

Parsing:

10/3/18

Given a context-free grammar G and a string w , find a sequence of productions of G that produces w , or find out that w can't be produced.

or derivation tree
or sequence of
leftmost productions

E.g.: Given

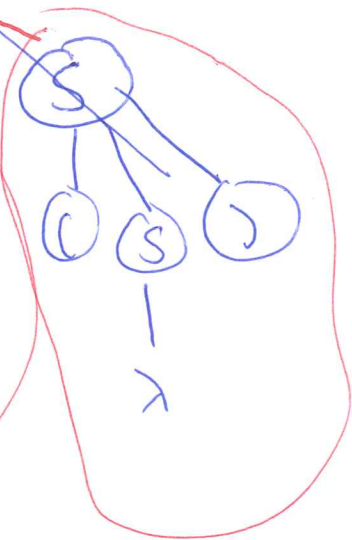
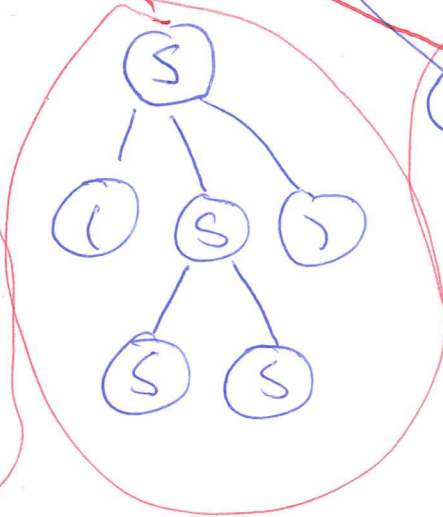
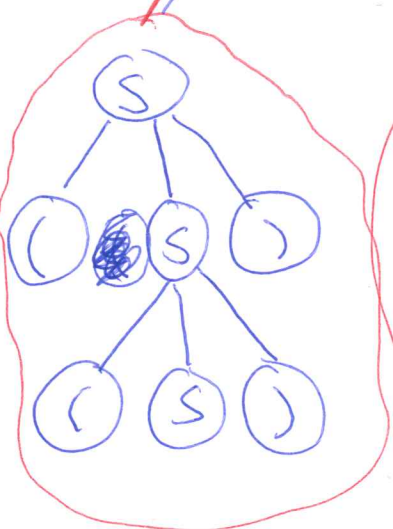
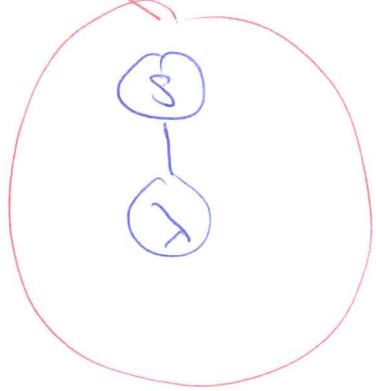
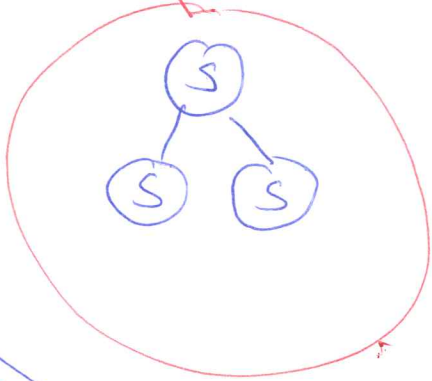
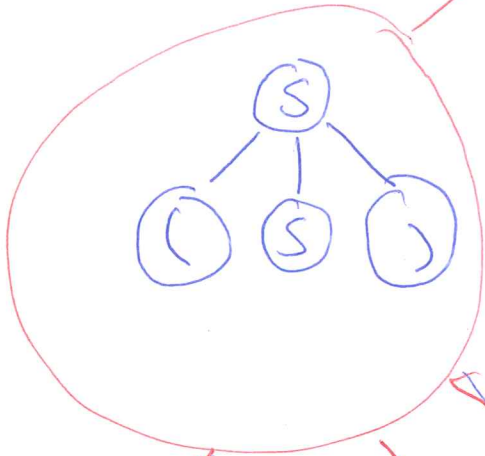
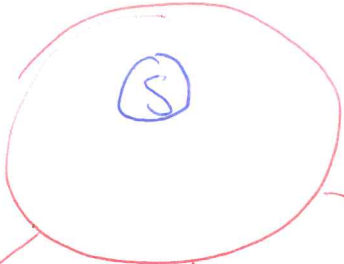
$$S \rightarrow (S) \mid SS \mid \lambda$$

find out if

$((()())(())())(())$ is produced by this grammar.

Stupid method: Try everything. Do a breadth-first search on the tree of possible derivation trees.

Problems: ① This takes ~~for~~ a long time.
② If the answer is that ~~the~~ w can't be produced, we never find out.



How can we solve ②?

Idea: Modify the grammar so that ~~it~~ it has no productions of the form

$$A \rightarrow \lambda \quad (\text{var} \rightarrow \lambda)$$

or

$$A \rightarrow B \quad (\text{var} \rightarrow \text{single var only})$$

If I don't have productions of these forms, I must "make progress" at every production - I either

a) increase # of letters in the string
or b) increase # of variables, in this case I can lose 1 var.

i.e. to generate a string w/ 20 letters, the slowest possible method is to add 19 vars in 19 steps and convert vars 1 by 1 to letters in 20 steps

total at most 39 steps.

This solves ② since we know we don't need to search further down than twice the length of the string. (i.e. once we have searched that much of the tree, we can declare failure)

Making ① a little better: pruning the search tree:

- a) ~~Always~~ Always do leftmost derivations
- b) Once the leftmost letters (before any variable) don't match the leftmost letters of your string, you can stop looking down that branch.

Some ~~the~~ restrictions on our grammar can make them easier to parse

- 1) Regular grammars are easy to parse
- 2) S-grammars:

~~a) "linear" every deriv.~~
a) Every derivation looks like
$$V \rightarrow l(\text{stuff})$$

first thing must be a letter. (not var)

b) Given any variable V and letter l , there is at most one derivation looking like

$$V \rightarrow l(\text{stuff})$$

i.e. if I have

$$A \rightarrow xyBC \mid xAyB,$$

this is not an s-grammar.

These are easy to parse b/c you can look at the first letter that is still in variable form and find out what derivation you use to get it (or find out none is possible)

e.g.: If my string is $w = xy \dots$, then the first prod I use must be the production

$$S \rightarrow x(\text{stuff})$$

Suppose this was

$$S \rightarrow xABAy.$$

Then the next prod would have to be

$A \rightarrow y$ (stuff) (and there is at most one of these)

Another way to think of this:

pruning takes off all but one branch.

Ambiguity

A grammar G is ambiguous if there is a string w produced by G in more than one way - i.e. w has two different derivation trees.

A language L is inherently ambiguous if every grammar for L is ambiguous (there is no unambiguous grammar for L)

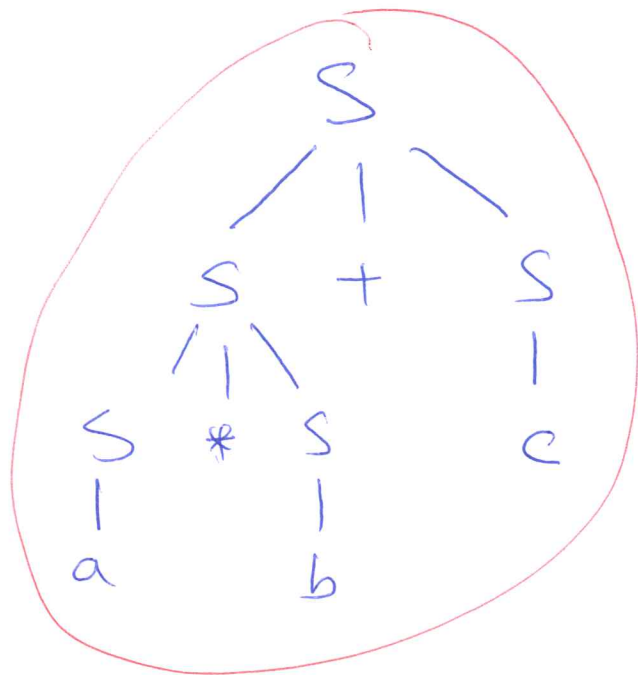
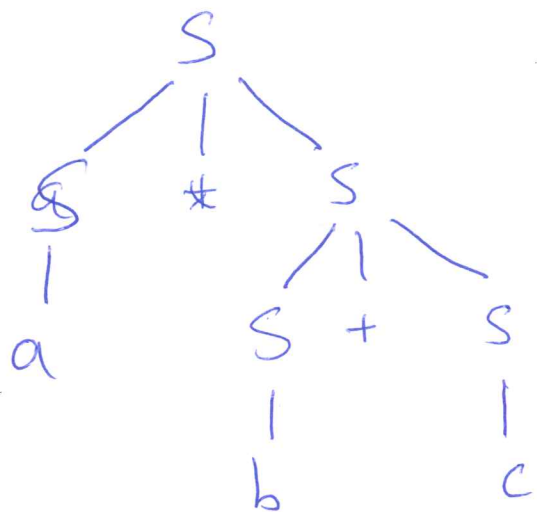
Example: Grammar for restricted arith. express

$$\Sigma = \{ (, +, *,), a, b, c, d \}$$

$$S \rightarrow (S) | \lambda | a | b | c | d | S + S | S * S$$

This is ambiguous since

$a * b + c$ has derivations:



We can fix this
at the parsing level
- specify which deriv.
to try first.

better b/c it
reflects our order
of operations.

We can also try to make an unambiguous grammar:

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

$$T \rightarrow F * T \mid F$$

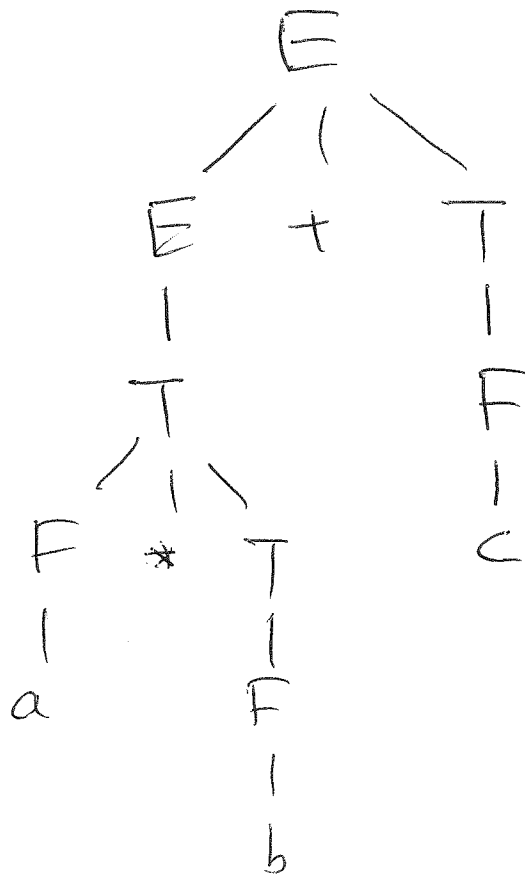
$$F \rightarrow a \mid b \mid c \mid d \mid (E)$$

E = expression ↙ start

T = term

F = factor

To parse $a * b + c$, we have to have



Something that is inherently ambiguous.

$$L = \{ a^x b^y c^z \mid x=y \text{ or } y=z \}$$

$$S \rightarrow \cancel{A} Z \mid X B$$

$$Z \rightarrow c Z \mid \lambda$$

$$X \rightarrow a X \mid \lambda$$

$$A \rightarrow a A b \mid \lambda$$

$$B \rightarrow b B c \mid \lambda$$

aabbcc