# More parsing

## Today:

- HW due Thursday
  (via git)

# Last time

- More parsing: CFGs.
  - removing ambiguity (or recognizing)
  - eliminating left recursion
  
  ↳ set of rules to apply to get rid of LR

$$S \rightarrow x AB$$

Goal: parse (apply prod. rules) until all non-terminals (which means it is valid)

or until stuck.
  ↳ invalid

# Back to the practical:

- Any CFG can be parsed
  - ↳ Chomsky Normal Form
  - CYK algorithm
  - Run time: $O(n^3)$  :n

This is too slow!

Most modern parsers look
for certain restricted
families of CFGs.

Result: 2 main families:

→ <u>LL</u> : more limited, faster

<u>LR</u>: more general, slower

Both: $O(n)$ parsers

# LL:

- left to right persing
- leftmost derivation

Anything accepted by this type of parser is called an LL grammar.

## Recall:

Left to right:
on input strings, try to force a rule to recognize leftmost terminal first

Leftmost der:
if nonterm, try to resolve left one first

# Top down parsing (for LLs)

Called predictive parsing.

Works well on LL(1) ~ scan 1 token at a time
Grammars.

° Table based in practice

Simple Ex:   $S \rightarrow cAd \,/\, aAa \,/\, cAAA$
$A \rightarrow ab \,/\, a$
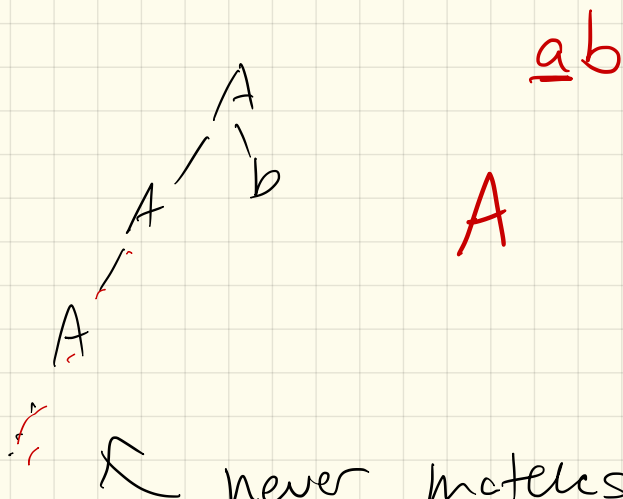
Parse   cad$

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

$S \rightarrow cA\underline{d}$

$\rightarrow cab\underline{d}$

$\rightarrow ca\underline{d}\$$

Note: Left recursion is
very bad on these!

$$A \rightarrow A b \,|\, a$$

<u>ab</u>

A

A
A   b
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.
- Scan 1 ~~token~~ ~~character~~ at a time, & guess which rule should match

## How?

- 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(s)
terminal )

LL(k):

k tokens to decide

(we'll just do LL(1))

Algorithm to construct table:
- Based upon listing "first"
+ "follow" sets for each
non-terminal.

(Essentially, these will encode
our predictions.)

# FIRST & FOLLOW Sets (for LL(1)):

FIRST $(\alpha)$ ← any string of non-terminals & terminals

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

So:

1) if $x$ is a terminal,

$$FIRST(x) = \{ \textcolor{red}{x} \}$$

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

3) If $X$ is a nonterminal:

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

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

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

Ex:
$$S \rightarrow E\$$$
$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \varepsilon$$

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

FIRST(S) = { (, id }
FIRST(E) = { (, id }
FIRST(E') = { +, ε }
FIRST(T) = { (, id }    F(F)
FIRST(T') = { *, ε }
FIRST(F) = { (, id }

# FOLLOW Sets:

(We'll assume any input ends in
$ , just to have an
end of file character)

## Rules:

1) Put $ in FOLLOW(s) ✓
    where S is start symbol.

2) Given a production: ✓
$$A \rightarrow \boxed{\alpha B \beta}$$
    everything in FIRST($\beta$) goes
    in FOLLOW(B)
    (except $\varepsilon$, if it is there).

3) Given a production:
$$A \rightarrow \alpha B$$
  or   $A \rightarrow \underline{\alpha B \beta}$  with $\varepsilon \in$ FIRST($\beta$)
    then everything in FOLLOW(A)
    also goes in FOLLOW(B)

Ex:  $S \to E\$$ ←

$E \to TE'$

$E' \to +TE' \mid \varepsilon$ ←

$T \to FT'$

$T' \to *FT' \mid \varepsilon$ ✓✓

$F \to (E) \mid id$ ✗

We had:

$FIRST(E) = FIRST(T) = FIRST(F)$
$\qquad = \{ (, id \}$

$FIRST(E') = \{ \underline{+}, \varepsilon \}$

$FIRST(T') = \{ \underline{*}, \varepsilon )$

So:

$FOLLOW(S) = \{ \$ \}$

$FOLLOW(E) = \{ +, ), *, \$ \}$

$FOLLOW(E') = \{ *, \$ \}$

$FOLLOW(T) = \{ *, \$$

$FOLLOW(T') = \{ *, \$$

$FOLLOW(F) = \{ *, \$$

Then, the Table: M: ~~(Next time)~~

For any production $X \rightarrow \alpha$, do

1) for each terminal $a$ in FIRST($\alpha$), add

$$X \rightarrow \alpha \quad \text{to} \quad M[A, a]$$

2) If $\varepsilon$ is in FIRST($\alpha$), add $X \rightarrow \alpha$ to $M[A, b]$

for each terminal $b$ in FOLLOW($A$).

If $\varepsilon$ is in FIRST($\alpha$) and $\$$ is in FOLLOW($A$), add $A \rightarrow \alpha$ to $M[A, \$]$.

Any other entries are errors.

(construct on board)

# End result:

Inputs

| Nonterminal | id | + | * | ( | ) | $ |
|---|---|---|---|---|---|---|
| E | E→TE′ | | | E→TE′ | | |
| E′ | | E′→+TE′ | | | E′→ε | E′→ε |
| T | T→FT′ | | | T→FT′ | | |
| T′ | | T′→ε | T′→*FT′ | | T′→ε | T′→ε |
| F | F→id | | | F→(E) | | |

# Then: parsing!

| Stack | Input | Action | Matched |
|-------|-------|--------|---------|
| E $ | id + id * id $ | | |