


Parsing

Grammars

Parse trees

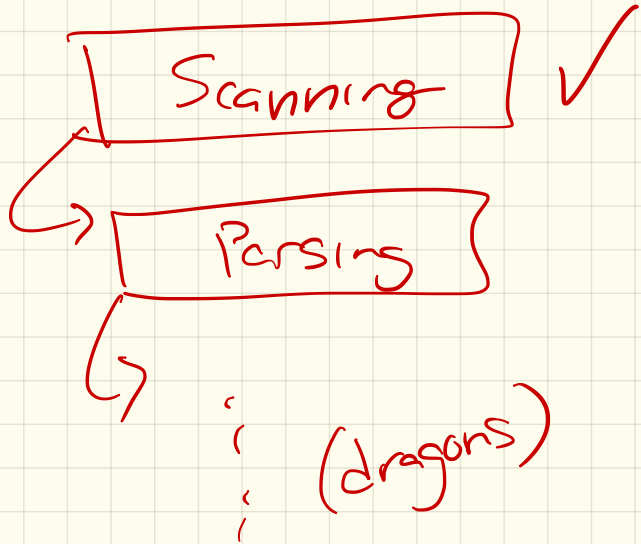
left vs right



Recap:

- HW due today
- Next HW - flex Scanner (due next Thursday)
- Git repo guide

Compilation



The limits of regular expressions:

Certain languages can never be recognized by regular expressions.

Simple example: $0^n 1^n$

$L = \{\epsilon, 01, 0011, 000111, \dots\}$

Only option: *

$0^* 1^*$ → this accepts 011

More real world: math equations

$((((x+7)*2)+3)-1)$

Why? At heart,
 $(\dots)^n \dots$

So, in our next level down, we need something stronger

Parsing:

- Given string of input tokens, a parser must determine if the tokens generate a valid program

The basis of these are context free grammars (CFGs):

- terminals: tokens (via scanner)
 $ids, +, (,), -$
- nonterminals (one a start symbol)
(usually capital letter)
- production rules

$$E \rightarrow (E)$$
$$\rightarrow \underline{id} + \underline{id}$$

Example: $0^n 1^n$

terminals: $0 + 1$ ($+ \epsilon$)

L : $S \rightarrow 0S1$ (\leftarrow)

$S \rightarrow \epsilon$ (\leftarrow)

(nonterminals: S
also start)

Show $000111 \in L$

$S \xrightarrow{(1^{st})} 0S1$

$\rightarrow 00S11$ (1^{st})

$\rightarrow 000S111$ (2^{nd})

$\rightarrow 000111$ (2^{nd})

term: id, -, +, /, ↑

Ex: (bigger)

non terms: E, A

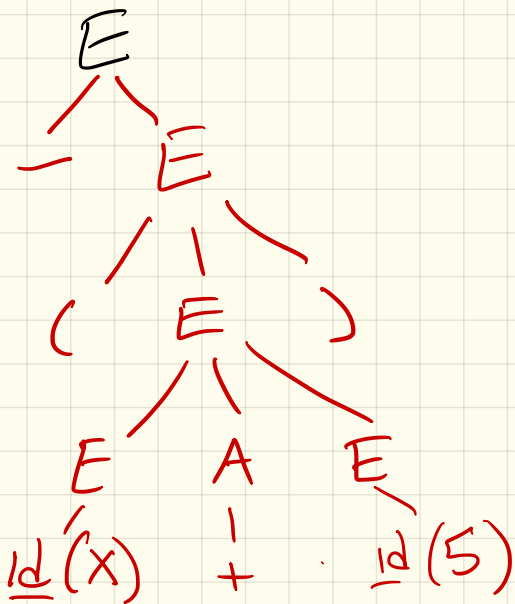
E	→	E A E	1
E	→	(E)	2
E	→	-E	3
E	→	id	4
A	→	+	5
A	→	-	6
A	→	*	7
A	→	/	8
A	→	↑	9

Derivation: The process by which a grammar parses & defines a language.

Ex: Show $\ominus(x+5)$ is accepted by the above grammar:

$$\begin{aligned}
 E &\Rightarrow -E \xRightarrow{(2)} -(E) \\
 &\xRightarrow{(1)} -(E A E) \xRightarrow{(5)} -(E + E) \\
 &\xRightarrow{(4)} -(\underline{\text{id}} + E) \xRightarrow{(4)} -(\underline{\text{id}} + \underline{\text{id}}) \xleftarrow{5}
 \end{aligned}$$

Parse tree: A graphical representation of this derivation:



Each parent/child^(ren) shows one step of the derivation

- leaves are terminals
- root is start non-terminal

Leftmost vs rightmost

Leftmost derivation: one where the leftmost nonterminal is replaced in each step

Ex: $-(id + id)$

$$E \Rightarrow -E \Rightarrow -(E)$$

$$\Rightarrow -(EAE) \Rightarrow -(id \underline{A}E)$$

$$\Rightarrow -(id + \underline{E}) \Rightarrow -(id + id)$$

- Rightmost: always the one on the right

(Often we'll fix one way, since each tree has a unique left & right most derivation.)

Note: Not necessarily unique!

Ex: $5 + 2 * 6$

① $E \Rightarrow \underline{E} A E$
 $\Rightarrow \underline{E A E} A E$
 $\Rightarrow \text{id}(5) \underline{A} E A E$
 $\Rightarrow \text{id}(5) + \underline{E} A E$
 $\Rightarrow \text{id}(5) + \text{id}(2) * \text{id}(6)$

② $E \Rightarrow E A \underline{E}$

Ambiguity:

- Any grammar that produces more than one parse tree for some sentence is ambiguous.

This can be very undesirable!

We'll spend time trying to rule this possibility in our grammars. ☺

Note: Any regular expression can also be written as a grammar (CFG) nonterm \rightarrow char

Ex: $(a|b)^*abb$

Grammar:

$S \rightarrow Aabb$

$A \rightarrow \epsilon$

$\rightarrow aA$

$\rightarrow bA$

(However, the reverse is not true!)

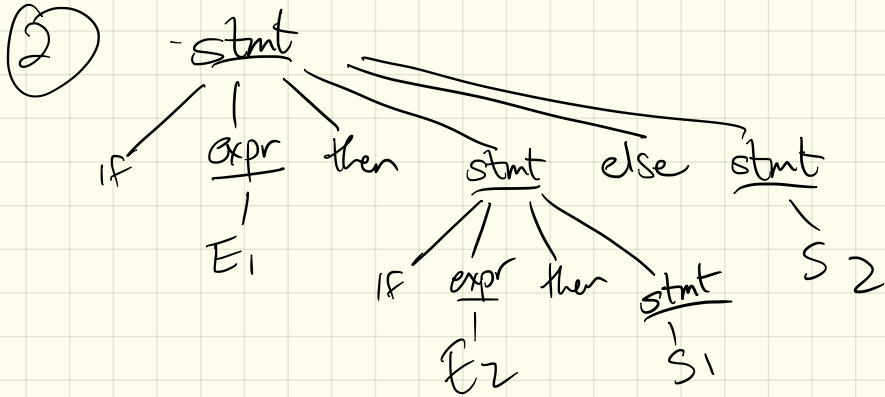
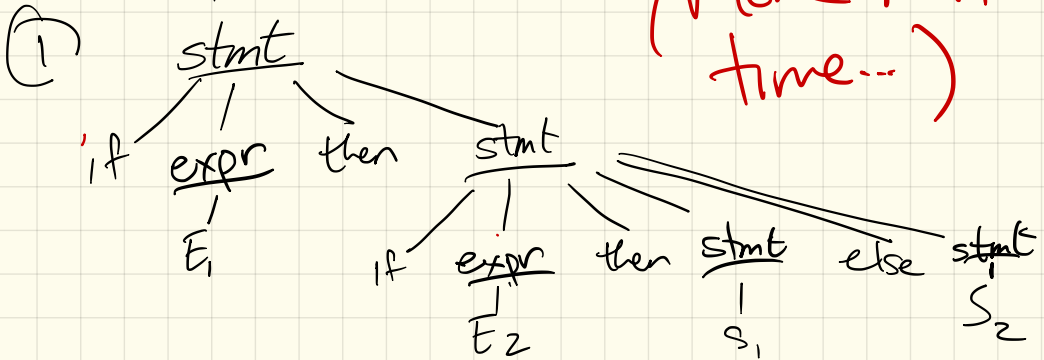
Still do regular expressions:
Simpler + faster

A more complex example: **if statements**
stmt \rightarrow if expr then stmt
 | if expr then stmt else stmt
 | other

Then: if E_1 then if E_2 then S_1 else S_2

2 parse trees:

(More next time...)



General rule:

Match each else w/ closest unmatched then

How?

- Rewrite so any statement between an "else" + a "then" must be matched (so no if-then w/ no else)

Grammar:

stmt \rightarrow matched_stmt
| unmatched_stmt

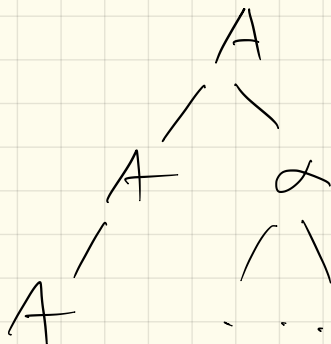
matched_stmt \rightarrow if expr then matched_stmt else matched_stmt
| other

unmatched_stmt \rightarrow if expr then stmt
| if expr then matched_stmt
else unmatched_stmt

Dfn: A grammar is left-recursive
if it has a non-terminal A
with some rule

$$A \rightarrow A \alpha$$

These are bad for parsers:



When scanning tokens &
trying to build a tree,
not sure when to stop!

$$\begin{aligned}\underline{\text{Ex:}} \quad E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id}\end{aligned}$$

$$\underline{\text{Parse:}} \quad x + y * 10$$

However, we do have left recursion!

To eliminate:

$$A \rightarrow A\alpha \mid \beta$$

$$\hookrightarrow A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

On

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

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

$$F \rightarrow (E) \mid id$$

Back to the practical:

- Any CFG can be parsed
 - ↳ Chomsky Normal Form
 - CYK algorithm
 - Run time:

This is too slow!

Most modern parsers look for certain restricted families of CFGs.

Result:

Top down parsing

Called predictive parsing.

Works well on LL(1) grammars.

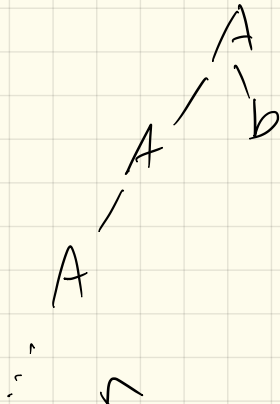
Ex: $S \rightarrow cAd$
 $A \rightarrow ab/a$

Parse cad:

Rule: string w/ S,
apply rules until
one matches the
next input
(back track if there
is a mistake)

Note: Left recursion is
very bad on these!

$$A \rightarrow Ab$$



∴ never matches an input or hits a conflict

So never forced to backtrack.

How predictive parsing works:

- the input string w is in an input buffer.

- Construct a predictive parsing table for G .

- if you can match a terminal, do it
(+ move to next character)

- otherwise, look in table for rule to get transition that will eventually match

Hard part (next time):
the table