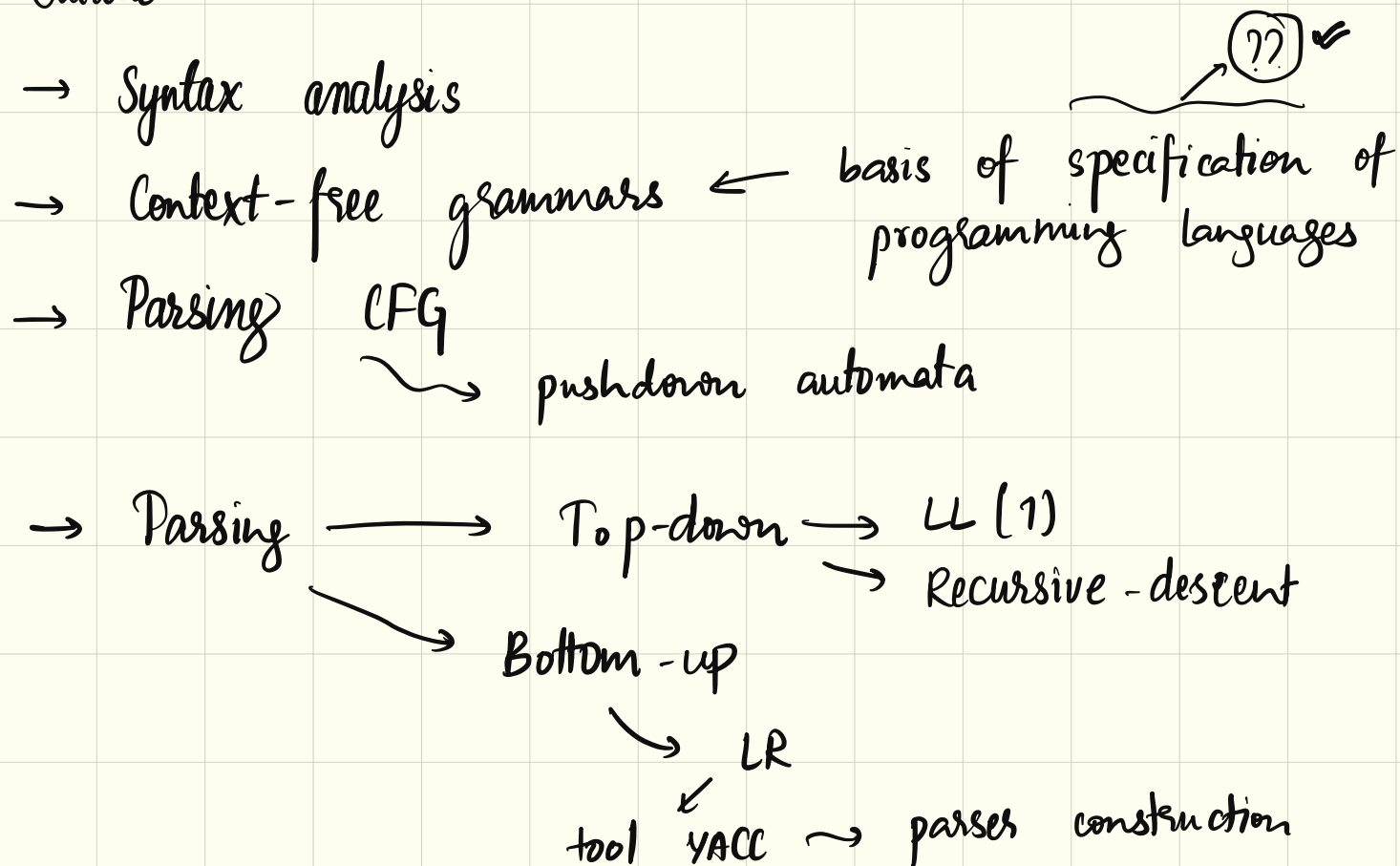


28 Apr 2025 - Compilers - I - Week 05

Context - Free Grammars, Pushdown Automata and Parsing

Outline



Grammars

→ used for precise description of rules

e.g.: Rules stating how functions are made out
of parameter lists, declarations and statements.
↓
expressions, etc.

→ Certain types of grammars



parsers can be automatically
constructed from the grammar

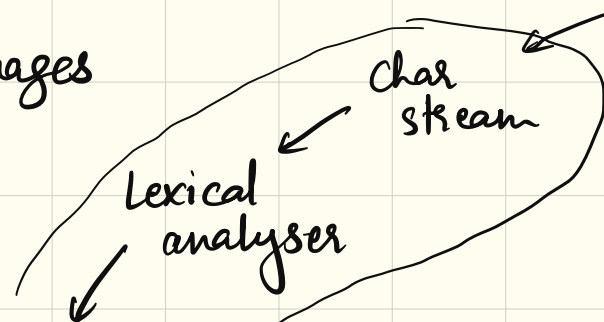
→ Parsers / Syntax Analysers are generated for a
a particular grammar

input
grammar → YACC → output
parser

→ CFG → subclass of languages

- regular
- CFG
- context-sensitive
- type-0 etc.

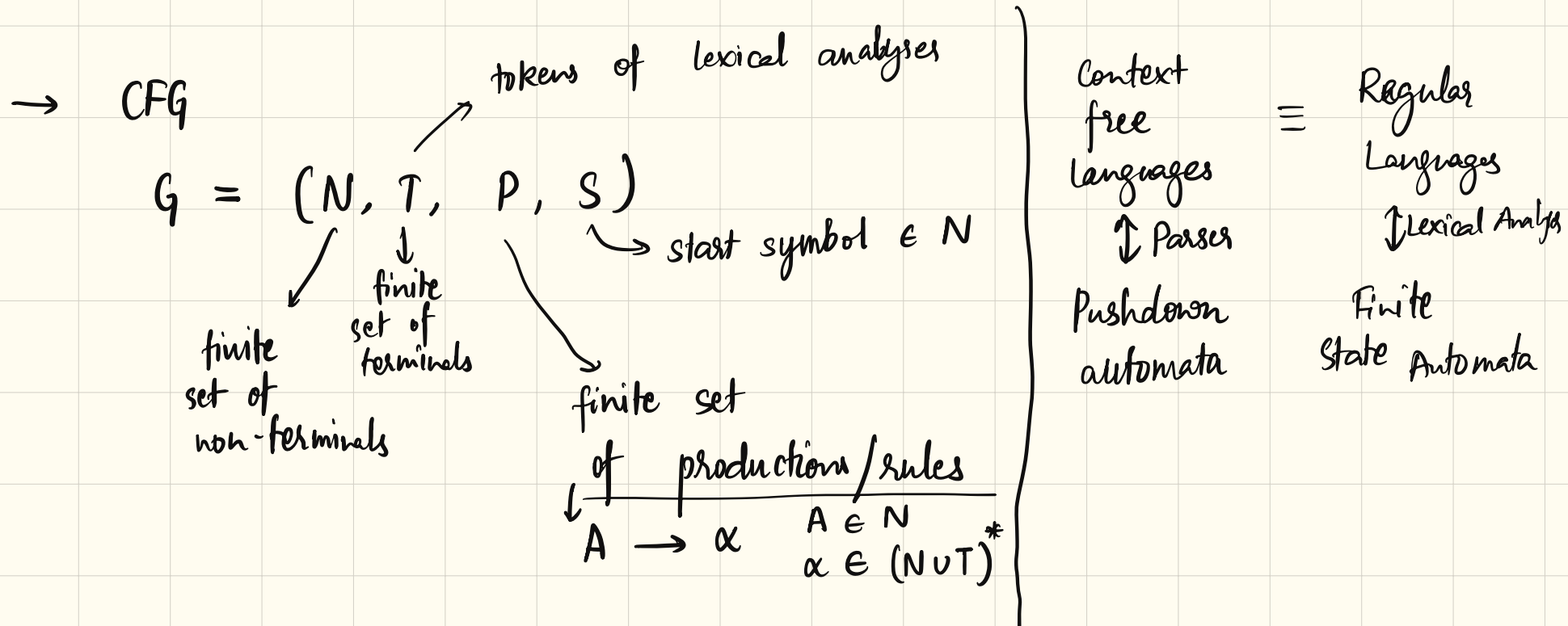
CFGs specify context-free languages



→ Parser

→ verifies that the string of tokens for a program in that language can be generated from that grammar

- reports syntax errors
- construct parse tree representation ← not always necessary
- usually calls lexical analyzer to supply tokens when necessary. (??)
- handwritten or automatically generated.



e.g.
start
(implicit)

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

$$S \rightarrow a S b$$

$$S \rightarrow \epsilon$$

loosely, $\epsilon \in T$

start

$$S \rightarrow aB \mid bA$$

$$A \rightarrow a \mid aS \mid bAA$$

$$B \rightarrow b \mid bS \mid aBB$$

equivalent

$$E' \rightarrow E + E \mid E * E \mid id$$

Derivations

$\underbrace{E}_{\text{start symbol}}$

$$E \rightarrow E + E$$

$$E + E$$

$$E \rightarrow id$$

$$id + E$$

depends on
what you
want to derive

Which production to choose

$$\Downarrow E \rightarrow id$$

$\frac{id + id}{\text{terminal string}}$

if $\{id, +\} \in T$

Short : $E \Rightarrow^* id + id$
 ~~~~~  
 0 or  
 more steps  
~~~~~  
 E derives $id + id$

Context - free Languages

CFGs $\xrightarrow{\text{generate}}$ CFLs
 $\xleftarrow{\text{specified by}}$

$$L(G) = \{ w \mid w \in T^* \text{ and } S \Rightarrow^* w \}$$

all strings derivable from start symbol and
 consist only of terminals

$$\left. \begin{array}{l}
 S \rightarrow 0S0 \\
 S \rightarrow 1S1 \\
 S \rightarrow 1 \\
 S \rightarrow \epsilon
 \end{array} \right\} L(G) = \text{palindromes over } 0 \text{ and } 1$$

Sentential form

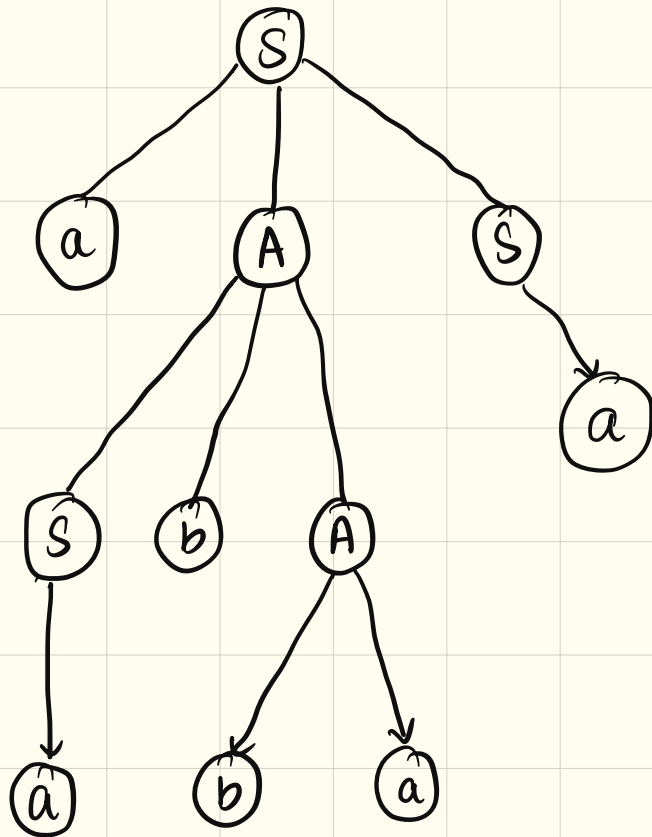
A string $\alpha \in (N \cup T)^*$

is a sentential form if $S \Rightarrow^* \alpha$

sentence \rightarrow only terminals

$$G_1 \equiv G_2 \iff L(G_1) = L(G_2)$$

Derivation Tree example



$$S \rightarrow aAS \mid a$$

$$A \rightarrow SbA \mid SS \mid ba$$

leaves = terminals

non-terminals = internal nodes

$$S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbas$$

$\Rightarrow \underline{aabbaa}$
yield of
derivation tree

Leftmost and Rightmost derivation

↓
expand
left-most
non-terminal
first in
each step

↓
expand
right-most
non-terminal
first in
each step

→ if some $w \in L(G)$ has more than one

parse tree \leadsto G is ambiguous

if not ambiguous \rightarrow unique leftmost and rightmost derivations

→ A CFL for which every G is
ambiguous \rightarrow inherently ambiguous language

Ambiguous Grammars

If the grammar is ambiguous \leadsto tools for generating parsers and parsing techniques will be in trouble.

\swarrow
2 parse trees
for same
grammar, same
word

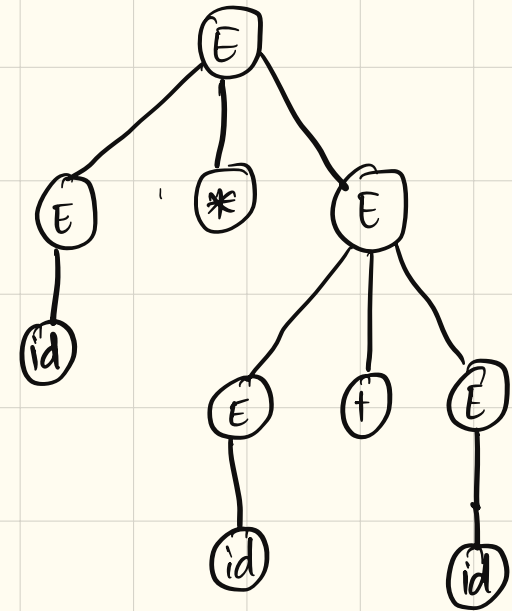
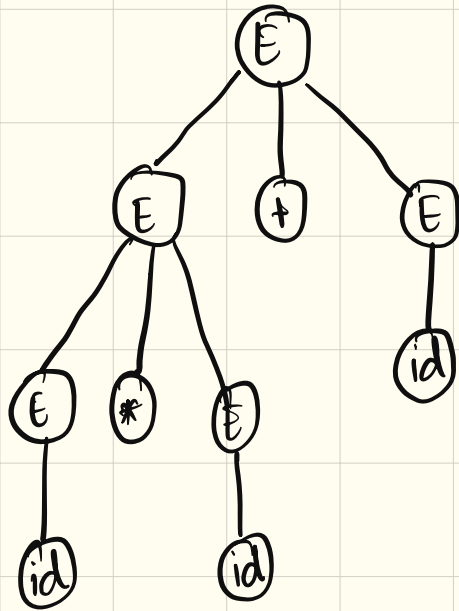
$E \rightarrow E + E \mid E * E \mid (E) \mid id$ $\left. \vphantom{E \rightarrow E + E \mid E * E \mid (E) \mid id} \right\} \begin{array}{l} + \text{ and } * \\ \text{have same} \\ \text{precedence} \end{array}$

\updownarrow same language

$\left. \begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid id \end{array} \right\} \begin{array}{l} * \text{ has more} \\ \text{precedence than } + \end{array}$

Example 1: $id * id + id$

$E \rightarrow E + E \mid E * E \mid (E) \mid id$



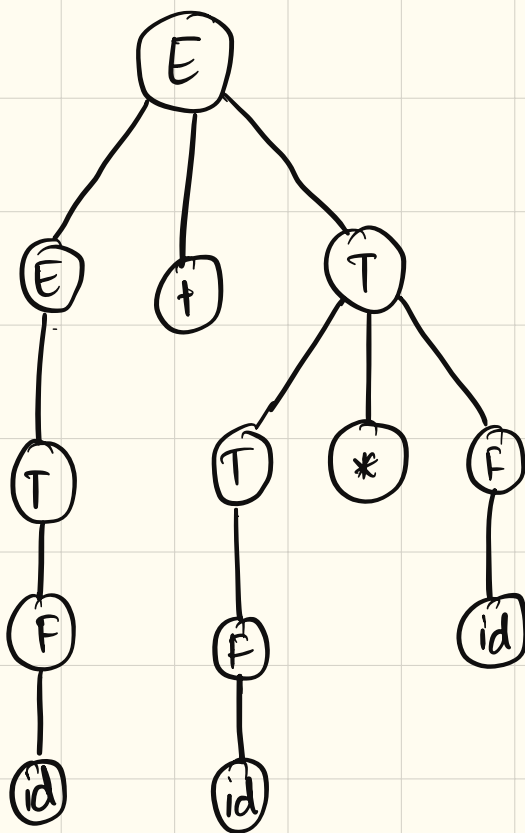
Equivalent Unambiguous Grammars

$id + id * id$

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

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

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



leftmost derivation (unique)

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \Rightarrow id + T \Rightarrow id + T * F \Rightarrow id + F * F \Rightarrow id + id * F \Rightarrow id + id * id$$

Trying first

leftmost derivation

$$E \Rightarrow T \Rightarrow T * F \Rightarrow F * F \Rightarrow (E) * F \Rightarrow (E + T) * F$$

cannot
go to
E
no such
rule

$$\Downarrow \\ (T + T) * F \\ \Downarrow$$

$$(id + F) * F \Leftarrow (id + T) * F \Leftarrow (F + T) * F$$

\Downarrow

$$(id + id) * F \Rightarrow \underline{(id + id) * id}$$

Ambiguity Example 2

$\text{stmt} \rightarrow \text{IF expr stmt} \mid \text{IF expr stmt ELSE stmt} \mid \text{other_stmt}$

↑
ambiguous

$\text{stmt} \rightarrow \text{IF expr stmt} \mid \text{IF expr matched_stmt ELSE stmt}$

$\text{matched_stmt} \rightarrow \text{IF expr matched_stmt ELSE}$
 $\text{matched_stmt} \mid \text{other_stmt}$

↙ equivalent
unambiguous (example later)

* No algorithm exists \rightarrow to convert ambiguous grammars to unambiguous

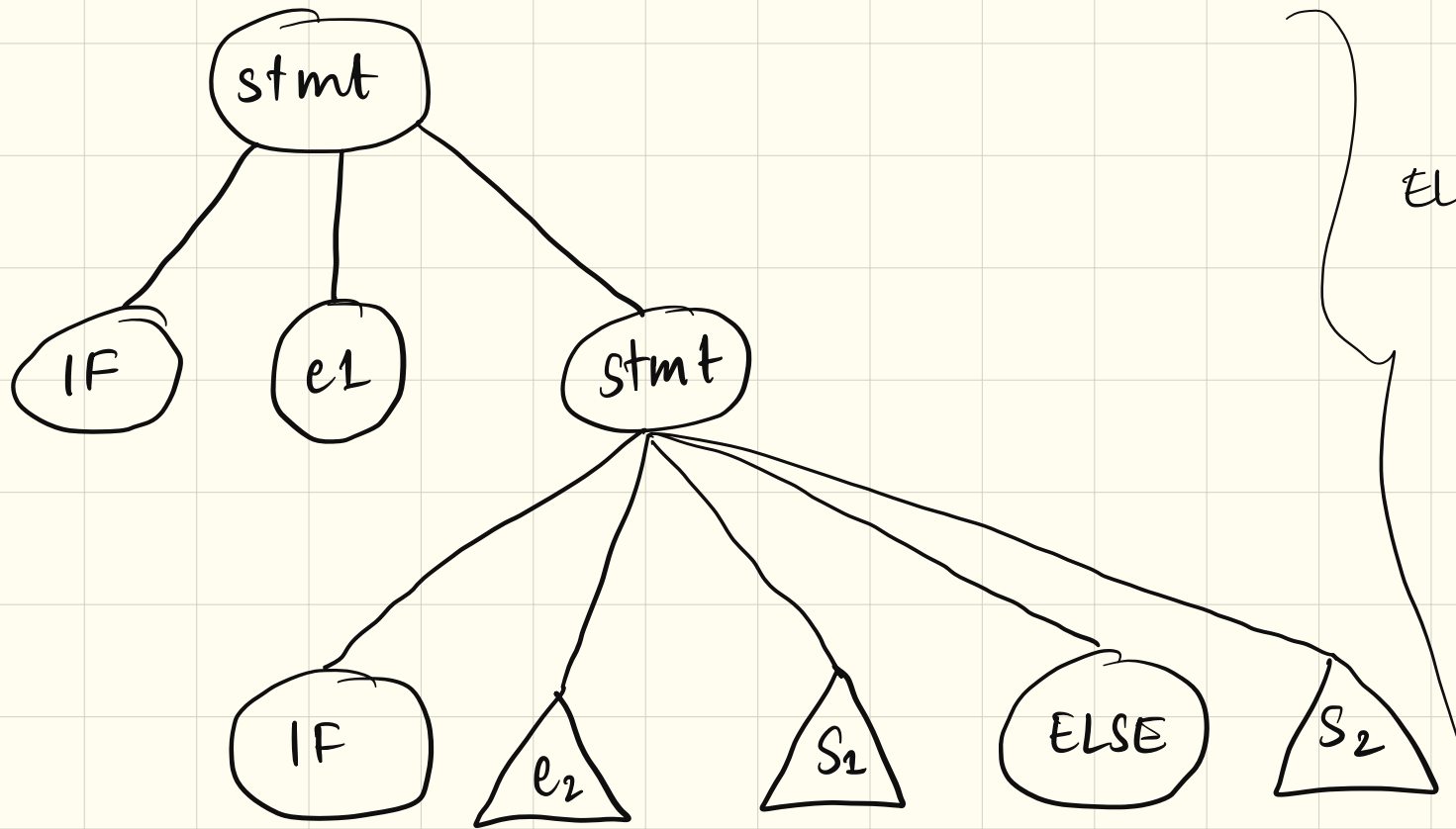
$$L = \{ a^n b^n c^m d^m \mid n, m \geq 1 \} \cup \{ a^n b^m c^m d^n \mid n, m \geq 1 \}$$

↳ inherently ambiguous

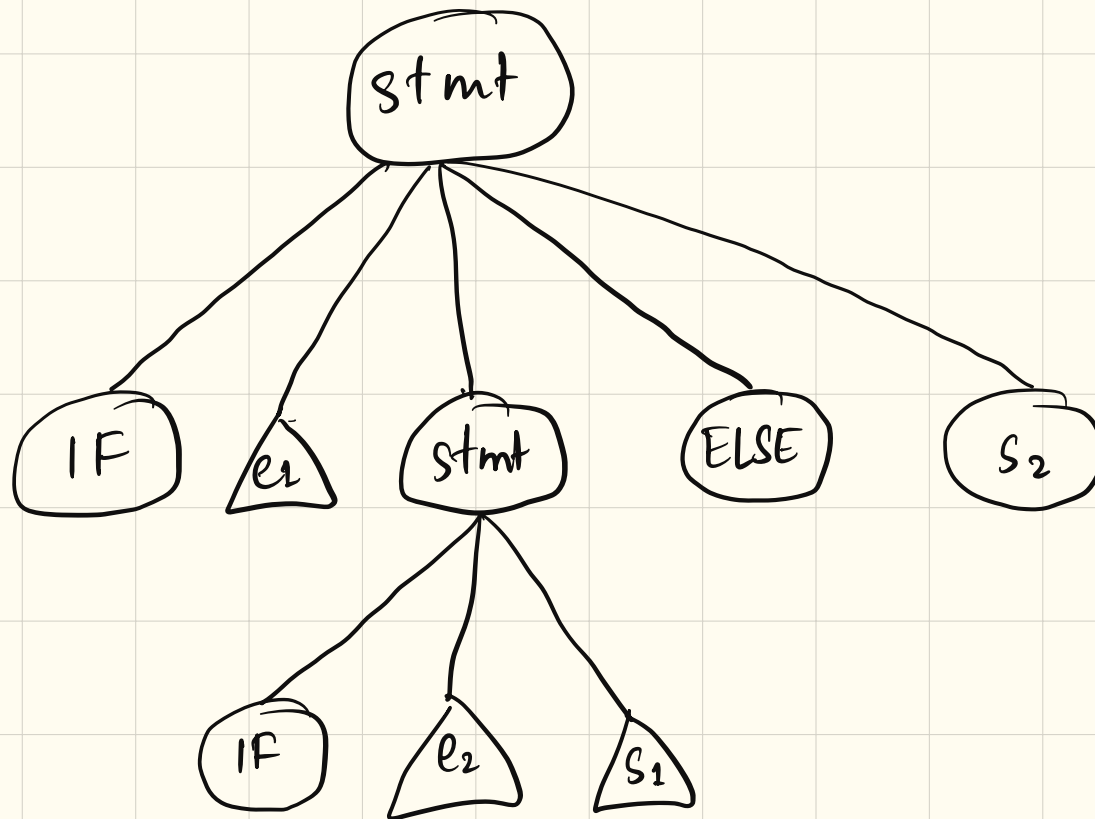
no grammar that
is not ambiguous

Example 2 : Two parse trees for the sentence

IF e1 IF e2 s1 ELSE s2

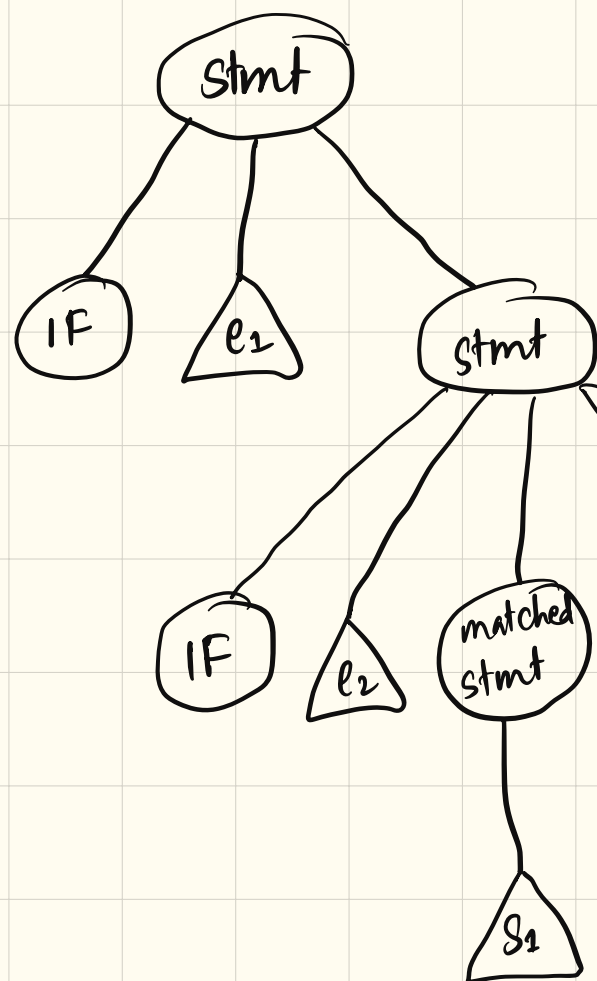


ELSE corresponds
to
inner if



ELSE corresponds
to outer if

$$S \rightarrow IF \quad e \quad S \quad | \quad IF \quad e \quad ms \quad ELSE \quad S$$

$$\underline{ms} \rightarrow IF \quad e \quad ms \quad ELSE \quad ms \quad | \quad other_S$$


ms = matched statement

cannot generate
if-then
statements

gives
rise to
ambiguity

always generates
an if-then-else
matched statement

Fragment of C grammar (Expressions)

??

logical_or_expr \rightarrow logical_and_expr

| logical_or_expr OR_OP

precedence { or is
at lower
level

logical_and_expr

AND is
at higher
level

logical_and_expr \rightarrow equality_expr

| logical_and_expr AND_OP

equality
expression

==

>

&&

lower priority \longrightarrow higher priority
| lower priority (operation) higher priority

Pushdown Automata

Machine \rightarrow can be automatically derived with a given grammar
 \rightarrow stack based system

PDA $M = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$

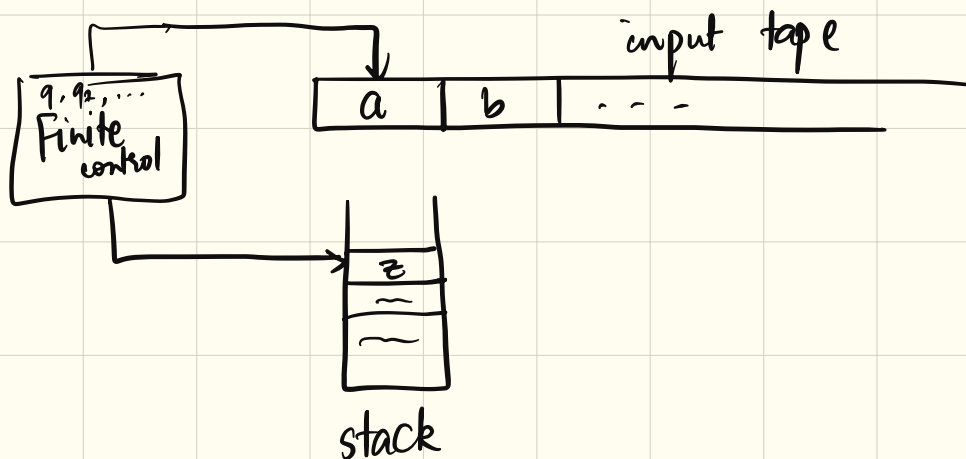
finite \swarrow set of states \swarrow input alphabet \downarrow stack alphabet \downarrow transition fn \swarrow start state \swarrow start symbol on stack \swarrow set of final states $F \subseteq Q$

$$\delta : Q \times \underbrace{\Sigma \cup \{\epsilon\}}_{\text{can be } \epsilon} \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma^*)$$

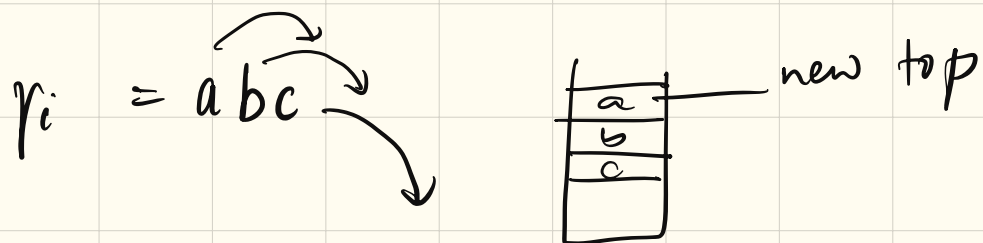
$$\delta(q, \overset{\substack{\uparrow \\ \text{input} \\ \text{symbol}}}{a}, \overset{\substack{\uparrow \\ \text{top of} \\ \text{stack} \\ \text{symbol}}}{z}) = \{ (q_1, \gamma_1), (q_2, \gamma_2), \dots, (q_m, \gamma_m) \}$$

PDA can enter any of the states q_i and replace the symbol z by string γ_i

Non-deterministic control



→ leftmost symbol of γ_i will be new top of stack



$L(M) \rightsquigarrow$ language accepted by M by final state

$$L(M) := \{ w \mid (q_0, w, Z_0) \xrightarrow{*} (p, \varepsilon, \gamma) \text{ for some } p \in F \text{ and } \gamma \in \Gamma^* \}$$

\uparrow start \uparrow input string \uparrow stack start \nwarrow series of moves

$N(M) \rightsquigarrow$ language accepted by M by empty stack

F becomes irrelevant
we usually set $F = \emptyset$
here

$$N(M) = \{ w \mid (q_0, w, z_0) \vdash^* (p, \varepsilon, \varepsilon) \text{ for some } p \in Q \}$$

Part 2

$$L = \{ ww^R \mid w \in \{a, b\}^+ \}$$

\downarrow
non-deterministically guess middle of input

Non deterministic and deterministic PDA

$NFA \equiv DFA$ finite automata

↓
accept the
same set of
languages

But

NPDA is more powerful than DPDA

↓
 $\{w w^R \dots\}$

$\{w c w^R\}$ → can be recognized by DPDA

In practice we need DPDA, since they have exactly one possible move at any instant.

→ Our parsers all are DPDA.

Parsing

Parsing is the process of constructing a parse tree
for a sentence generated by a given grammar.

↓
using a DPDA
+ some actions

no restrictions on language and the form of grammar \Rightarrow parsers for CFLs require $O(n^3)$ time (n = length of the string parsed)

- * CYK algorithm \rightarrow DP
- * Earley's algo

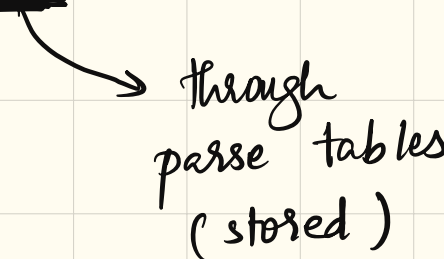
Subsets of CFLs \rightarrow $O(n)$ time
we are interested in these

- ① Predictive parsing \rightarrow on class of grammars called LL(1) grammars
 - } uses top-down parsing
- ② Shift-Reduce parsing
 - \downarrow using LR(1) grammars (bottom-up)

Top-down parsing using LL grammars

→ Basic idea: trace the leftmost derivation of the string while constructing the parse tree

$$\left. \begin{array}{l} S \rightarrow aAS \mid c \\ A \rightarrow ba \mid SB \\ B \rightarrow bA \mid S \end{array} \right\} \begin{array}{l} \text{considers } acbbac \\ \text{leftmost derivation} \end{array}$$
$$\begin{array}{ccccccc} \underset{1}{S} & \Rightarrow & \underset{2}{a} \underset{2}{A} S & \Rightarrow & \underset{3}{a} \underset{3}{S} B S & \Rightarrow & \underset{4}{ac} \underset{4}{B} S \Rightarrow \underset{5}{acb} \underset{5}{A} S \\ & & & & & & \Downarrow \\ & & & & & & acbba \underset{6}{S} \\ & & & & & & \Downarrow \\ & & & & & & acbbac \quad \underset{7}{} \end{array}$$

→ start from the start symbol, predicts the next production used in the derivation.  through parse tables (stored)

→ Next production to be used in the derivation is determined using the next input symbol to lookup the parsing table (look-ahead symbol)

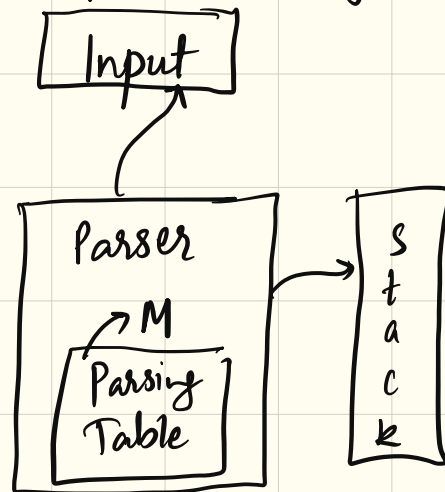
→ Restriction on the grammar ^{ensures} → no slot in the parsing table contains more than one production.

→ if more than one, we cannot decide which to use next

Parse table construction

→ if two productions become eligible to be placed in the same slot → grammar is declared unfit for predictive parsing.

LL(1) parsing algorithm



Initial configuration :

stack = S

Input = w \$

S = start symbol

\$ \equiv end of file marker

repeat {

let X be the top stack symbol;

let a be the next input symbol \rightarrow $*$ may be $\$ \quad * \quad /$

if X is a terminal symbol or $\$$ then

if $X == a$ then { $\left. \begin{array}{l} \text{pop } X \text{ from the stack;} \\ \text{remove } a \text{ from input;} \end{array} \right\}$ top of stack matches input

} else ERROR(); } stack symbol and input symbol do not match

else $/^*$ X is a non-terminal symbol $^*/$

if $M[X, a] == X \rightarrow Y_1 Y_2 \dots Y_k$ then {

pop X from stack;

push Y_k, Y_{k-1}, \dots, Y_1 onto stack }
(Y_1 on top)

Single unique production

} until stack has emptied.

Example:

Grammar

$S' \rightarrow S\$$

$S \rightarrow aAS \mid c$

$A \rightarrow ba \mid SB$

$B \rightarrow bA \mid S$

LL(1) Parsing Table
(No slot can have more than 1 entry)

construction
later

| | a | b | c | \$ |
|------|----------------------|--------------------|----------------------|----|
| S' | $S' \rightarrow S\$$ | | $S' \rightarrow S\$$ | |
| S | $S \rightarrow aAS$ | | $S \rightarrow c$ | |
| A | $A \rightarrow SB$ | $A \rightarrow ba$ | $A \rightarrow SB$ | |
| B | $B \rightarrow S$ | $B \rightarrow bA$ | $B \rightarrow S$ | |

Rows indexed by non-terminals

Columns indexed by terminals

String : a c b b a c

a

| |
|----|
| |
| |
| |
| S' |

a

| |
|----|
| |
| |
| S |
| \$ |

a

| |
|----|
| a |
| A |
| S |
| \$ |

c

| |
|----|
| |
| A |
| S |
| \$ |

c

| |
|----|
| S |
| B |
| S |
| \$ |

c

| |
|----|
| c |
| B |
| S |
| \$ |

b

| |
|----|
| |
| B |
| S |
| \$ |

b

| |
|----|
| b |
| A |
| S |
| \$ |

b

| |
|----|
| |
| A |
| S |
| \$ |

b

| |
|----|
| b |
| a |
| S |
| \$ |

| a | c | c | \$ | stack empty |
|----|----|----|----|-------------|
| | | | | |
| a | | | | |
| S | S | c | | |
| \$ | \$ | \$ | \$ | |

How is the parsing table constructed ?

LL(1) grammars \leadsto sub class of CFGs
some restrictions

Strong LL(k) grammars

Let the given grammar be G .

→ Input is extended with k EOF symbols
 $\k

$k \rightarrow$ lookahead of the grammar

$k = 1 \rightarrow$ allowed to see 1 symbol at a time

$k = 2 \rightarrow$ allowed to see 2 symbols at a time

New non-terminal S' and production : $S' \rightarrow S \k

→ Consider leftmost derivations only

assume grammar has no useless symbols

(terminals / non-terminals
that are never used)

① not part of
any production

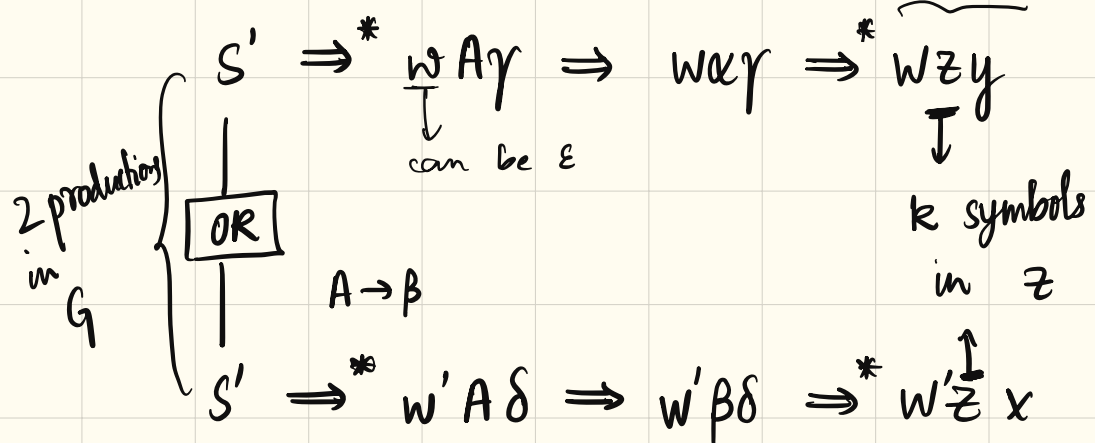
OR

② LHS of a
production
can never
be reached

such
productions
and
corresponding
symbols are
useless

→ A production $A \rightarrow \alpha$ in G is called a strong $LL(k)$ production,

if in G



$|z| = k, z \in \Sigma^*, w \text{ and } w' \in \Sigma^*, \text{ then } \alpha = \beta$

$\gamma \neq \delta$

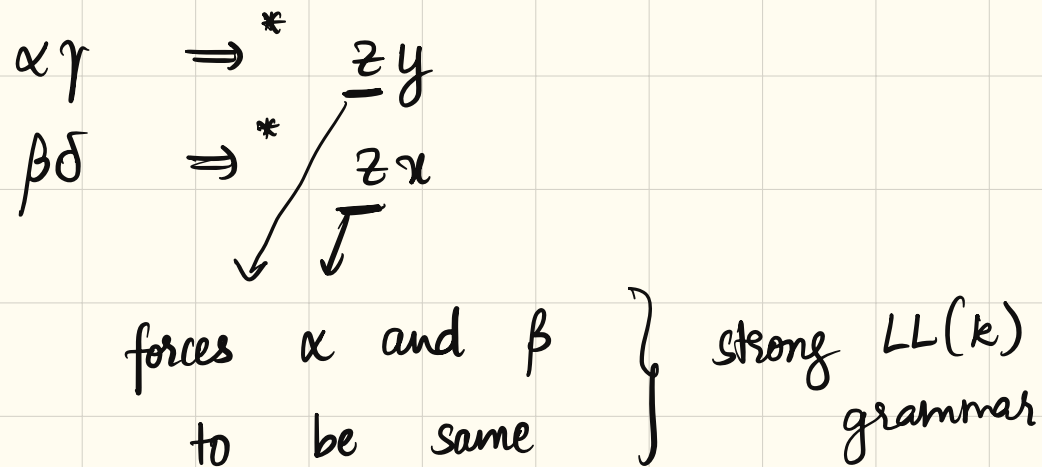
$z = \text{lookahead}$

Can we look at string z at some point, and determine whether $A \rightarrow \alpha$ was applied, or $A \rightarrow \beta$??

not necessary

Strong $LL(k)$ condition: If the lookahead (z) is same at some point, then we know exactly which production was applied at a point

* A grammar (non-terminal) is strong $LL(k)$ if all its productions are strong $LL(k)$.



e.g: $S \rightarrow Abc \mid aAc b$

$A \rightarrow \epsilon \mid b \mid c$

S is a strong $LL(1)$ non-terminal
 $k=1$

* $S' \Rightarrow S \$ \Rightarrow \frac{Abc}{\alpha} \$ \Rightarrow \begin{matrix} bc \$ \\ A \rightarrow \epsilon \\ z = b \end{matrix} \text{ or } \begin{matrix} bbc \$ \\ A \rightarrow b \\ z = b \end{matrix} \text{ or } \begin{matrix} cbc \$ \\ A \rightarrow c \\ z = c \end{matrix}$

$w = \epsilon$ here

1st symbol = lookahead = z
($k=1$)

* $S' \Rightarrow S \$ \Rightarrow \frac{aAc b}{\beta} \Rightarrow \begin{matrix} acb \$ \\ A \rightarrow \epsilon \end{matrix} \text{ or } \begin{matrix} abcb \$ \\ A \rightarrow b \end{matrix} \text{ or } \begin{matrix} accb \$ \\ A \rightarrow c \end{matrix}$
 $z = a$ in all cases

In this case, $w = w' = \epsilon$,

$$\alpha = Abc$$

$$\beta = aAc b$$

but z is different in the two derivations,
in all the derived strings

\therefore vacuously true $\rightarrow LL(1)$ \checkmark

$$\left. \begin{array}{l} S \rightarrow Abc \mid aAc b \\ A \rightarrow \epsilon \mid b \mid c \end{array} \right\} A \text{ is } \underline{\text{NOT}} \text{ strong } LL(1)$$

$$\begin{array}{lll}
 S' \Rightarrow^* A b c \$ \Rightarrow b c \$ & A \rightarrow \underbrace{\epsilon}_{\alpha} & z = b \\
 & A \rightarrow \underbrace{b}_{\beta} & z = b \\
 & & \left. \begin{array}{l} z = b \\ z = b \end{array} \right\} \rightarrow \therefore \text{Not LL(1)} \\
 c b c \$ & & w = \epsilon \\
 & & w' = \epsilon
 \end{array}$$

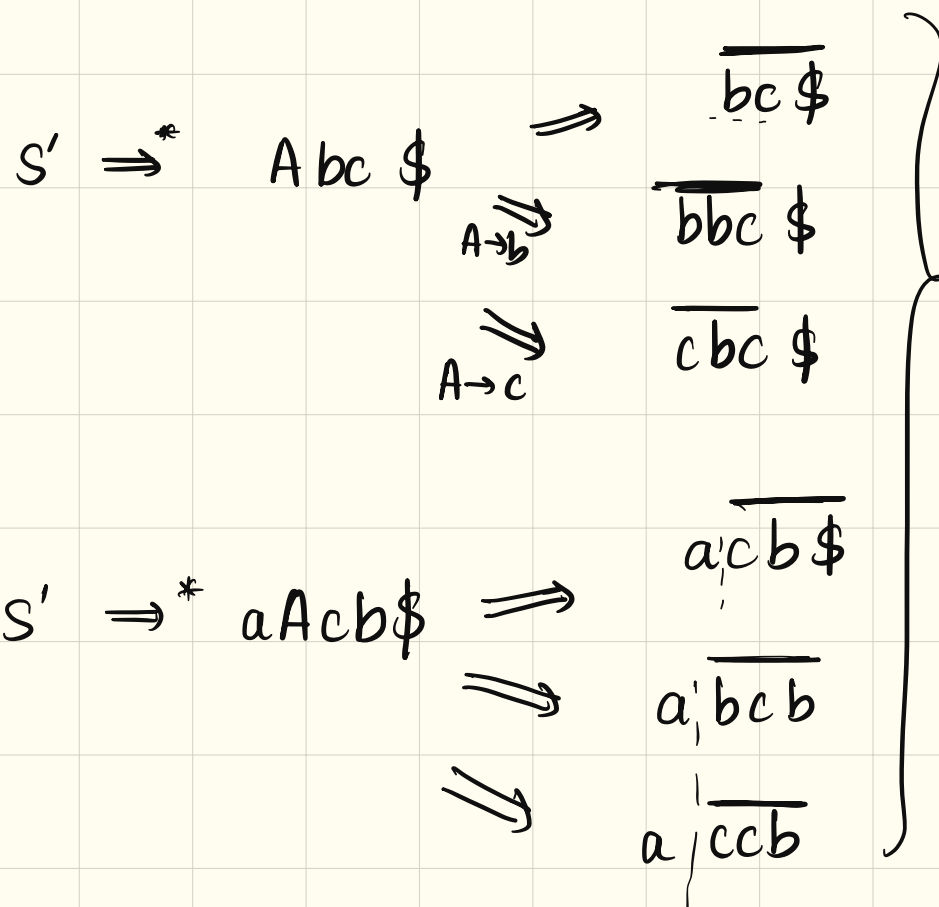
Even though the lookaheads are same ($z = b$),
 $\alpha \neq \beta$, \therefore grammar is not strong LL(1).

A is not strong LL(2)

$$w = \epsilon, w' = a$$

$$\begin{array}{lll}
 S' \Rightarrow^* A b c \$ \Rightarrow \underline{b c} \$ & A \rightarrow \epsilon & \left. \begin{array}{l} A \rightarrow \epsilon \\ A \rightarrow b \end{array} \right\} z = bc \\
 S' \Rightarrow^* a A c b \$ \Rightarrow a \underline{b c} b \$ & A \rightarrow b & \\
 \hline
 \therefore A \text{ is not LL(2)}
 \end{array}$$

A is strong LL(3)



vacuously true

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

$S' \Rightarrow^* w A \gamma \Rightarrow^* w z y$
 $S' \Rightarrow^* w' A \delta \Rightarrow^* w' z \alpha$

if A is LL(|z|) nonterminal

same
 \Downarrow
 $\alpha = \beta$

Testable conditions for LL(1)

→ Lookaheads longer than 1 will not be considered from now on

→ $\left. \begin{array}{l} \text{ordinary} \\ \text{weak LL(1)} \\ \text{strong LL(1)} \end{array} \right\} \rightarrow \begin{array}{l} \text{for lookahead} = 1 \\ \text{strong} = \text{weak} \\ \text{for others (LL(2), strong} \neq \text{weak)} \end{array}$

The classical condition for LL(1) property uses FIRST
and FOLLOW sets

$\text{def}^n \rightarrow \text{cond}^n$
next

FIRST

If α is any string of grammar symbols
 $\alpha \in (N \cup T)^*$

$$\text{FIRST}(\alpha) = \{ a \mid a \in T \text{ and } \alpha \Rightarrow^* a\alpha, \alpha \in T^* \}$$

↳ ① collect all (terminal) strings (sentence)
derivable from α

② add first symbol of each string to $\text{FIRST}(\alpha)$

By defⁿ : $\text{FIRST}(\epsilon) = \{ \epsilon \}$

FOLLOW

defined only for a non-terminal

} $\overbrace{\text{FIRST}}^{\downarrow}$
 $(NUT)^*$

If A is a nonterminal, then

$$\text{FOLLOW}(A) = \{ a \mid S \Rightarrow^* \alpha A a \beta, \alpha, \beta \in (NUT)^* \}$$

$$a \in T \cup \{\$ \}$$

① Start with S , derive all sentential in which A appears } \rightarrow requires a context

② take the first terminal that appears after A . add it to $\text{FOLLOW}(A)$

If $\alpha A a \beta$ cannot be reached from S

 \rightarrow a not included in $\text{FOLLOW}(A)$

Example:

$$S' \rightarrow S\$$$

$$S \rightarrow aAS \mid c$$

$$A \rightarrow ba \mid SB$$

$$B \rightarrow bA \mid s$$

$$* \text{ FIRST}(S') = \text{FIRST}(S) = \{a, c\} \quad \text{because}$$

$$S' \Rightarrow S\$ \Rightarrow \underline{c}\$$$

$$\text{FIRST}(\alpha) = \{a \mid a \in T \text{ and } \alpha \Rightarrow^* a\alpha, \alpha \in T^*\}$$

$$S' \Rightarrow S\$ \Rightarrow \underline{a}AS\$ \Rightarrow \underline{a}baS\$ \Rightarrow \underline{a}bac\$$$

NOTE: Do the complete derivation
 $\alpha \Rightarrow^* a\alpha, \alpha \in T^*$ } check if terminal string possible

* $\text{FIRST}(A) = \{a, b, c\}$ because

$$A \Rightarrow \underline{b}a$$

$$A \Rightarrow SB$$

$$\therefore \text{FIRST}(S) \subseteq \text{FIRST}(A)$$

" $\{a, c\}$

Incomplete method?
(??)

* $\text{FOLLOW}(S) = \{a, b, c, \$\}$ because

$$S' \Rightarrow S \underline{\$}$$

$$S' \Rightarrow aAS\$ \Rightarrow aSBS\$ \xRightarrow{B \rightarrow bA} aS\underline{b}AS\$$$

$$\xRightarrow{b \rightarrow S} aSSS\$$$

$\swarrow \quad \searrow$
 $\underline{a}AS \quad \underline{c}$

* FOLLOW (A) = { a, c } because

$$S' \Rightarrow^* a A S \$ \Rightarrow a A \underline{a} A S \$$$

$$S' \Rightarrow^* a A S \$ \Rightarrow a A \underline{c} \$$$

Algorithms to compute FIRST

- ① terminals and non-terminals
- ② general string $\in (N \cup T)^*$

① Terminals and nonterminals

{

for each $(a \in T)$ $\text{FIRST}(a) = \{a\}$; $\text{FIRST}(\epsilon) = \{\epsilon\}$; // terminals

for each $(A \in N)$ $\text{FIRST}(A) = \emptyset$; // initialize for nonterminals

// fixed point computation : see also \rightarrow Ford-Fulkerson Algorithm

while (FIRST sets are still changing) {

for each production p {

Let p be the production $A \rightarrow X_1 X_2 \dots X_n$

$$\text{FIRST}(A) = \underbrace{\text{FIRST}(A)}_{\text{could have elements from other production}} \cup \underbrace{(\text{FIRST}(X_1) - \{\epsilon\})}_{\text{// if } X_1 \Rightarrow^* \epsilon}$$

$i = 1;$

while $(\epsilon \in \text{FIRST}(X_i) \ \&\& \ i \leq n-1) \{$

$\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(X_{i+1}) - \{\epsilon\});$

$i++;$

$\}$

// $X_1 \rightarrow \epsilon \leadsto$ consider $X_2 \dots X_n$
// $X_1 X_2 \rightarrow \epsilon \leadsto$ consider next
and so on

if $(i == n \ \&\& \ (\epsilon \in \text{FIRST}(X_n)))$ // reached end
// no symbol

$\text{FIRST}(A) = \text{FIRST}(A) \cup \{\epsilon\}$ \leadsto // $A \Rightarrow^* \epsilon$ (✓)

$\}$

$\}$

② $\text{FIRST}(\beta)$: β , a string of grammar symbols

{ /* Assume $\text{FIRST}(\text{terminal})$, $\text{FIRST}(\text{nonterminal})$ computed */

$\text{FIRST}(\beta) = \phi$

while (FIRST sets are still changing) {

Let β be the string $x_1 x_2 \dots x_n$

$\text{FIRST}(\beta) = \text{FIRST}(\beta) \cup (\text{FIRST}(x_1) - \{\epsilon\})$;

$i = 1$

while ($\epsilon \in \text{FIRST}(x_i)$ && $i \leq n-1$) {

$\text{FIRST}(\beta) = \text{FIRST}(\beta) \cup (\text{FIRST}(x_i) - \{\epsilon\})$;

$i++$;

}

if ($i == n$ && ($\epsilon \in \text{FIRST}(x_n)$)

$$\text{FIRST}(\beta) = \text{FIRST}(\beta) \cup \{\epsilon\}$$

}

Example:

$$S' \rightarrow S\$$$

$$S \rightarrow aAS \mid \epsilon$$

$$A \rightarrow ba \mid SB$$

$$B \rightarrow cA \mid S$$

Initially : $\text{FIRST}(S) = \text{FIRST}(A) = \text{FIRST}(B) = \emptyset$

Iteration 1 :

$$* \text{FIRST}(S) = \{a, \epsilon\} \quad \text{from } S \rightarrow aAS \mid \epsilon$$

$$\begin{aligned} * \text{FIRST}(A) &= \{b\} \cup (\text{FIRST}(S) - \{\epsilon\}) \cup (\text{FIRST}(B) - \{\epsilon\}) \\ &= \{b, a\} \end{aligned}$$

$\downarrow_{a, \epsilon}$ $\downarrow_{\emptyset \text{ atom}}$ $\downarrow_{\text{Because } S \rightarrow \epsilon}$

$$\begin{aligned} * \text{FIRST}(B) &= \{c\} \cup (\text{FIRST}(S) - \{\epsilon\}) \cup \{\epsilon\} \\ &= \{a, c, \epsilon\} \end{aligned}$$

\downarrow
included
because $\epsilon \in \text{FIRST}(S)$

Iteration 2 (values stabilize and do not change in iteration 3)

* $\text{FIRST}(S) = \{a, \epsilon\}$ (no change from iteration 1)

$$* \text{ FIRST}(A) = \{b\} \cup (\text{FIRST}(S) - \{\epsilon\}) \cup (\text{FIRST}(B) - \{\epsilon\}) \cup \{\epsilon\}$$

$\underbrace{\hspace{1.5cm}}_{\text{contains } \epsilon} \quad \underbrace{\hspace{1.5cm}}_{\text{contains } \epsilon}$

$$= \{b, a, c, \epsilon\}$$

* $\text{FIRST}(B) = \{c, a, \epsilon\}$ no change from iteration 1.

Computation of FOLLOW

{

for each $(x \in N \cup T)$ FOLLOW $(x) = \phi$;

initialize follow of every nonterminal

FOLLOW $(S) = \{\$ \}$; /* S is the start symbol of the grammar */

repeat {

for each production $A \rightarrow X_1 X_2 \dots X_n$ { /* $X_i \neq \epsilon$ */

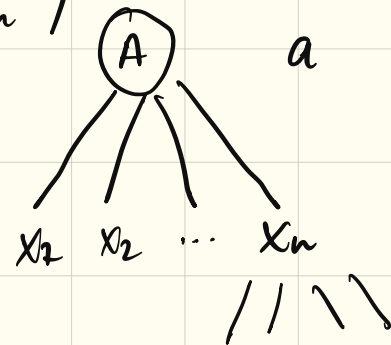
$A \rightarrow \epsilon$ not allowed

/* To compute the follow of A , we need to look at the production in which the RHS contains an A */

/* \therefore follow (A) requires a context $S' \Rightarrow^* \alpha A \beta$ */

/ * symbols which follow A will also follow x_n */

$$\text{FOLLOW}(x_n) = \underbrace{\text{FOLLOW}(x_n)}_{\substack{\text{from prev} \\ \text{iterations}}} \cup \text{FOLLOW}(A);$$



REST = FOLLOW(A);

shouldn't this be x_n ??

for $i = n \rightarrow 2$ {

if ($\epsilon \in \text{FIRST}(x_i)$) {

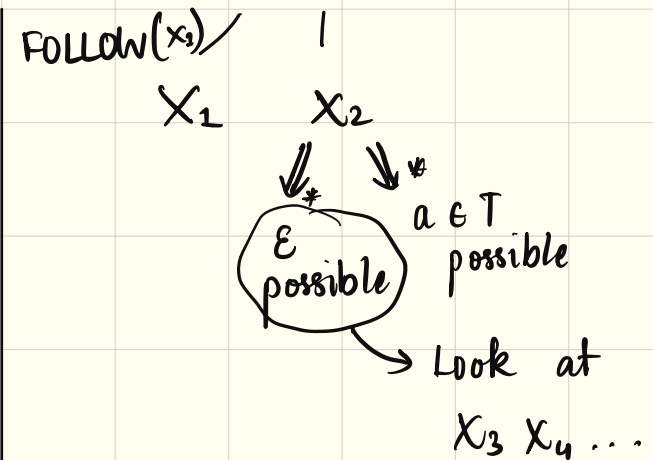
$\text{FOLLOW}(x_{i-1}) = \text{FOLLOW}(x_{i-1}) \cup$

$\text{FIRST}(x_i) - \{\epsilon\} \cup \text{REST};$

REST = FOLLOW(x_{i-1})

}

else {



$$\text{FOLLOW}(x_{i-1}) = \text{FOLLOW}(x_{i-1}) \cup \text{FIRST}(x_i);$$

$$\text{REST} = \text{FOLLOW}(x_{i-1});$$

}

}

} until no FOLLOW set has changed.

Example:

$$S' \rightarrow S\$$$

$$S \rightarrow aAS \mid \epsilon$$

$$A \rightarrow ba \mid SB$$

$$B \rightarrow cA \mid S$$

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

$\text{FOLLOW}(A) = \text{FOLLOW}(B) = \emptyset$

$\text{FIRST}(S) = \{a, \epsilon\}$

$\text{FIRST}(A) = \{a, b, c, \epsilon\}$

$\text{FIRST}(B) = \{a, c, \epsilon\}$

Iteration 1

* $S \rightarrow aAS$

$\text{FOLLOW}(S) \cup \{\$ \} ;$

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

$\text{FOLLOW}(A) \cup (\text{FIRST}(S) - \{\epsilon\}) \cup \text{REST} = \{a, \$ \}$

* $A \rightarrow SB$

$$\text{FOLLOW}(B) \cup \text{FOLLOW}(A) = \{a, \$\}$$

$$\text{REST} = \text{FOLLOW}(A) = \{a, \$\}$$

$$\text{FOLLOW}(S) \cup \text{FIRST}(B - \{\epsilon\}) \cup \text{REST} = \{a, c, \$\}$$

* $B \rightarrow cA$

$$\text{FOLLOW}(A) \cup \text{FOLLOW}(B) = \{a, \$\}$$

* $B \rightarrow S$

$$\text{FOLLOW}(S) \cup \text{FOLLOW}(B) = \{a, c, \$\}$$

iteration 1 : $\text{FOLLOW}(S) = \{a, c, \$\}$; $\text{FOLLOW}(A) = \text{FOLLOW}(B) = \{a, \$\}$

Iteration 2

$$* S \rightarrow aAS$$

$$\text{FOLLOW}(S) \quad U = \text{FOLLOW}(S) = \{a, c, \$\};$$

$$\text{REST} = \text{FOLLOW}(S) = \{a, c, \$\};$$

$$\text{FOLLOW}(A) \quad U = (\text{FIRST}(S) - \{\epsilon\}) \cup \text{REST} = \{a, c, \$\} \quad (\text{changed!})$$

$$* A \rightarrow SB$$

$$\text{FOLLOW}(B) \quad U = \text{FOLLOW}(A) = \{a, c, \$\} \quad (\text{changed!})$$

$$\text{REST} = \text{FOLLOW}(A) = \{a, c, \$\}$$

$$\text{FOLLOW}(S) \quad U = (\text{first}(B) - \{\epsilon\}) \cup \text{REST} = \{a, c, \$\}$$

(no change)

* At the end of iteration 2

$$\text{FOLLOW}(S) = \text{FOLLOW}(A) = \text{FOLLOW}(B) = \{a, c, \$\}$$

* FOLLOW sets do not change any further

LL(1) conditions

→ based on FIRST and FOLLOW

why? → we want an algorithm
to compute the parsing table

for LL(1) grammars.

↳ based on FIRST and FOLLOW

* Let G be a context-free grammar.

* G is $LL(1)$ iff for every pair of productions

$$A \rightarrow \alpha$$

$$A \rightarrow \beta,$$

the following condition holds:

$$\text{dirsymbol}(\alpha) \cap \text{dirsymbol}(\beta) = \emptyset$$

first
symbol in
the
derivation
≠ ϵ

$\text{dirsymbol}(\gamma) =$ if $(\epsilon \in \text{first}(\gamma))$ then

α, β

$((\text{first}(\gamma) - \{\epsilon\}) \cup \text{follow}(A))$

else $\text{first}(\gamma)$

direction symbol
set

* Equivalent formulation ALSU \rightarrow Dragon book } Ali Sethi when
 Ravi Sethi walks
 in ...

$$\text{first}(\alpha \cdot \underbrace{\text{follow}(A)}_{\substack{\text{set of} \\ \text{terminals} \\ \{a, b, c, d\}}}) \cap \text{first}(\beta \cdot \text{follow}(A)) = \phi$$

$$\{ \beta a, \beta b, \dots \}$$

first of $\{ \alpha a, \alpha \beta, \dots \}$

same as disymb

Given this definition for LL(1) grammar,

we now give an algorithm for constructing
 the parsing table.

for each production $A \rightarrow \alpha$

for each symbol $s \in \text{dirtsymb}(\alpha)$

/* s may be either a terminal symbol or $\$$ */

add $A \rightarrow \alpha$ to $\text{LLPT}[A, s]$

Make each undefined entry of LLPT as error.

OR

for each production $A \rightarrow \alpha$

for each terminal symbol $a \in \text{first}(\alpha)$

add $A \rightarrow \alpha$ to $\text{LLPT}[A, a]$

if ($\epsilon \in \text{first}(\alpha)$)

for each terminal symbol $b \in \text{follow}(A)$

add $A \rightarrow \alpha$ to $\text{LLPT}[A, b]$

if $\$ \in \text{follow}(A)$

add $A \rightarrow \alpha$ to $\text{LLPT}[A, \$]$

Make each undefined entry as error.

* After the construction of the $\text{LL}(1)$ table is complete \rightarrow using any method

if any slot in the table has two or more productions $\left. \vphantom{\begin{matrix} \text{if any slot in the} \\ \text{table has} \\ \text{two or more} \\ \text{productions} \end{matrix}} \right\} \rightarrow$ then the grammar is not $\text{LL}(1)$

Example:

$P_1: S \rightarrow \text{if } (a) \text{ else } S \mid \text{while } (a) S \mid \text{begin } SL \text{ end}$

$P_2: SL \rightarrow SS'$

$P_3: S' \rightarrow ; SL \mid \epsilon$

$\{\text{if, while, begin, end, a, (,), ;}\} \rightarrow \text{all terminal symbols}$

\rightarrow clearly, all alternatives of P_1 start with distinct symbols (if, while, begin)

$\rightarrow \therefore$ no problem $LL(1) \checkmark$
check

→ P2 has no choices LL(1) check ✓

→ P3 → $\text{dirsymb} (; SL) = \{ ; \}$ → $= \{ \epsilon \}$

$\text{dirsymb} (\epsilon) = \text{if } (\epsilon \in \text{first}(\epsilon))$

$(\text{first}(\epsilon) - \{ \epsilon \}) \cup \text{FOLLOW}(S')$

else $\text{first}(\epsilon)$

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

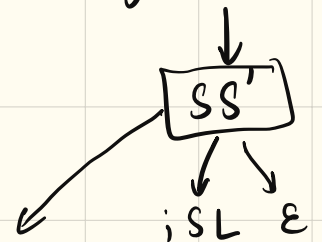
$= \{ \text{end} \}$

∴ $\text{dirsymb} (; SL) \cap \text{dirsymb} (\epsilon) = \emptyset$

No common symbols

LL(1) check ✓

$S \rightarrow \text{begin } SL \text{ (end)}$



$S' \rightarrow \text{last symbol}$

∴ $\text{FOLLOW}(S') = \text{FOLLOW}(SL)$

Hence the grammar is $LL(1)$