

CS 536 Announcements for Thursday, February 6, 2025

Programming Assignment 2

- has been released
- due Tuesday, February 18

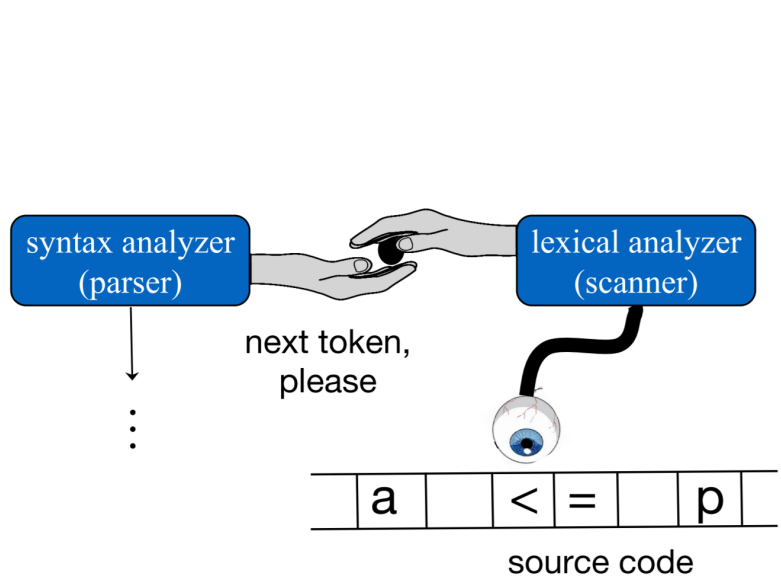
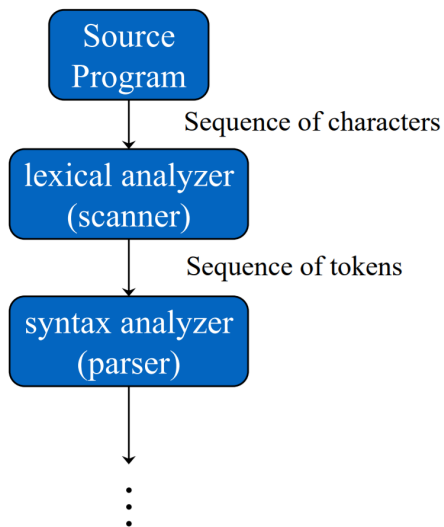
Last Time

- non-deterministic FSMs
- equivalence of NFAs and DFAs
- regular expressions
- regular languages
- regular expressions → DFAs
- language recognition → tokenizers
- scanner generators
- JLex

Today

- CFGs
- Makefiles
- resolving ambiguity
- expression grammars
- list grammars

Recall big picture



Conceptual organization

Why regular expressions are not good enough

Regular expression wrap-up

- + perfect for tokenizing a language
- limitations
 - define only limited family of languages
 - can't be used to specify all the programming constructs we need
 - no notion of structure

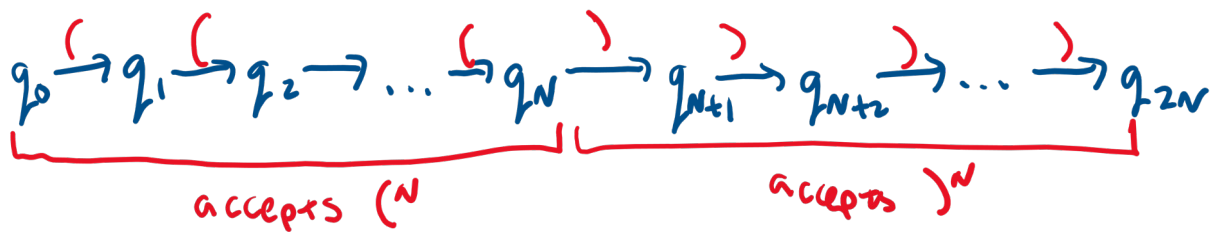
Regexs cannot handle "matching"

Example: $L_{()} = \{ (^n \text{ where } n > 0) \} = \{ "()", "(())", "((()))", \dots \}$

Theorem: No regex/DFA can describe the language $L_{()}$

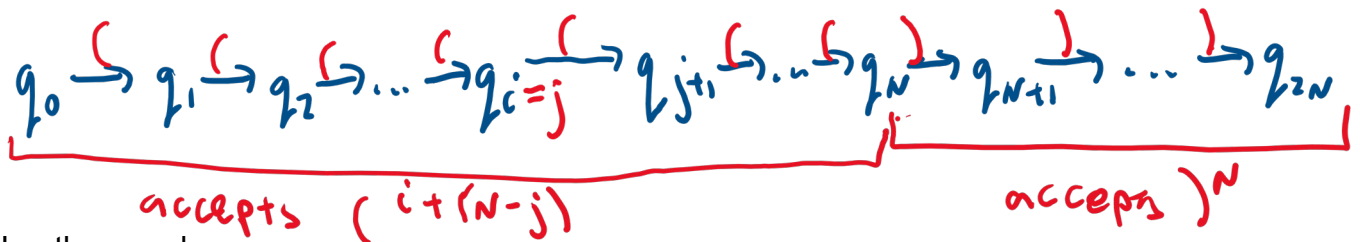
Proof by contradiction: Suppose there exists a DFA A for $L_{()}$ where A has N states.

Then A has to accept the string $(^N)^N$ with some sequence of states



By the pigeonhole principle, there exists $i, j \leq N$ where $i < j$ such that $q_i = q_j$

So



In other words,

A accepts $(^{N-j+i})^N$ but $(^{N-j+i})^N \notin L_{()}$
which is a contradiction

No notion of structure

$$x = y + z$$

Consider the following stream of tokens: ID ASSIGN ID PLUS ID

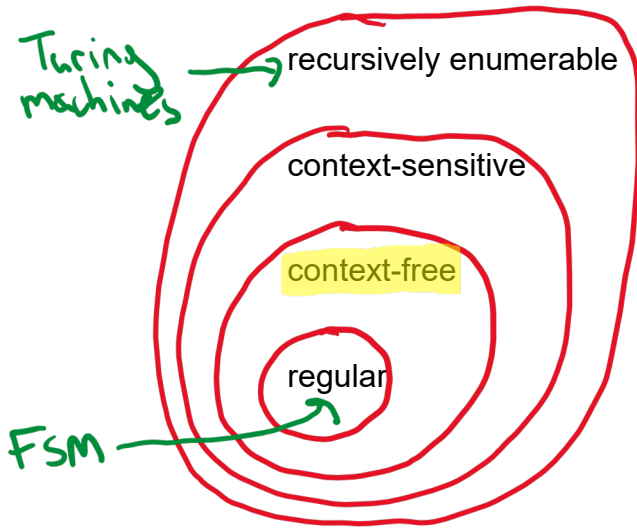
What should be done 1st? $(x=y)+z$ or $x=(y+z)$

What about precedence & associativity

(Noam)

The Chomsky Language Hierarchy

Language class:



power

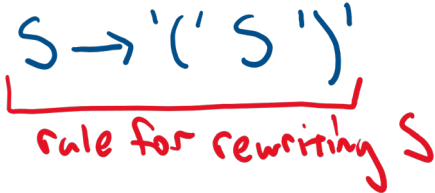


efficient



Context-free grammar (CFG)

= a set of recursive rewriting rules to generate patterns of strings



Before applying
S



Formal definition: A CFG is a 4-tuple (N, Σ, P, S)

- N = set of **non-terminals** - placeholders (interior nodes in parse tree)
- Σ = set of **terminals** - tokens from scanner
- P = set of **productions** - rules for re-writing nonterminals (ie for deriving)
- S = initial non-terminal symbol ("**start symbol**"), $S \in N$
- if not otherwise specified, use non-term on LHS of 1st production as start symbol

Productions

Production syntax : LHS \rightarrow RHS

Single non-terminal \uparrow \uparrow expression (seq of terms & non-terms) or ϵ

non-term \rightarrow expression
non-term $\rightarrow \epsilon$

$S \rightarrow '(S)'$
 $S \rightarrow \epsilon$

or (assuming non-terms on LHS are the same)

non-term \rightarrow expression
| ϵ

$S \rightarrow '(S)'$
| ϵ

or
non-term \rightarrow expression | ϵ

$S \rightarrow '(S)'$ | ϵ

Language defined by a CFG

= set of strings (i.e., sequences of terminals) that can be **derive**d from the **start** non-terminal

To **derive** a string (of terminal symbols):

- set Curr_Seq to start symbol
- repeat
 - find a non-terminal **x** in Curr_Seq
 - find production of the form $x \rightarrow \alpha$
 - "apply" production: create new Curr_Seq by replacing x with α
- until Curr_Seq contains no non-terminals

Derivation notation

- derives \Rightarrow
- derives in one or more steps \Rightarrow^+ or $\stackrel{+}{\Rightarrow}$
- derives in zero or more steps \Rightarrow^* or $\stackrel{*}{\Rightarrow}$

$L(G)$ = language defined by CFG G

= $\{w \mid S \stackrel{+}{\Rightarrow} w \text{ where } S \text{ is start non-terminal \& } w \text{ is seq of terms or } \epsilon \}$
 \uparrow "such that"

Example grammar

Terminals

BEGIN } program
END } boundary
SEMICOLON - ";" to separate statements
ASSIGN - "=" in asg stmts
ID - identifier (variable name)
PLUS - "+" operator in expression

Non-terminals

prog - start non-term (root of parse tree)
stmts - list of statements
stmt - a single statement
expr - a (mathematical expression)

Productions - define syntax of legal programs

- 1) prog → BEGIN stmts END
- 2) stmts → stmts SEMICOLON stmt
- 3) | stmt
- 4) stmt → ID ASSIGN expr
- 5) expr → ID
- 6) | expr PLUS ID

Example derivation

Productions

- 1) prog \rightarrow BEGIN stmts END
- 2) stmts \rightarrow stmts SEMICOLON stmt
- 3) | stmt
- 4) stmt \rightarrow ID ASSIGN expr
- 5) expr \rightarrow ID
- 6) | expr PLUS ID

Derivation

prog

\Rightarrow BEGIN stmts END ^①

\Rightarrow BEGIN stmts SEMICOLON stmt END ^②

\Rightarrow BEGIN stmt SEMICOLON stmt END ^③

\Rightarrow BEGIN ID ASSIGN expr SEMICOLON stmt END ^④

\Rightarrow BEGIN ID ASSIGN expr SEMICOLON ID ASSIGN expr END ^④

\Rightarrow BEGIN ID ASSIGN ID SEMICOLON ID ASSIGN expr END ^⑤

\Rightarrow BEGIN ID ASSIGN ID SEMICOLON ID ASSIGN expr PLUS ID END ^⑥

\Rightarrow BEGIN ID ASSIGN ID SEMICOLON ID ASSIGN ID PLUS ID END ^⑤

BEGIN ID = ID ; ID = ID + ID END

prog $\stackrel{+}{\Rightarrow}$ BEGIN ID = ID ; ID = ID + ID END

Q: is BEGIN ID = ID ; END in language?

No - but would be if we added stmt $\rightarrow \epsilon$

Start again at 7:35 pm

Parse trees

= way to visualize a derivation

To derive a string (of terminal symbols):

- set **root** of parse tree to **start** symbol
- repeat
 - find a **leaf** non-terminal **x**
 - find production of the form $x \rightarrow \alpha$
 - "apply" production: **symbols in α** become the **children of x**
- until there are no more leaf non-terminals

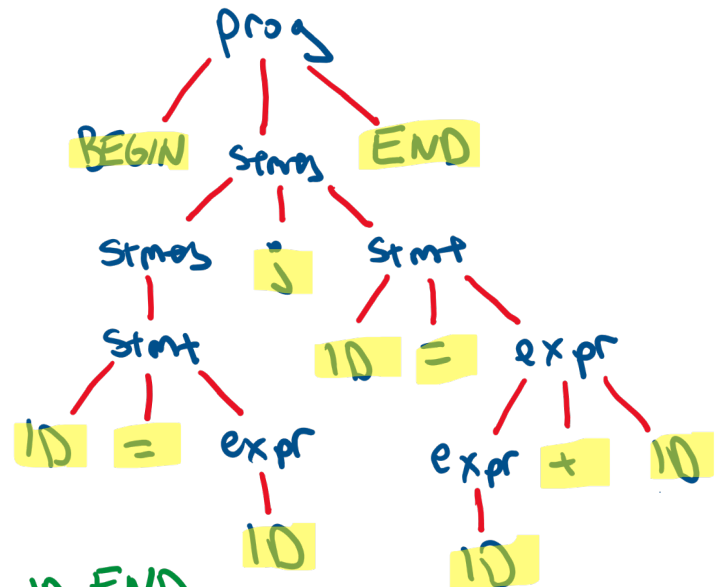
Derived sequence determined from leaves, from left to right

Productions

- 1) prog \rightarrow BEGIN stmts END
- 2) stmts \rightarrow stmts SEMICOLON stmt
- 3) | stmt
- 4) stmt \rightarrow ID ASSIGN expr
- 5) expr \rightarrow ID
- 6) | expr PLUS ID

↑
this notation is
BNF (or extended BNF)
Backus-Naur Form

BEGIN ID = ID; ID = ID + ID END



Makefiles

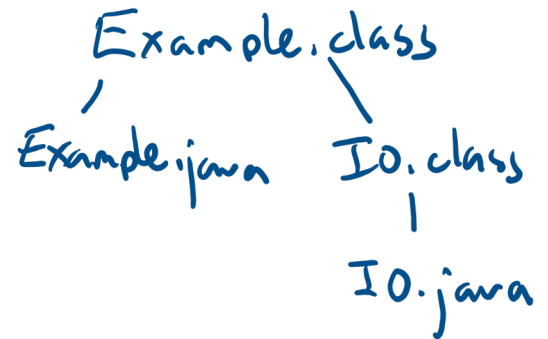
Basic structure

```
<target>: <dependency list>  
  <command to satisfy target>
```

tab

Example

```
Example.class: Example.java IO.class  
  javac Example.java  
  
IO.class: IO.java  
  javac IO.java
```



Make creates an internal **dependency graph**

- a file is rebuilt if one of its dependencies changes

Variables – for common configuration values to use throughout your makefile

Example

```
JC = /s/std/bin/javac  
JFLAGS = -g ← build for use with debugger  
  
Example.class: Example.java IO.class  
  $(JC) $(JFLAGS) Example.java  
  
IO.class: IO.java  
  $(JC) $(JFLAGS) IO.java
```

Phony targets

- target with no dependencies = "phony"
- use make to run commands:

Example

```
clean:  
  rm -f *.class
```

force remove

```
test:  
  java Example inFile.txt outFile.txt  
  java Example inErrFile.txt outErrFile.txt
```


Programming Assignment 2

Modify:

- bach.jlex
- P2.java
- Makefile

Makefile

```
###
# testing - add more here to run your tester and compare
# its results to expected results
###
test:
    java -cp $(CP) P2
    diff allTokens.in allTokens.out

###
# clean up
###

clean:
    rm -f *~ *.class bach.jlex.java

cleantest:
    rm -f allTokens.out
```

Run make to compile
(by default make
does 1st target in
Makefile)

Running the tester

```
vm-instunix-07(53)% make test
java -cp ./deps:. P2
3:1 ****ERROR**** ignoring illegal character: a
diff allTokens.in allTokens.out
3d2
< a ] → output of diff command
make: *** [Makefile:40: test] Error 1
```

Commands
from
Makefile

error msg produced by
bach scanner when
P2 is run

← from running make

terminals \equiv tokens

CFG review

formal definition: CFG $G = (N, \Sigma, P, S)$

CFG generates a string by applying productions until no non-terminals remain

\Rightarrow^+ means "derives in 1 or more steps"

$q \Rightarrow (q) \Rightarrow (\epsilon)$ ie $()$

language defined by a CFG G

$L(G) = \{ w \mid s \Rightarrow^+ w \}$ where

s = start is the start non-terminal of G , an

w = sequence consisting of (only) terminal symbols or ϵ

$L(G) = \{ \epsilon, (), (()), ((())), \dots \}$

Example: nested parens

$N = \{ q \}$

$\Sigma = \{ '(', ')' \}$

$P = q \rightarrow (q) \mid \epsilon$

$q \Rightarrow (q)$

$\Rightarrow ((q))$

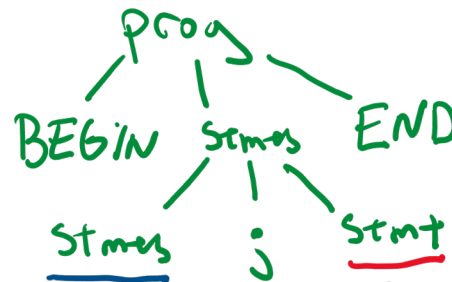
$\Rightarrow (((q)))$

$\Rightarrow (((\epsilon)))$

$q \Rightarrow^+ (((\epsilon)))$

Derivation order

- 1) prog \rightarrow BEGIN stmts END
- 2) stmts \rightarrow stmts SEMICOLON stmt
- 3) | stmt
- 4) stmt \rightarrow ID ASSIGN expr
- 5) expr \rightarrow ID
- 6) | expr PLUS ID



Leftmost derivation : leftmost non-terminal is always expanded

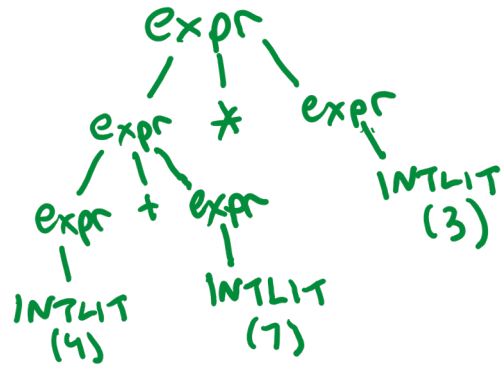
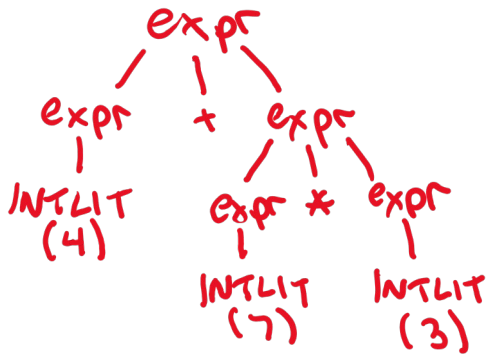
Rightmost derivation : rightmost non-terminal is always expanded

Expression Grammar Example

- 1) $\text{expr} \rightarrow \text{INTLIT}$
- 2) $\text{expr} \rightarrow \text{expr PLUS expr}$
- 3) $\text{expr} \rightarrow \text{expr TIMES expr}$
- 4) $\text{expr} \rightarrow \text{LPAREN expr RPAREN}$

Goal: create a CFG for arithmetic expressions involving only +, *, parens, & integer literals

Derive: $4 + 7 * 3$



ambiguous grammar!

For grammar G and string w , G is **ambiguous** if there is

> 1 leftmost derivation of w

OR

> 1 rightmost derivation of w

OR

> 1 parse tree for w

these
are all
equivalent

Grammars for expressions

Goal: write a grammar that correctly reflects precedences and associativities

$$a + b * c \leftrightarrow a + (b * c)$$

$$a + b + c \leftrightarrow (a + b) + c$$

$$a = b = c \leftrightarrow a = (b = c)$$

Precedence

- use different non-terminal for each precedence level
- start by re-writing production for lowest precedence operator first

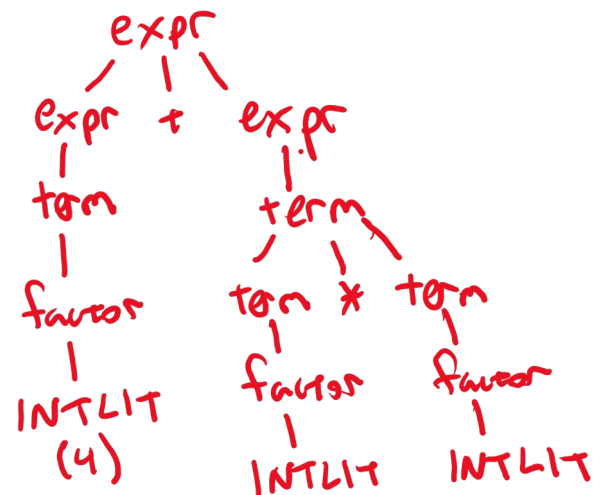
Example

- 1) $\text{expr} \rightarrow \text{INTLIT}$
- 2) $\text{expr} \rightarrow \text{expr PLUS expr}$
- 3) $\text{expr} \rightarrow \text{expr TIMES expr}$
- 4) $\text{expr} \rightarrow \text{LPAREN expr RPAREN}$

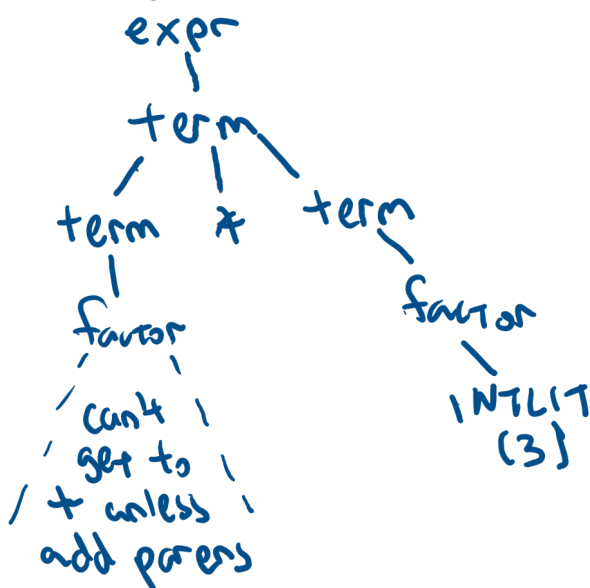
$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{expr} \\ &\quad | \text{term} \\ \text{term} &\rightarrow \text{term} * \text{term} \\ &\quad | \text{factor} \\ \text{factor} &\rightarrow \text{INTLIT} \\ &\quad | (\text{expr}) \end{aligned}$$

+ has lowest precedence

$$4 + 7 * 3$$



Consider



Grammars for expressions (cont.)

What about associativity? Consider $1 + 2 + 3$ equiv $(1+2)+3$



Definition: recursion in grammars

A grammar is **recursive** in non-terminal x if $x \Rightarrow^+ \alpha x \gamma$ for non-empty strings of symbols α and γ

A grammar is **left-recursive** in non-terminal x if $x \Rightarrow^+ x \gamma$ for non-empty string of symbols γ

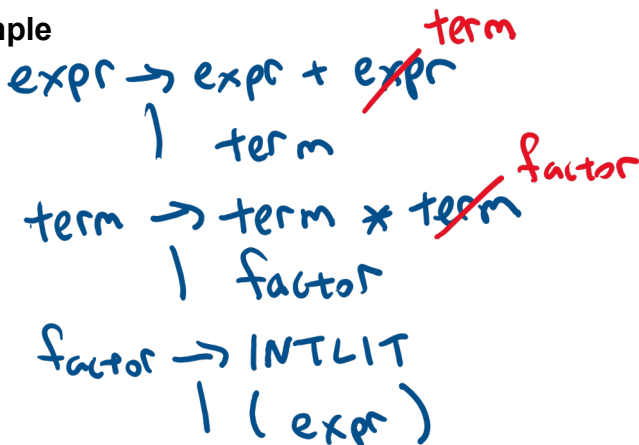
A grammar is **right-recursive** in non-terminal x if $x \Rightarrow^+ \alpha x$ for non-empty string of symbols α

In expression grammars

for left associativity, use left recursion

for right associativity, use right recursion

Example



AST for $1+2+3$



left associative

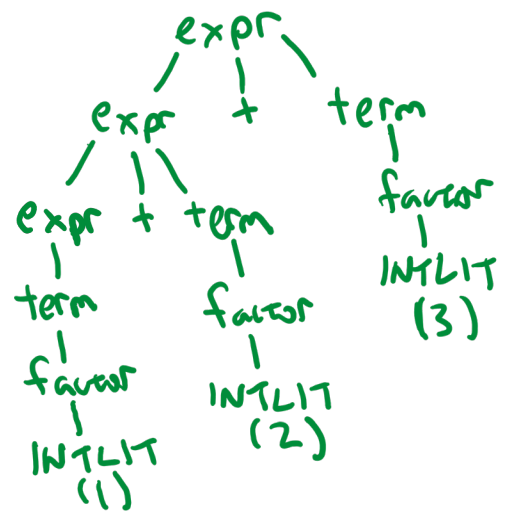
$+ - * /$

right associative

$= \wedge \leftarrow \text{exponentiation}$

$$2^3^4 \equiv 2^{3^4}$$

$1+2+3$



Extend this grammar to add exponentiation (POW)

Add exponentiation (POW) to this grammar, with the correct precedence and associativity.

expr \rightarrow expr PLUS term
| term

term \rightarrow term TIMES factor
| factor

~~factor \rightarrow INTLIT~~
~~| LPAREN expr RPAREN~~

factor \rightarrow exponent POW factor
| exponent

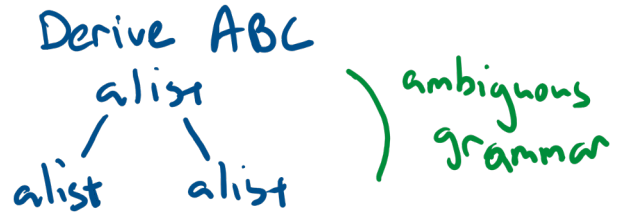
exponent \rightarrow INTLIT
| LPAREN expr RPAREN

\hookrightarrow highest \hookrightarrow right

List grammars

Example a list with no separators, e.g., A B C D E F G

$alist \rightarrow ITEM$
 $| alist\ alist$

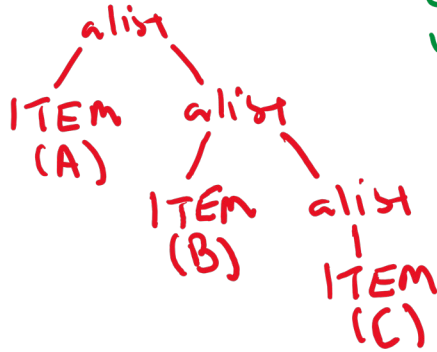


$alist \rightarrow ITEM$
 $| ITEM\ alist$

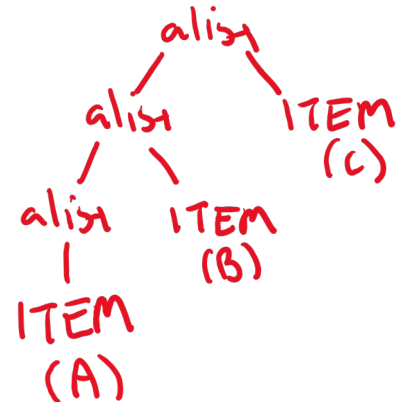
OR

$alist \rightarrow ITEM$
 $| alist\ ITEM$

Derive ABC



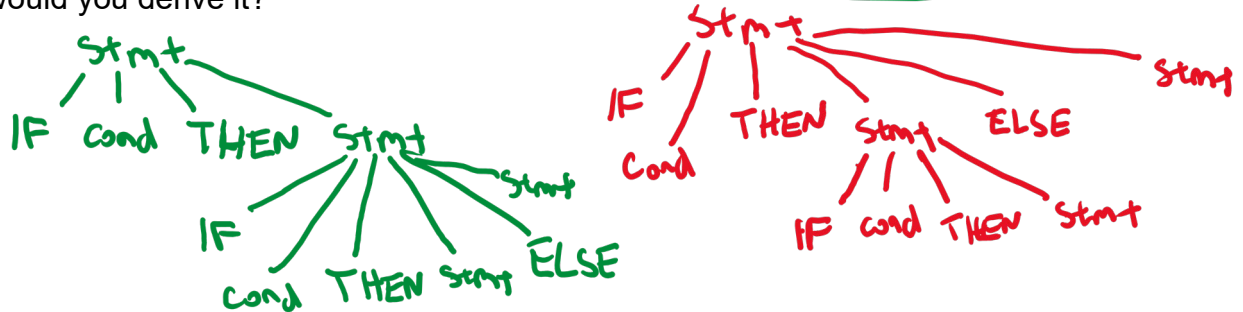
Associativity
 doesn't matter
 with lists so
 either grammar
 is fine



Another ambiguous example

$stmt \rightarrow$ IF cond THEN stmt
 $|$ IF cond THEN stmt ELSE stmt
 $|$...

Given this word in this grammar: if a then if b then s1 else s2
 How would you derive it?



$stmt \rightarrow$ IF cond THEN LCHRLY stmt RCHRLY
 $|$ IF cond THEN LCHRLY stmt RCHRLY ...
 ELSE LCHRLY stmt RCHRLY