

CS143: Parsing I

David L. Dill

Stanford University

Parsing (and lexical analysis wrapup)

- Practical Issues in Lexical Analysis (finish)
- Parsing
 - Context-Free Grammars
 - Parse Trees and Languages
 - Ambiguity
 - Mathematical Definitions
 - Precedence and Associativity
 - Extended BNF

Start Conditions

A better way to handle comments

When you see "/*"

Use new DFA

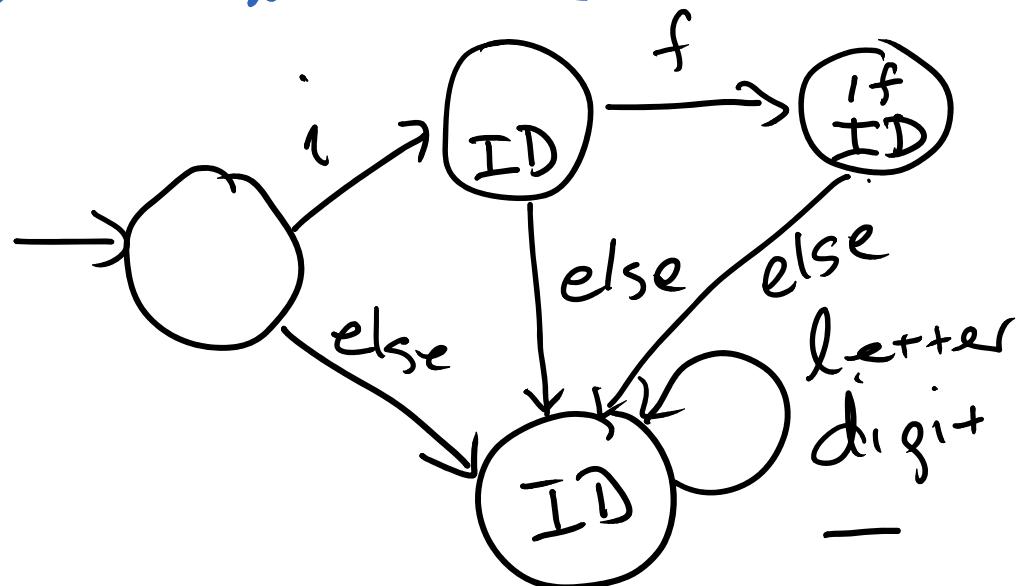
Patterns:

- any single char - discard
- "*/" - discard, return to main DFA

Rule Ordering

if if

ID letter(letter | digit | -)



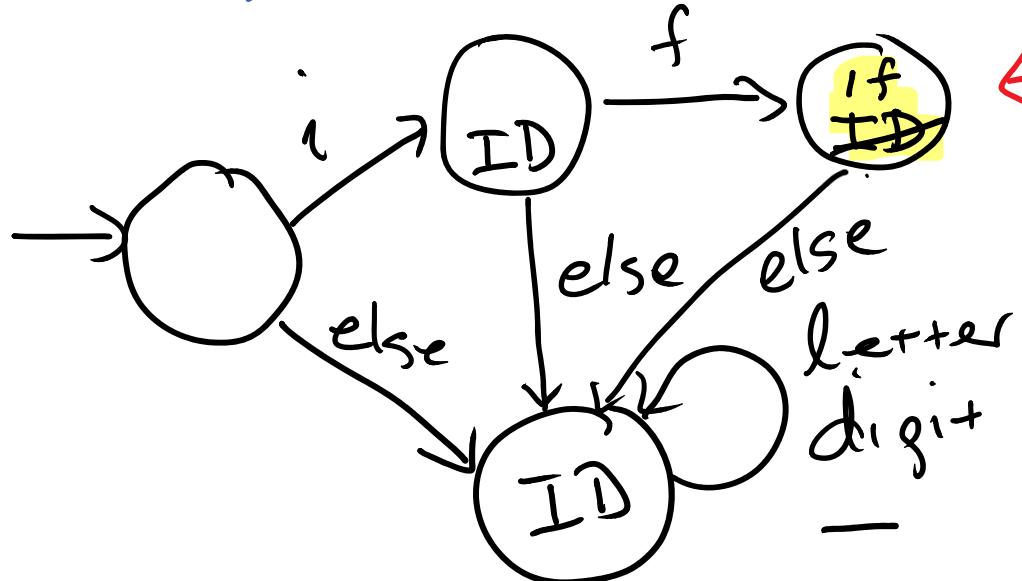
Rule Ordering

if

if ← special case first

ID

letter (letter | digit | -) ← default



← label for
first rule in
file wins.

Parsing

Context-Free Grammars

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

↑ token type (any number)

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

production

LHS \rightarrow RHS

(left-hand side) (right-hand side)

CFG for arithmetic expressions

$$\begin{array}{l} E \rightarrow E + E \\ E \rightarrow E - E \\ E \rightarrow E \times E \\ E \rightarrow -E \\ E \rightarrow (E) \\ E \rightarrow \text{num} \end{array}$$

LHS of production
must be a single symbol.

These symbols are called
"non-terminals"

This grammar only has one,
but a grammar can have
many.

CFG for arithmetic expressions

$$\begin{array}{l} E \rightarrow E + E \\ E \rightarrow E - E \\ E \rightarrow E \times E \\ E \rightarrow -E \\ E \rightarrow (E) \\ E \rightarrow \text{num} \end{array}$$

The sentence symbol is one of the non-terminals.

It is sort of like the start state of a DFA

In this case, it is E .

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

RHS is a string with
0 or more symbols.

Symbols that are not
non terminals are
terminals.

Here: +, -, ×, (,), num

Parse Trees and Languages

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$$1 + 2 \times 3$$

E \leftarrow sentence symbol

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

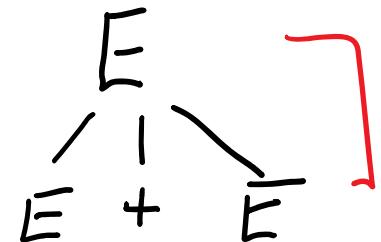
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$$1 + 2 \times 3$$



Parent = LHS
Children = RHS

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

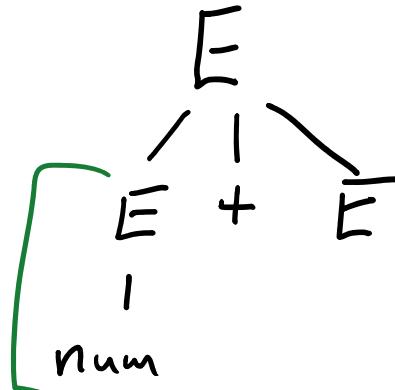
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

1 + 2 × 3



CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

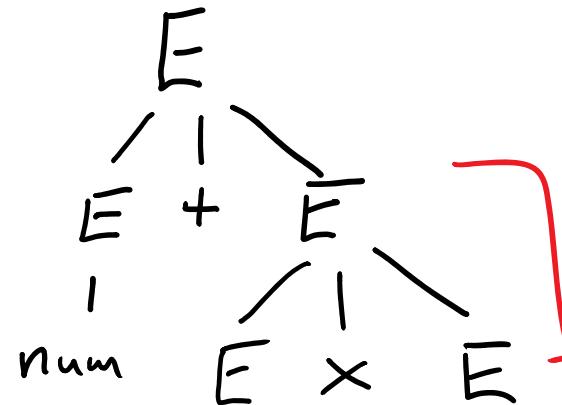
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$$1 + 2 \times 3$$



CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

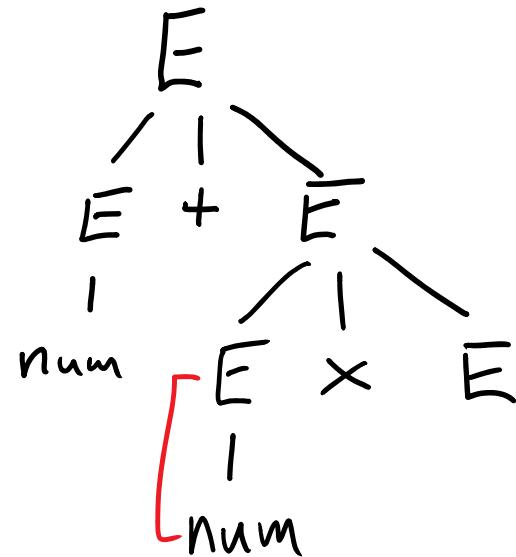
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

1 + 2 × 3



CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

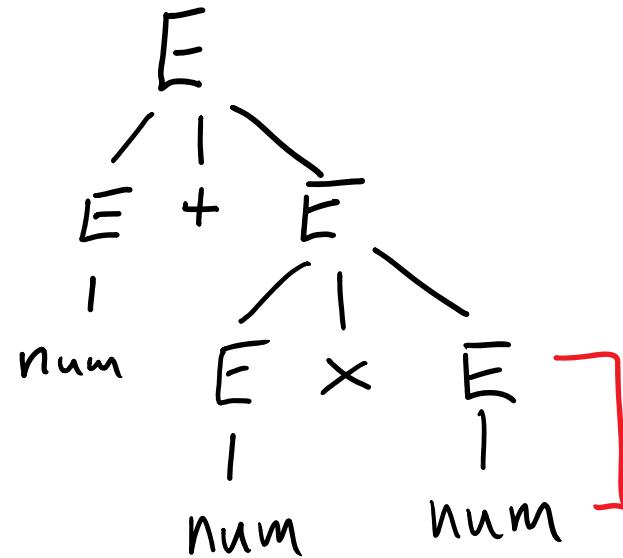
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$$1 + 2 \times 3$$



CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

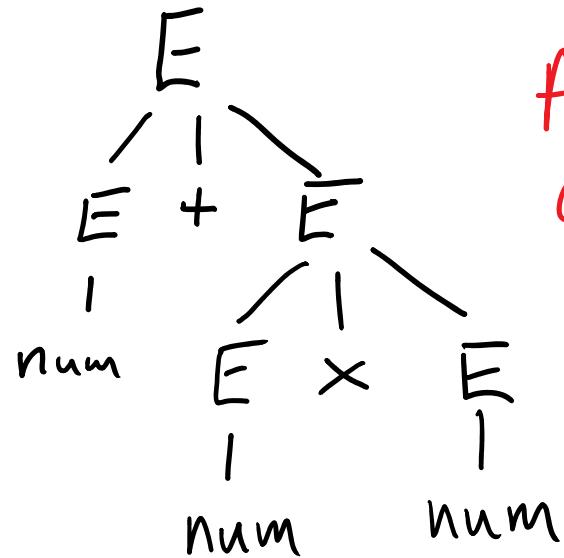
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$$1 + 2 \times 3$$



All leaves
are terminals

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

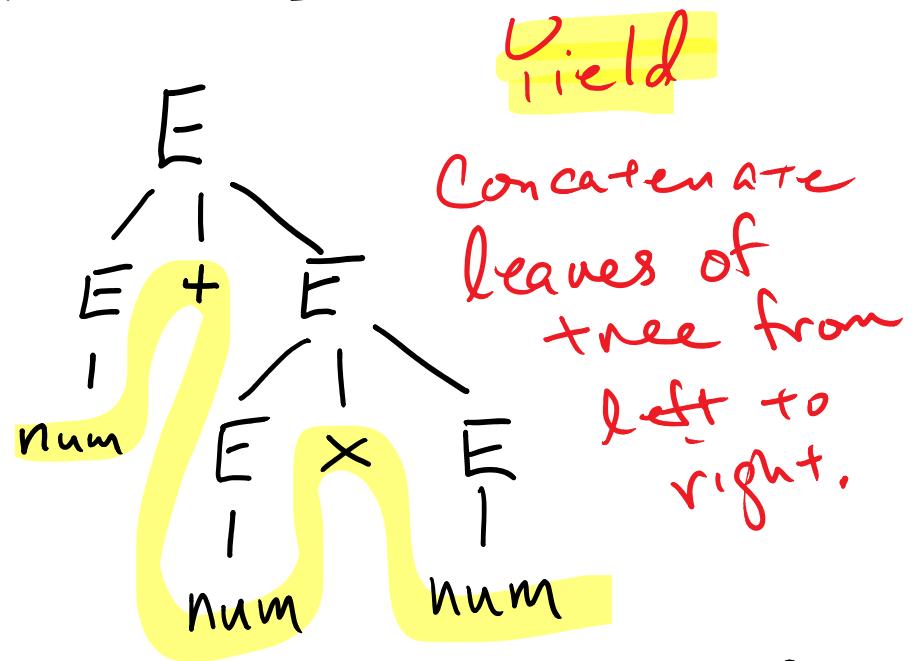
$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Parse tree

$1 + 2 \times 3$



num + num \times num is in language of grammar.

Def The language of a CFG, G , is the set of all terminal strings that are the yield of some parse tree for the CFG.

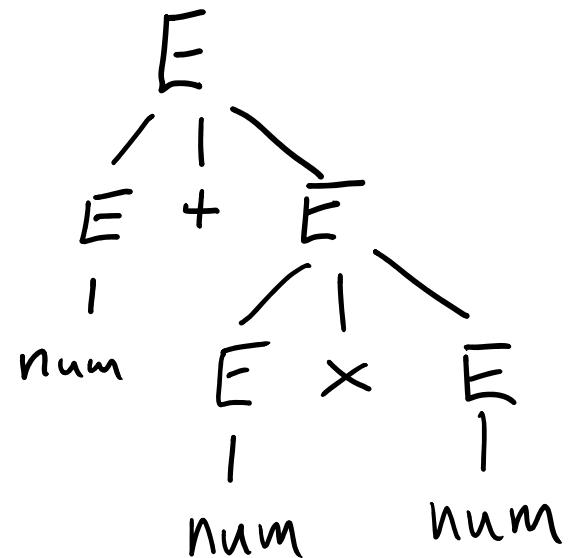
Parse tree for string X
= proof that x is in $L(G)$.

Def: The language of a CFG is a Context-Free language (CFL).

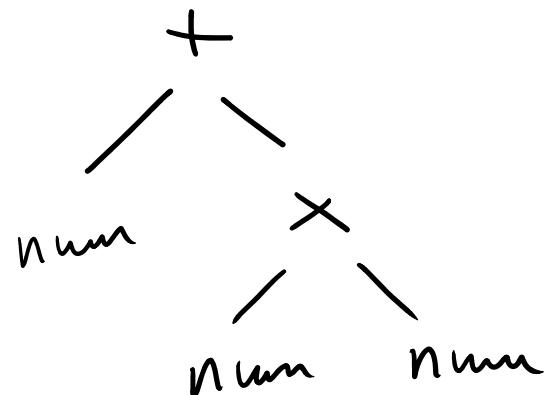
Parse Trees vs Abstract Syntax Trees

Parse tree

$1 + 2 \times 3$



AST

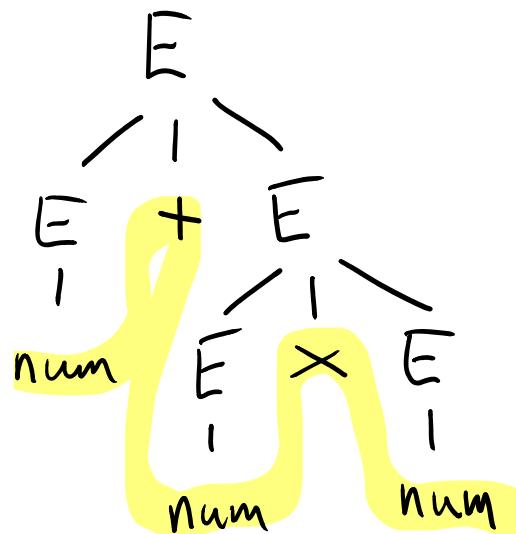


"Abstract" = details suppressed

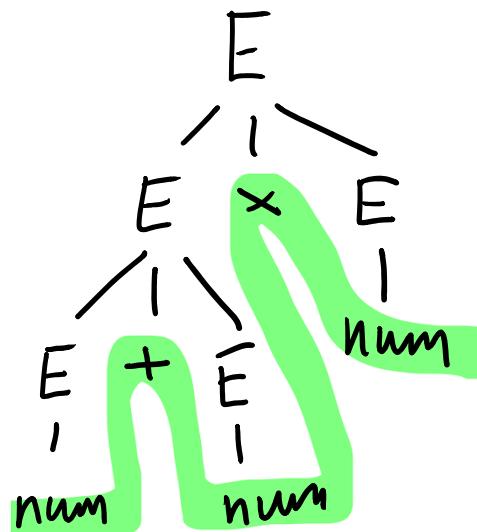
Ambiguity

Ambiguity

Def: A CFG is ambiguous if there is some terminal string with two parse trees.



$\text{num} + \text{num} * \text{num}$



Problems

1. Efficiency

Multiple parses

2. Interpretation

$$1 + \underline{2 \times 3} = 7$$

$$\underline{1 + 2} \times 3 = 9$$

Ambiguity in CFGs is a hard problem.

- Above examples are easy - come from failure to specify precedence.
- There are many ways to be ambiguous
- Some CFLs are inherently ambiguous
- Ambiguity of a CFG is undecidable!

Mathematical Definitions

A CFG is (V, Σ, R, S)

V - finite set of non-terminal symbols ("variables")

Σ - finite set of terminal symbols

$R \subseteq V \times (V \cup \Sigma)^*$ - finite set of productions

$S \in V$ - sentence symbol.

Notation

$a, b, c, \emptyset, 1$ — usually terminals

A, B, C, S — usually non terminals

x, y, z, w — usually strings in Σ^*

α, β, σ — usually strings in $(\Sigma \cup V)^*$

Derivation step

$$\alpha A \beta \Rightarrow \alpha \gamma \beta \text{ if } A \rightarrow \gamma$$

Derivation: Sequence of 0 or more steps

$$E \Rightarrow E + E \Rightarrow E + E \times E \Rightarrow E + \text{num} \times E$$

$$\Rightarrow \text{num} + \text{num} \times E \Rightarrow \text{num} + \text{num} \times \text{num}$$

$\alpha \xrightarrow{*} \beta$ (α "derives" β) if there
is a derivation $\alpha \Rightarrow \alpha_0 \dots \Rightarrow \beta$.

$\alpha \xrightarrow{*} \alpha$ always

$\alpha \xrightarrow{*} \beta$ (α "derives" β) if there
is a derivation $\alpha \Rightarrow \alpha_0 \dots \Rightarrow \beta$.

$\alpha \xrightarrow{*} \alpha$ always

Def The language of a CFG G ($L(G)$)

is $\{x \mid x \in \Sigma^* \wedge S \xrightarrow{*} x\}$

"sentences" terminal
string

derived from sentence
symbol

leftmost derivation

There are many derivations for one parse tree

$$E \Rightarrow E + E \Rightarrow \underline{\text{num}} + E \Rightarrow \dots$$

$$E \Rightarrow E + E \Rightarrow E + E \times E \Rightarrow \dots$$

Def: A **leftmost derivation** is a derivation in which the **leftmost nonterminal** is replaced at each step.

left most derivation

$$\bar{E} \Rightarrow E + \bar{E} \Rightarrow \text{num} + \bar{E} \Rightarrow \text{num} + \bar{E} \times \bar{E}$$

$$\Rightarrow \text{num} + \text{num} \times \bar{E} \Rightarrow \text{num} + \text{num} \times \text{num}$$

A CFG is ambiguous iff there exists
a string $x \in \Sigma^*$ that has

- 2 leftmost derivations,
equivalently,
- 2 rightmost derivations
equivalently
- 2 parse trees.

Two Leftmost Derivations for Same String

$$\bar{E} \Rightarrow \bar{E} + \bar{E} \Rightarrow \text{num} + \bar{E} \Rightarrow \text{num} + \bar{E} \times \bar{E}$$

$$\Rightarrow \text{num} + \text{num} \times \bar{E} \Rightarrow \text{num} + \text{num} \times \text{num}$$

$$\bar{E} \Rightarrow \bar{E} \times \bar{E} \Rightarrow \bar{E} + E \times E \Rightarrow \text{num} + \bar{E} \times \bar{E}$$

$$\Rightarrow \text{num} + \text{num} \times \bar{E} \Rightarrow \text{num} + \text{num} \times \text{num}$$

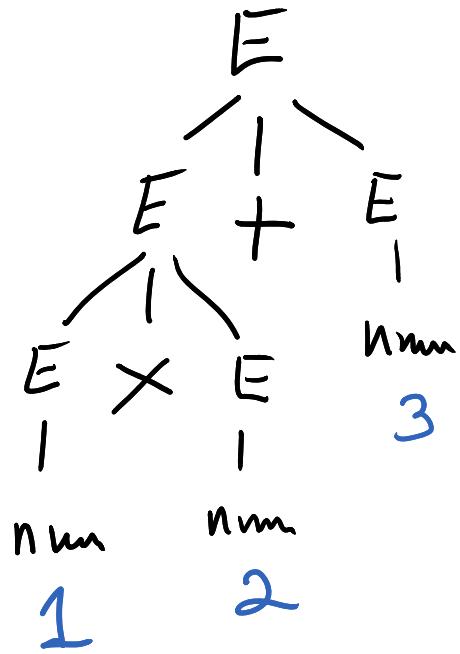
Announcements

- Correction: PA1 is due 4/16.
- WA1 assigned today, due 4/16 (1 week).
 - “Automatic extension” without penalty until Monday 4/20
 - Other assign/due dates will remain the same.

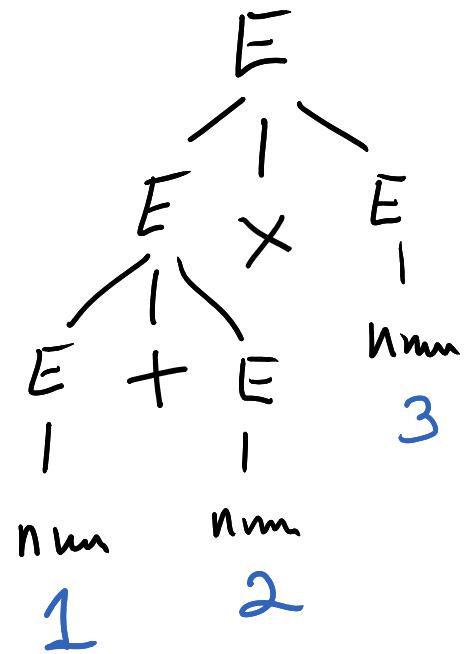
Precedence and Association

$$1 \times 2 + 3$$

Which operator wins?

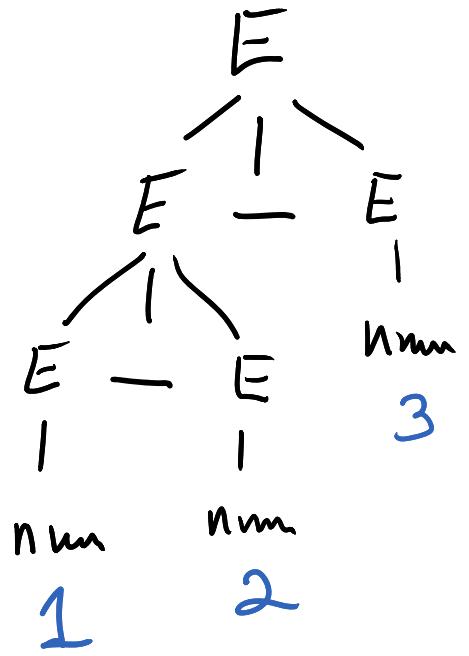


$(1 \times 2) + 3$
 $\uparrow \times$ "stickier"



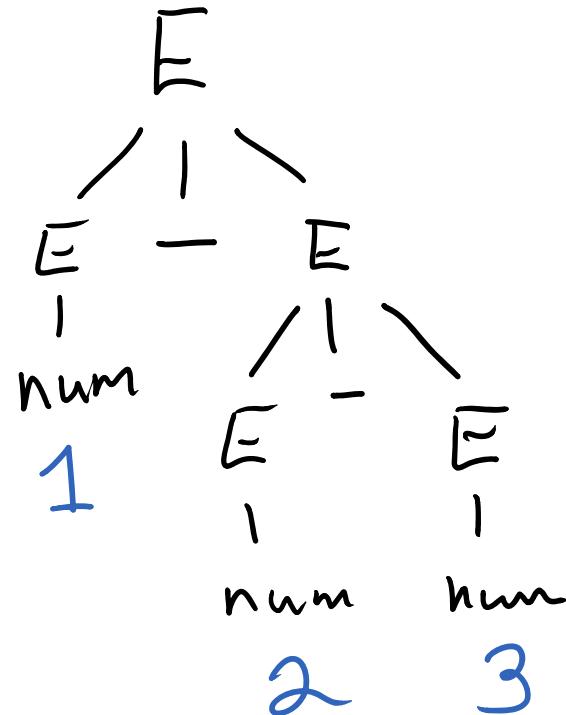
$(1+2) \times 3$
 $\uparrow +$ stickier

"higher precedence" means "stickier"



$(1 - 2) - 3$

left associative



$1 - (2 - 3)$

right associative

CFG for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

This CFG has ambiguities
for relative precedence of
operators and
left-right associativity

Abbreviated Notation.

$$E \rightarrow E + \bar{E}$$

$$E \rightarrow E - E$$

$$E \rightarrow E \times E$$

$$E \rightarrow -E$$

$$E \rightarrow (E)$$

$$E \rightarrow \text{num}$$

Group productions with
same LHS together

$$\begin{aligned} E \rightarrow & E+E \mid E-E \mid E \times E \\ & -E \mid (E) \mid \text{num} \end{aligned}$$

Changing a CFG to fix precedence problems

Not always possible

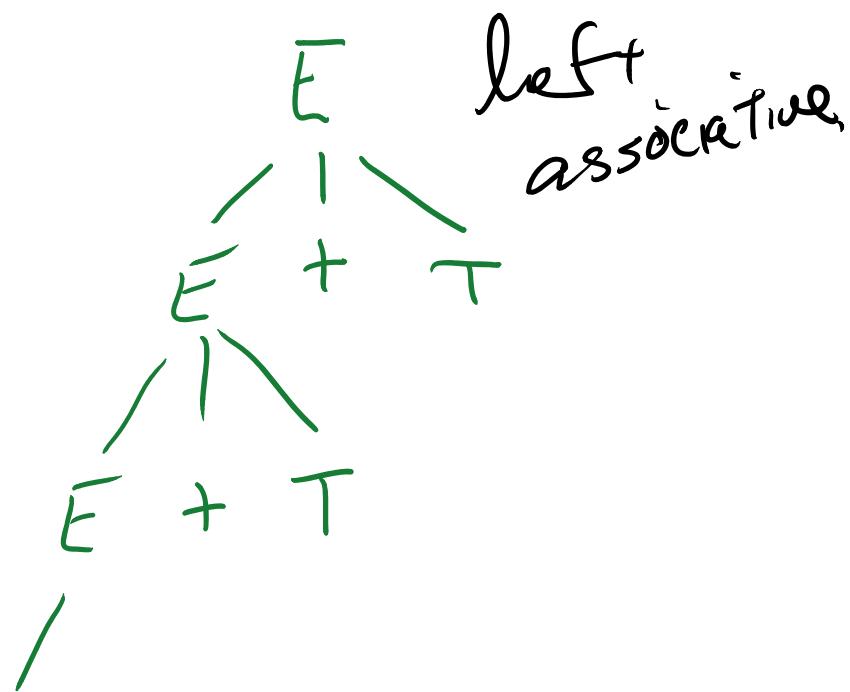
- Rank operators by increasing precedence
- New non-terminal for each precedence level
- Sentence symbol goes to lowest precedence operator
- Productions for generating more of same precedence operators, going to higher precedence
- left/right recursion for left/right associativity.

$,$ — same precedence (lowest)

\times next higher precedence

unary — highest precedence.

$$E \rightarrow E + T \mid E - T \mid T$$



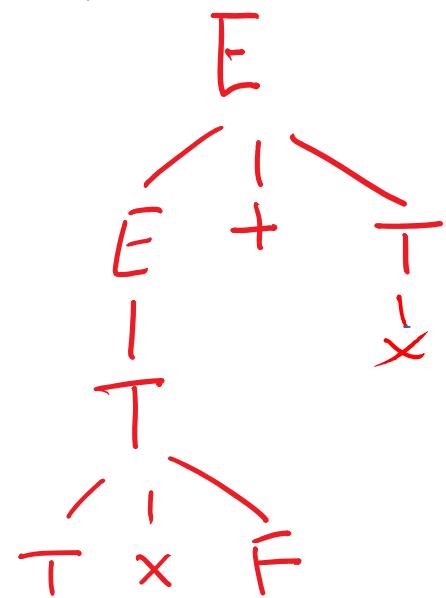
$+$, $-$ same precedence (lowest)

\times next higher precedence

unary - highest precedence.

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T \times F \mid F$$



$+$, $-$ same precedence (lowest)

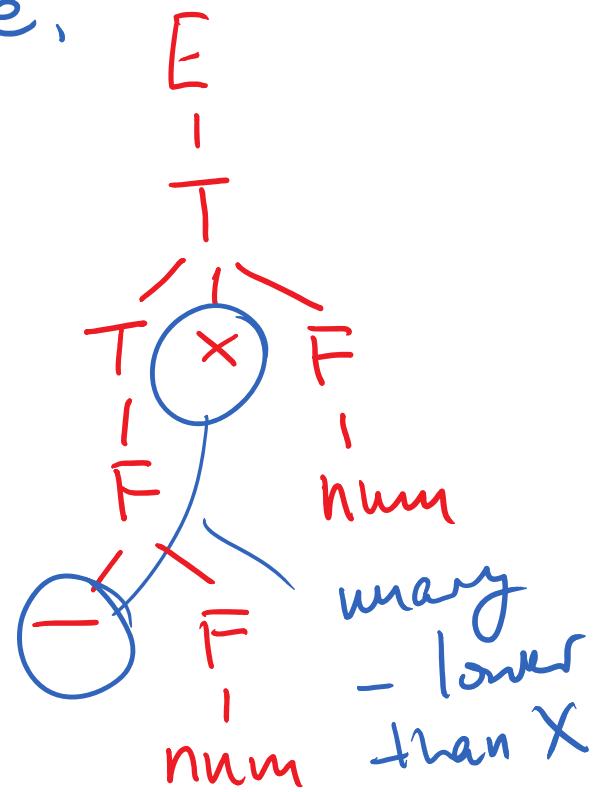
\times next higher precedence

unary - highest precedence.

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow -F \mid (E) \mid \text{num}$$

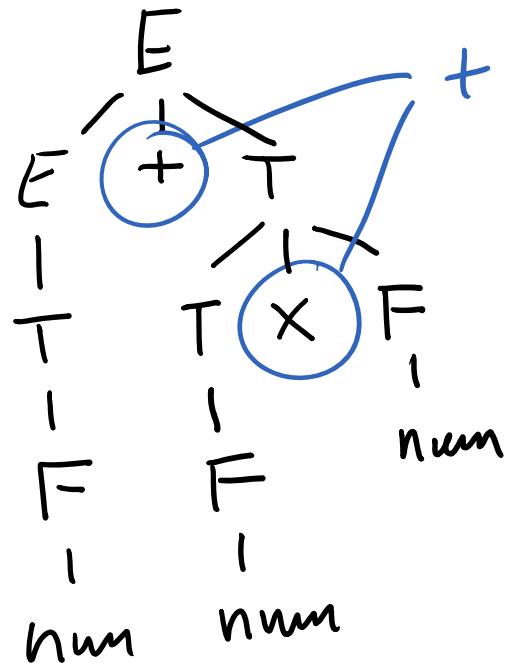


$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow -F \mid (E) \mid \text{num}$$

num + num × num



+ always above
×
(unless-
parenthesized)

Extended BNF

EBNF

BNF - "Backus-Naur Form" = CFG

Idea: Allow regular expressions on RHS
of productions.

Widely used (esp. in programming language
definitions)

Parser generators usually require "plain CFG"
Need to convert EBNF \rightarrow CFG

$A \rightarrow B^* d$ \leftarrow meta-symbol, not a terminal
 \emptyset or more B 's on RHS

α^* abbreviates

$C \rightarrow C\alpha \mid \epsilon$

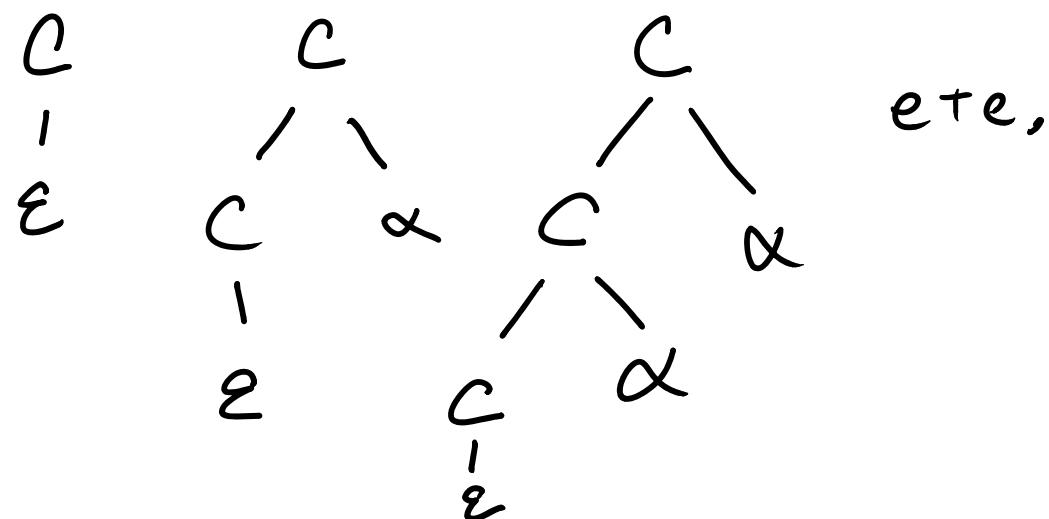


new non-terminal, not already in V.

$A \rightarrow B^* d$ ← meta-symbol, not a terminal
∅ or more B's on RHS

α^* abbreviates

$$C \rightarrow C\alpha \mid \epsilon$$



$$A \rightarrow B^* d$$
$$A \rightarrow c d$$
$$c \rightarrow C B \mid \varepsilon$$

α^*

$$C \rightarrow C\alpha | \epsilon$$



left-recursive

$$C \rightarrow \alpha C | \epsilon$$



right-recursive

May choose one or
the other because

- Necessary for
parsing method
- Grammar
properties
- Tree shape.

α^+

$C \rightarrow C\alpha/\alpha$

$C \rightarrow \alpha C/\alpha$

α^+

$C \rightarrow C\alpha/\alpha$ $C \rightarrow \alpha C/\alpha$

$\alpha?$ ($[\alpha]$ in COOL manual)

α^+

$C \rightarrow C\alpha|\alpha$ $C \rightarrow \alpha C|\alpha$

$\alpha?$ ($[\alpha]$ in COOL manual)

$C \rightarrow \alpha|\varepsilon$