

Parsing (cont)



Today:

- Flex HW due Thursday

- Via git

(only need personal
git repo)

- Started parsing -

Sec. 2.3

supplemental text:

dragon book

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: for, +, { → lowercase
- nonterminals (one a start S symbol) typically uppercase or underlined
- production rules
↳ tell transition

Notation: ↙ start

expr → expr op expr
| (expr)
| id (variable)

op → + | - | * | /

Ex: ^{capital, so non-terminal}

$$\begin{aligned} E &\rightarrow E A E \\ &\rightarrow (E) \\ &\rightarrow -E \\ &\rightarrow \text{id} \\ A &\rightarrow + \\ &\rightarrow - \\ &\rightarrow * \\ &\rightarrow / \\ &\rightarrow \uparrow \end{aligned}$$

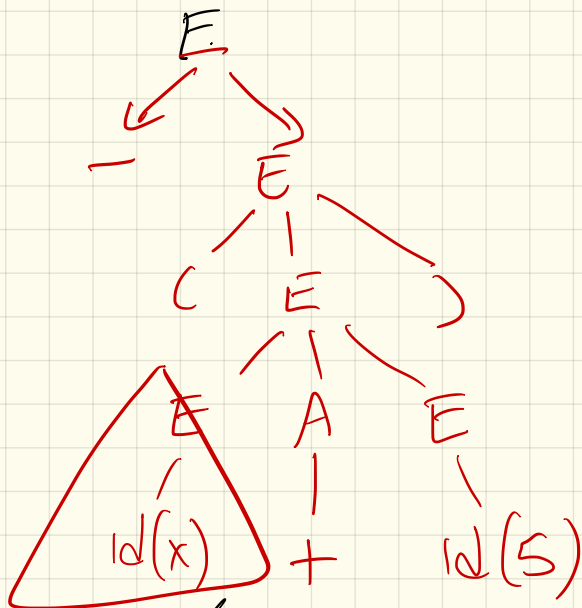
^{terminals}

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

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

$$\begin{aligned} E &\Rightarrow -E \Rightarrow -(E) \\ &\Rightarrow -(E A E) \Rightarrow -(id(x) A E) \\ &\Rightarrow -(id(x) + E) \\ &\Rightarrow -(id(x) + id(5)) \end{aligned}$$

Parse tree: A graphical representation of this derivation:



Each parent/child shows one step of the derivation

- leaves are terminals
- root is start non-terminal

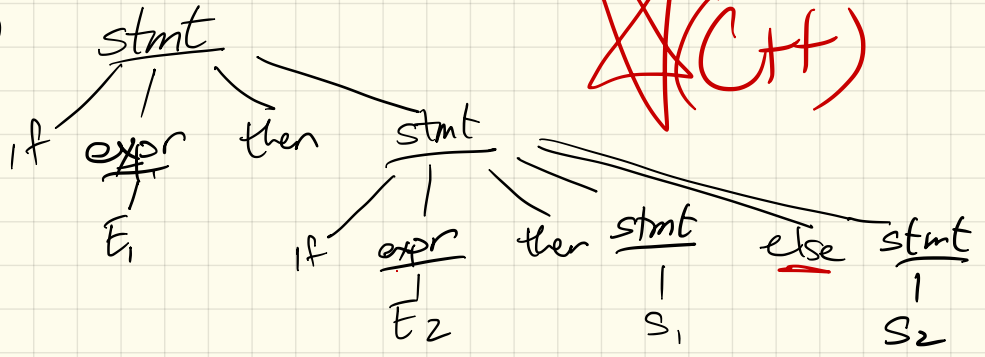
Other things

- Left most vs rightmost
- Ambiguity

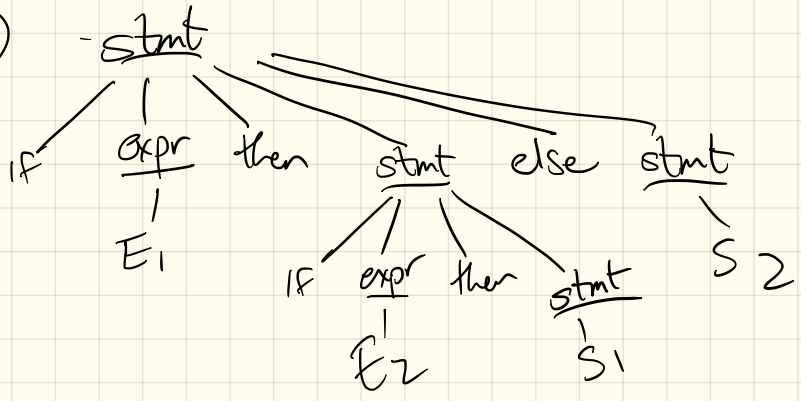
Ex: if E_1 then if E_2 then S_1 else S_2

2 parse trees:

①



②



General rule:

Match each else w/ closest 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 (other exprs)
| stmt

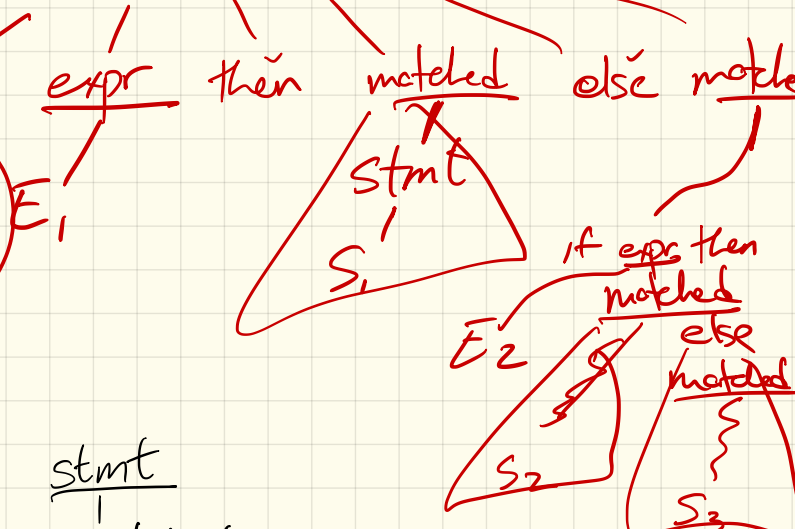
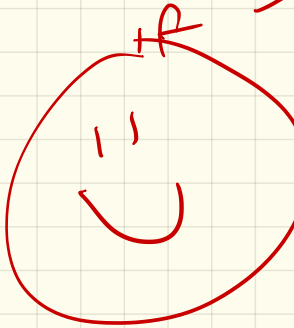
unmatched_stmt \rightarrow if expr then stmt
| if expr then matched_stmt
else unmatched_stmt

Example:

~~if~~ $\underline{if\ E_1}$ then $\underline{S_1}$ else $\underline{if\ E_2}$ $\underline{S_2}$ else $\underline{S_3}$)

①

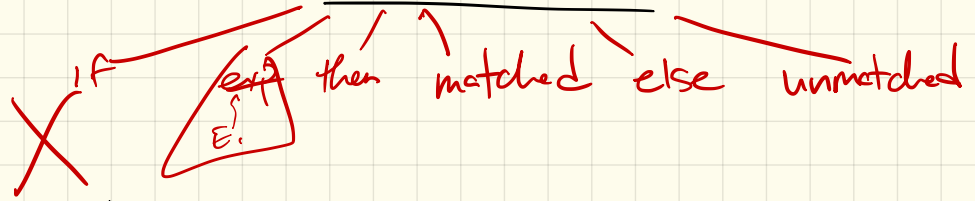
stmt
|
matched_stmt



②

error

stmt
|
unmatched_stmt



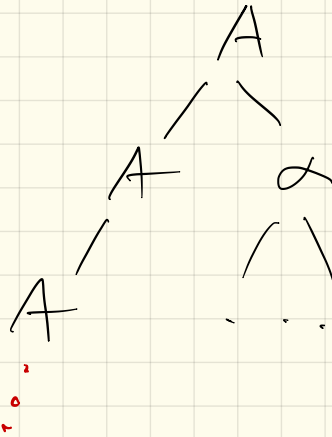
unmatched_stmt \rightarrow if expr then stmt
| if expr then matched_stmt
else unmatched_stmt

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

$$A \rightarrow A \alpha$$

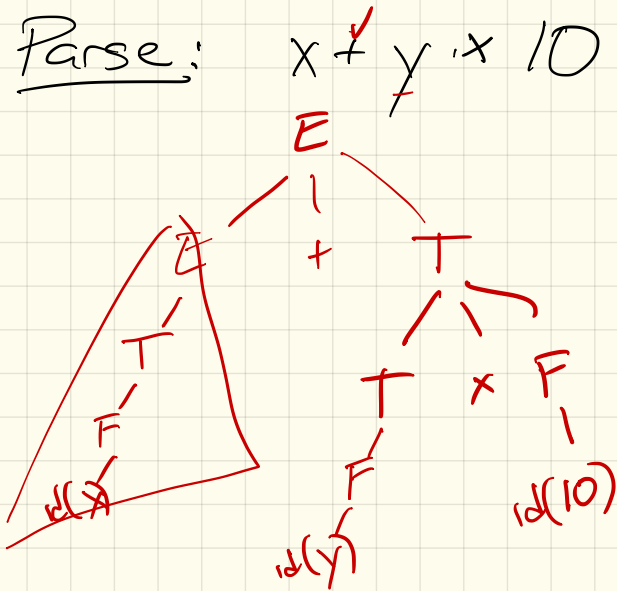
$$A \rightarrow \underline{B}A \text{ (not left recursion!)}$$

These are bad for parsers:



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

Ex: $E \rightarrow E + T \mid T$
 $T \rightarrow T \times F \mid F$
 $F \rightarrow (E) \mid \text{id}$



This deals nicely w/ precedence.
 However, we do have left recursion!

To eliminate:

$$\underline{A} \rightarrow \underline{A} \alpha \mid \beta$$

$$\left[\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array} \right]$$

On

$$\begin{array}{l} \textcircled{1} E \rightarrow E + T \mid T \\ \textcircled{2} T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid \text{id} \end{array}$$

$$\textcircled{1} \left[\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +TE' \mid \epsilon \end{array} \right]$$

$$\textcircled{2} \left[\begin{array}{l} T \rightarrow FT' \\ T' \rightarrow *FT' \mid \epsilon \\ F \rightarrow (E) \mid \text{id} \end{array} \right]$$

Back to the practical:

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

Run time: $O(n^3)$

(dynamic programming)

This is too slow!

Most modern parsers look
for certain restricted
families of CFGs.

Result: $O(n)$

LL or LR

Top down parsing

Called predictive parsing.

Works well on LL(1) grammars.

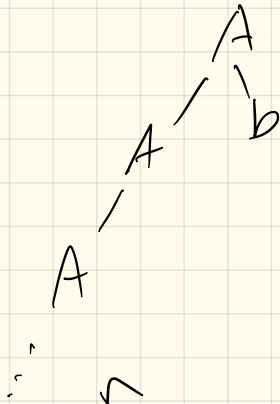
Ex: $S \rightarrow cAd$
 $A \rightarrow ab/a$

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:

• build the table

(need to decide a transition if at a nonterminal based on the next input terminal)

FIRST & Follow Sets:

FIRST (α) \leftarrow any string of non-terminals & terminals

$\hat{=}$ set of possible first terminals in any derivation of α by the grammar

So:

1) if x is a terminal,

$$\text{FIRST}(x) =$$

2) if $X \rightarrow \epsilon$ is a production,
add ϵ to $\text{FIRST}(x)$

3) If X is a nonterminal:

If $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production:

add a if a is in $\text{FIRST}(Y_i)$ and
 ϵ is in $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$

add ϵ if ϵ is in $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_k)$

Ex: $E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$

$$T \rightarrow FT'$$
$$T' \rightarrow *FT' \mid \epsilon$$

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

$$\text{FIRST}(E) =$$

$$\text{FIRST}(E')$$

$$\text{FIRST}(T)$$

$$\text{FIRST}(T')$$

$$\text{FIRST}(F)$$