# Parsing (cont)

# Today:

- HW due Friday
- Today: - after class
         or 1:30-2
- Next HW - due next
           Friday
     over flex
     submit via git

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

Notation: ↙ start

$$expr \rightarrow expr\ op\ expr$$
$$| \ (expr)$$
$$| \ id\ (variable)$$
$$op \rightarrow +|-|*|/$$

Ex :  *capitel, so* non-terminal

$$E \longrightarrow \boxed{E \quad A \quad E}$$
$$\longrightarrow (E)$$
$$\longrightarrow -E$$
$$\longrightarrow id$$

^c *terminals*

$$A \longrightarrow +$$
$$\longrightarrow -$$
$$\longrightarrow *$$
$$\longrightarrow /$$
$$\longrightarrow \uparrow$$
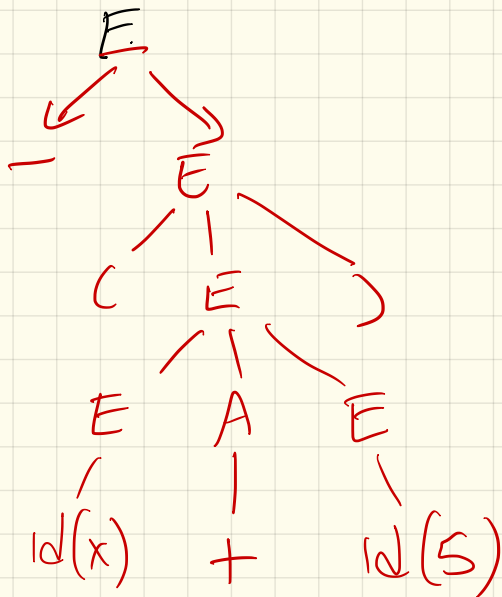
Derivation: The process by which a grammar parses & defines a language.

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

$$E \Rightarrow -E \Rightarrow -(E) :$$
$$\Rightarrow -(E \, A \, E) \Rightarrow -(id(x) \, A \, E)$$
$$\Rightarrow -(id(x) + E)$$
$$\Rightarrow -(id(x) + id(5))$$

Parse tree: A graphical representation
             of this derivation:

$$
\begin{array}{c}
E \\
\diagup \\
- \quad E \\
\diagup \mid \diagdown \\
C \quad E \quad ) \\
\diagup \mid \diagdown \\
E \quad A \quad E \\
\mid \quad \mid \quad \mid \\
Id(x) \quad + \quad Id(5)
\end{array}
$$

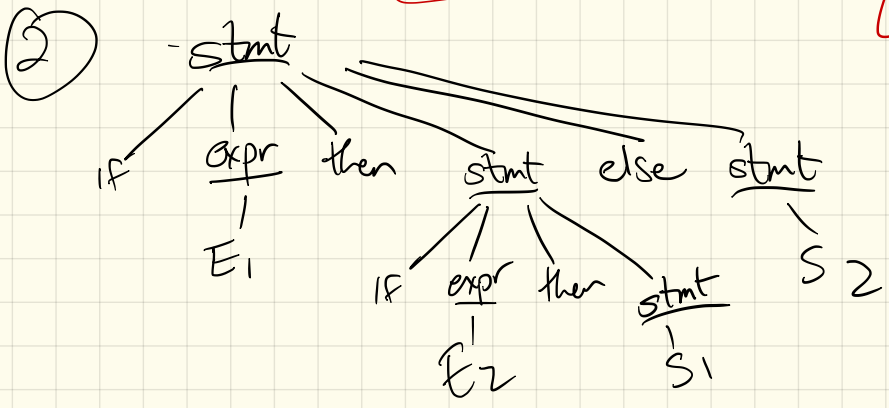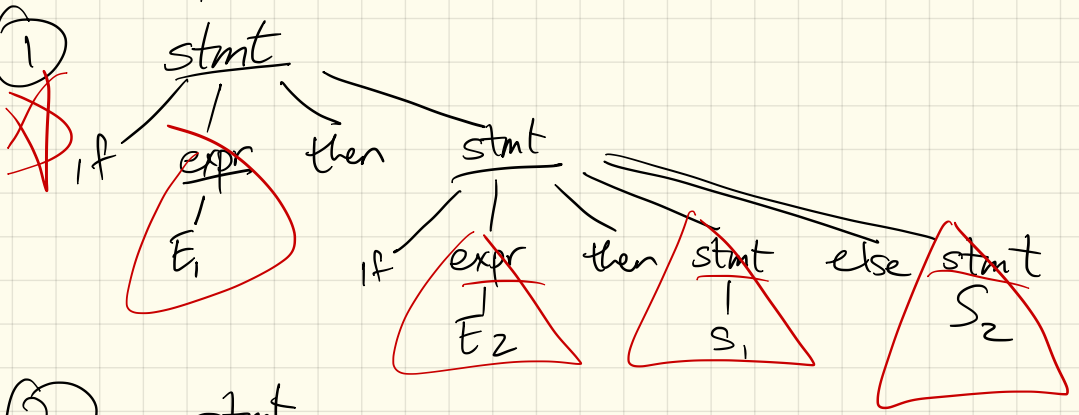Each parent/child shows one
   step of the derivation

- 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 

```
            stmt
         /   |    \
       if   expr  then        stmt
             |              /   |    \        \
            E_1          if  expr  then  stmt   else   stmt
                              |              |            |
                             E_2            S_1          S_2
```

② 

```
            stmt
         /   |    \         \         \
       if   expr  then      stmt     else    stmt
             |            /   |   \             \
            E_1         if  expr  then  stmt    S_2
                             |            |
                            E_2          S_1
```

# 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 → matched_stmt
 | unmatched_stmt

matched_stmt → if expr then matched_stmt else matched_stmt
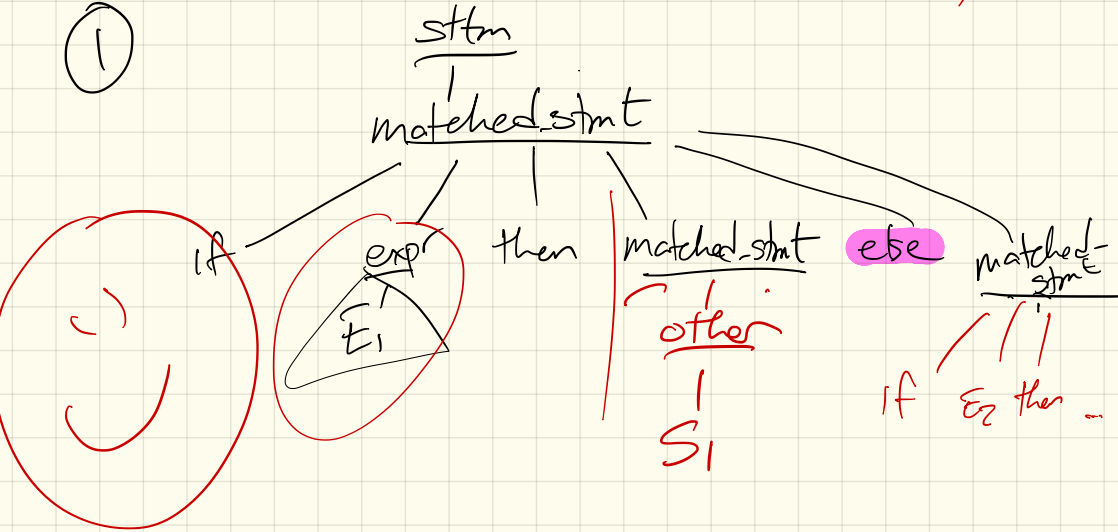 | other ← not an if statement

unmatched_stmt → if expr then stmt
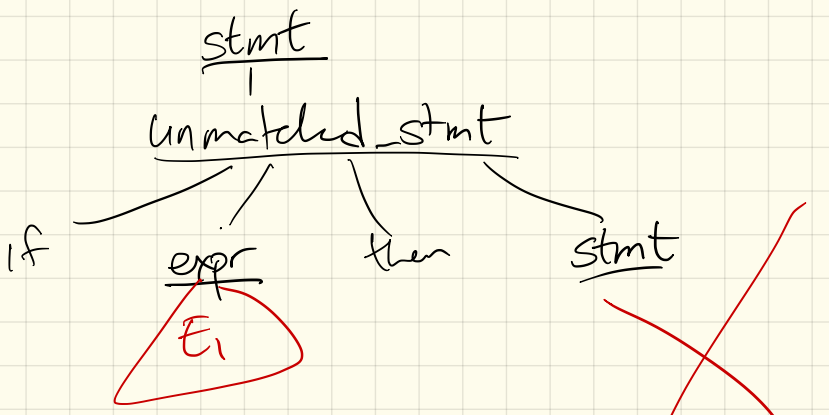 | if expr then matched_stmt else unmatched_stmt

next time — parsing

# Example:

if $E_1$ then ( $S_1$ else if $E_2$
then $S_2$ else $S_3$ )

① 

stm
|
matched_stmt

if — expr — then — matched_stmt — else — matched_stmt

$E_1$

other
|
$S_1$

if $E_2$ then ...

② 

stmt
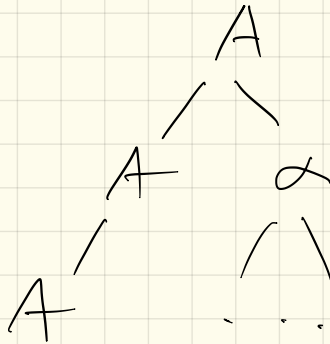|
unmatched_stmt

if — expr — then — stmt

$E_1$

⇒ only 1 parsing :)

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

$$A \rightarrow A \alpha$$

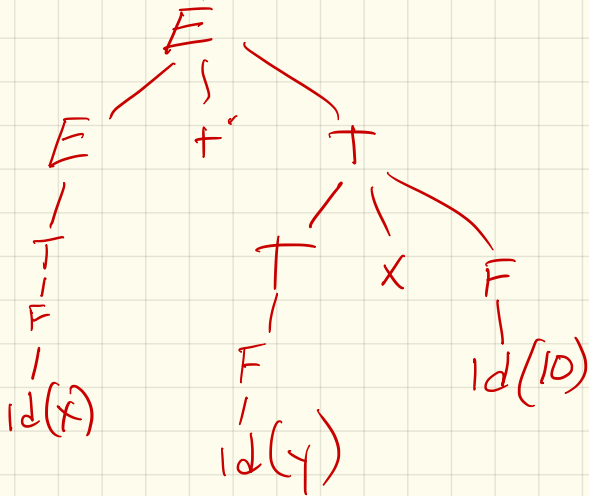These are bad for persers:



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

Ex:    $E \longrightarrow E + T \;/\; T$   (start)

$T \longrightarrow T \times F \;/\; F$

$F \longrightarrow (E) \;/\; id$

Parse:   $x + y \times 10$



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

To eliminate:

any other term/
non-term/lst

Rule:  $A \rightarrow A\alpha \mid \beta$

$$A \rightarrow \beta A'$$
$$A' \rightarrow \alpha A' \mid \varepsilon$$

Ex: On

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T \times F \mid F$$
$$F \rightarrow (E) \mid id$$

$\alpha$   $\beta$

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \varepsilon$$
$$T \rightarrow FT'$$
$$T' \rightarrow \times FT' \mid \varepsilon$$
$$F \rightarrow (E) \mid id$$

# 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

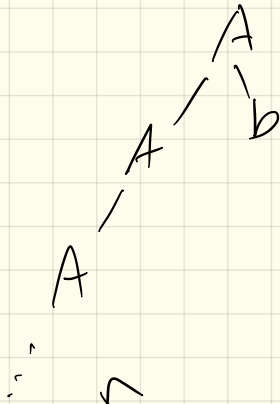Called predictive parsing.
Works well on LL (1)
Grammars.

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


Parse cad:

Rule: Stering 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!

$$A \rightarrow A b$$

A
A'  b
A'
A'
A
⋰

↙ 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 $(\alpha \curvearrowleft$ any string of non-terminals
$\qquad$ & terminals

$\qquad$ := set of possible first
terminals in any derivation
of $\alpha$ by the grammar

So:

1) if $x$ is a terminal,

$\qquad$ FIRST$(x)=$

2) if $X \to \varepsilon$ is a production,
$\qquad$ add $\varepsilon$ to FIRST $(x)$

3) If $X$ is a nonterminal:

$\qquad$ if $X \to Y_1 Y_2 \dots Y_k$ is a production:

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

add $\varepsilon$ if $\varepsilon$ is in FIRST$(Y_1), \dots$
$\qquad \qquad \qquad$ FIRST $(Y_k)$

**Ex:**

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \varepsilon$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT' \mid \varepsilon$$
$$F \rightarrow (E) \mid id$$

FIRST $(E) =$

FIRST $(E')$

FIRST $(T)$

FIRST $(T')$

FIRST $(F)$