsing starts at the leaves (here, the tokens), & tries to build the tree upward.

Continues scanning & shifting tokens onto a forest, then builds up when it finds a valid production.

Never predicts - when it recognizes right hand side of a rule, simplifies to left hand side.

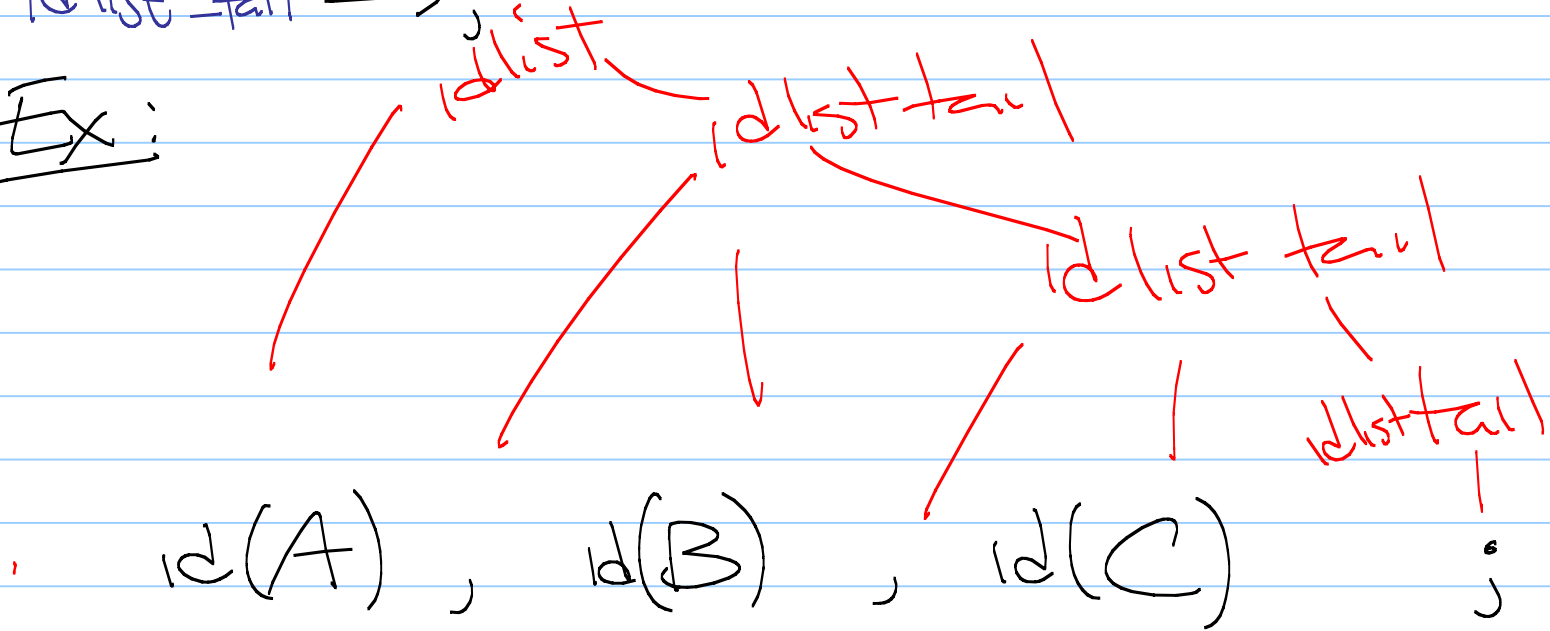
# Bottom-up parsing

$idlist \rightarrow id \quad idlist\_tail$

$idlist\_tail \rightarrow , \quad id \quad idlist\_tail$

$idlist\_tail \rightarrow ;$

Ex:







## Shift-reduce:

- Bottom up parsers are also called shift-reduce:
  - Shift token onto stack (in a forest)
  - when a rule is recognized, reduce to left-hand side
- Problem with last example:  
must shift all tokens onto the forest before reducing.  
What could happen in a large program?  
slow, might crash/overflow
- Sometimes unavoidable. However, sometimes other options...

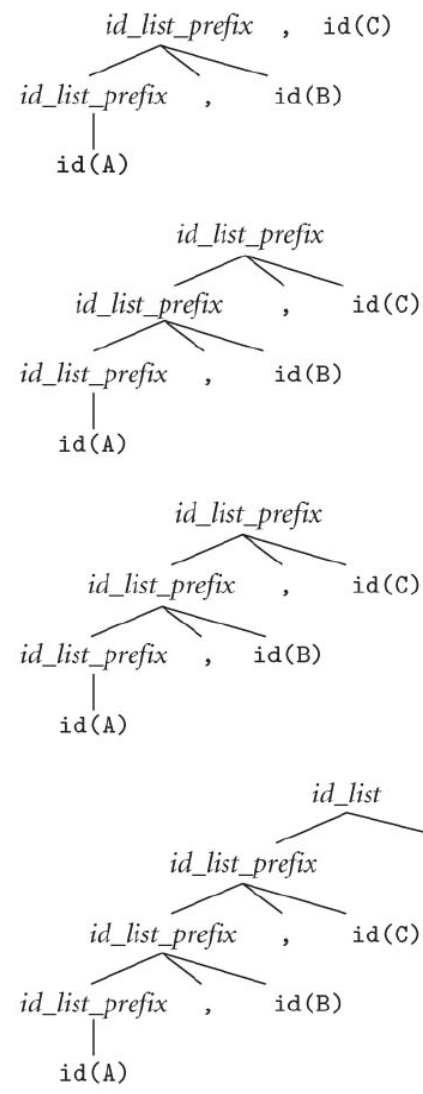
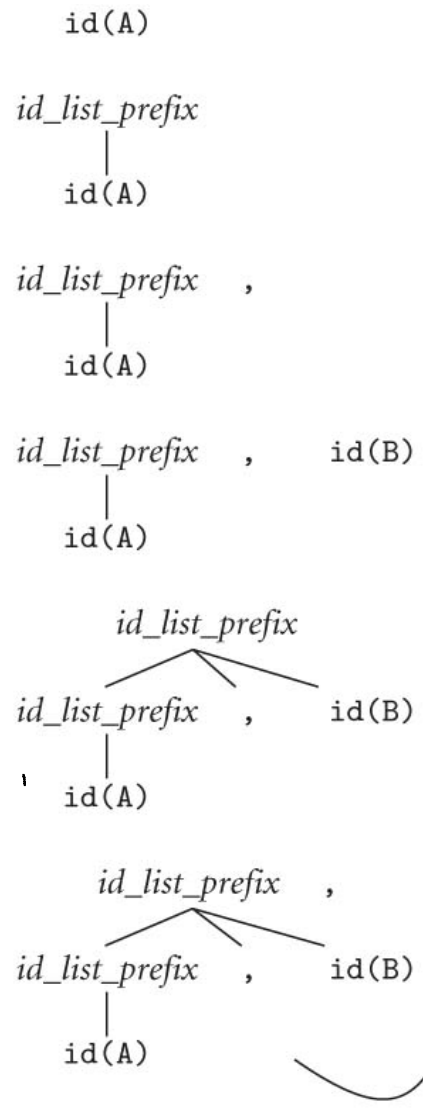
Bottom-up parsings: another example

$\text{id\_list} \rightarrow \text{id\_list\_prefix};$

$\text{id\_list\_prefix} \rightarrow \text{id\_list\_prefix}, \text{id}$   
 $\quad \quad \quad \rightarrow \text{id}$

Parse A, B, C; again, bottom-up:





## Bottom-up parsing: some notes

- The previous example cannot be parsed top-down. Why?  
not LL: left recursion, ...
- Note that it also is not an LL grammar, although the language is LL.
- There is a distinction between a language & a grammar.  
Remember, any language can be generated by an infinite number of grammars.

## LR grammars : An old example

$\text{expr} \rightarrow \text{term} \mid \text{expr add\_op term}$

$\text{term} \rightarrow \text{factor} \mid \text{term multop factor}$

$\text{factor} \rightarrow \text{id} \mid \text{number} \mid - \text{factor} \mid (\text{expr})$

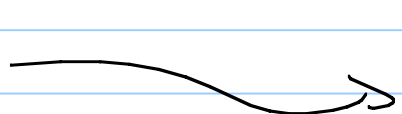
$\text{add\_op} \rightarrow + \mid -$

$\text{multop} \rightarrow * \mid /$

This grammar is not LL!

- If we get an `id` as input when expecting an `expr`, no way to choose between the 2 possible productions.

- It suffers from the common prefix issue we saw before.

(We can fix this )

Another LL-example:

$\text{expr} \rightarrow \text{term term\_tail}$

$\text{term\_tail} \rightarrow \text{add\_op term term\_tail}$   
 $\rightarrow \epsilon$

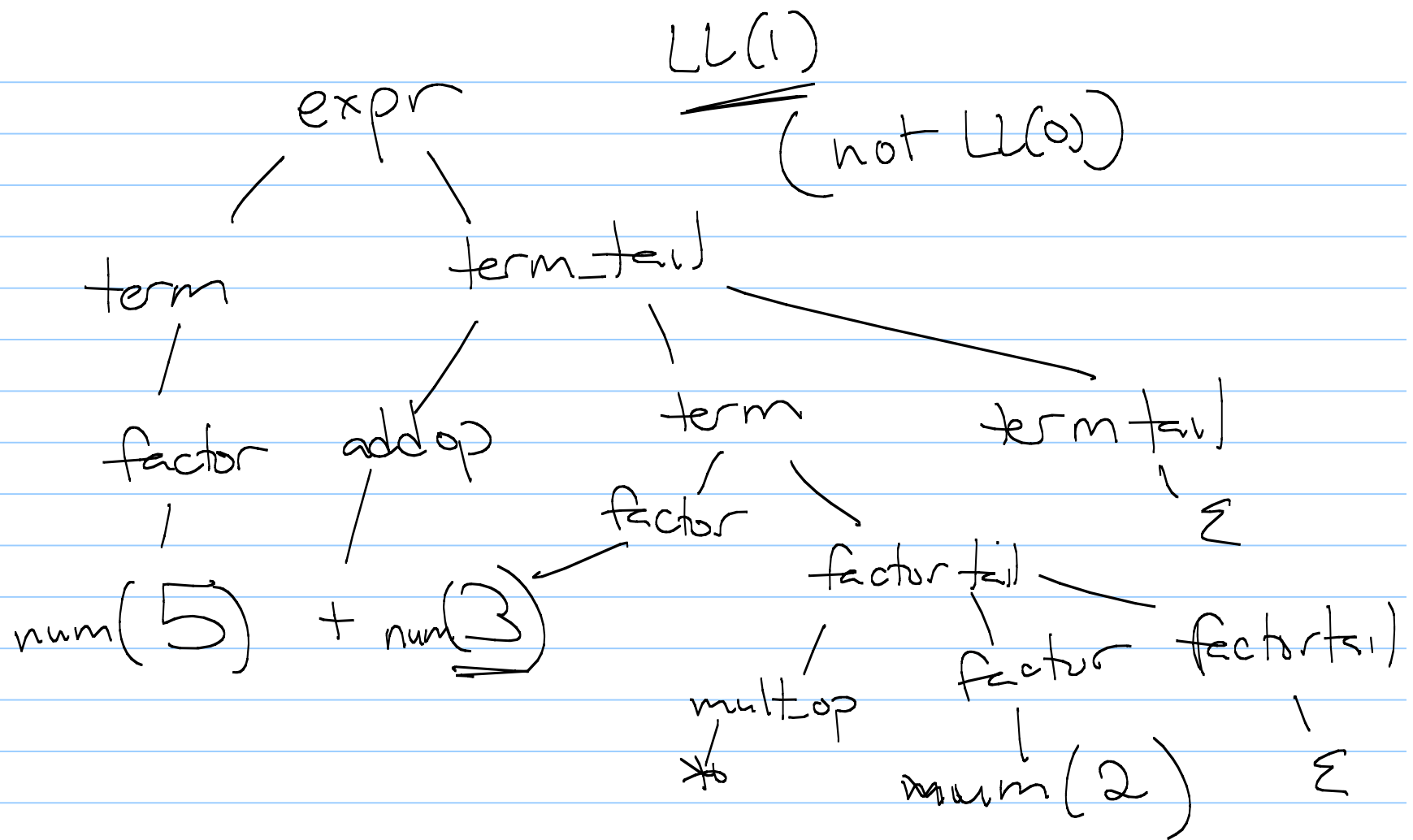
$\text{term} \rightarrow \text{factor factor\_tail}$

$\text{factor\_tail} \rightarrow \text{mult\_op factor factor\_tail}$   
 $\rightarrow \epsilon$

$\text{factor} \rightarrow (\text{expr}) \mid \text{id} \mid \underline{\text{number}}$

$\text{add\_op} \rightarrow + \mid -$   
 $\text{mult\_op} \rightarrow * \mid /$





Now can add this as part  
of a simple calculator  
language:

program  $\rightarrow$  stmt\_list \$\$  $\leftarrow$  end of file

stmt\_list  $\rightarrow$  stmt stmt\_list  
 $\rightarrow$   $\epsilon$

stmt  $\rightarrow$  id := expr  
 $\rightarrow$  read id  
 $\rightarrow$  write expr

Program: What does it do?

read A

read B

sum := A + B

write sum

write sum / 2

How to parse?

pgm  $\rightarrow$  stmt\_list

$\rightarrow$  stmt stmt\_list

$\rightarrow$  . stmt\_list

