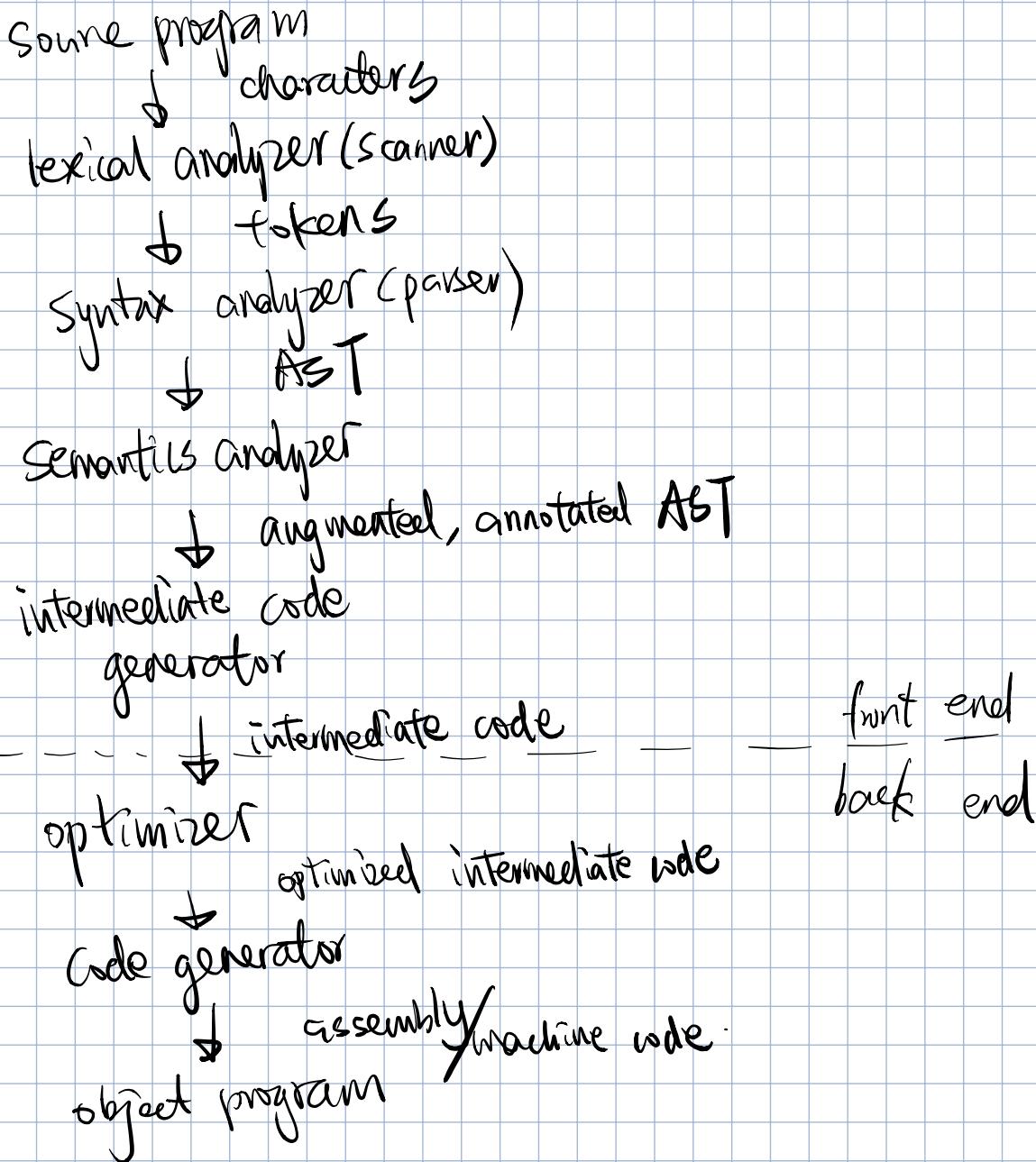


Lee I Compiler Overview



Errors to catch:

scanner

parser

Semantic
analyzer

illegal character

Syntax error

Semantic error

unterminated string

Semantic error

Type checking: wrong type

too long int literal

Name analysis: undeclared var

Errors can only be found by scanner, parser and semantic analyzer. After that, no errors would be found.

Lee 2 DFA

Scanner = sequence of chars \rightarrow sequence of tokens

Find longest sequence of chars corresponding to a token

Finite State Machine (FSM), a.k.a. finite automata

Deterministic Finite Automata (DFA)

$M = (Q, \Sigma, \delta, q_0, F)$, where

Q : finite set of states

Σ : set of characters

δ : $Q \times \Sigma \rightarrow Q$, transition functions

q_0 : start state, $q_0 \in Q$.

F : final states, $F \subseteq Q$.

Language of FSM M , $L(M)$. FSM accepts string $x_1 x_2 \dots x_n$ iff

$(\delta(\dots \delta(\delta(q_0, x_1), x_2) \dots), x_n) \in F$

Deterministic v.s. Non-deterministic

DFA : no state has > 1 outgoing edge with the same label

NFA : may have > 1 outgoing edges with same label.
edges may be labelled with ϵ

Lee 3 DFA vs. NFA ; RegEx

NFA: $M \equiv (Q, \Sigma, \delta, q_1, F)$

Q : finite set of states

Σ : set of characters

$\delta: Q \times \Sigma \rightarrow \underline{2^Q}$

q_1 : start state $q_1 \in Q$

F : final state(s) $F \subseteq Q$

For DFA: $\delta: Q \times \Sigma \rightarrow Q$

For NFA: $\delta: Q \times \Sigma \rightarrow 2^Q$. We have 2^Q as a state can have multiple outgoing edges with the same label.

Transition table

DFA

	C_1	C_2	\dots
s_1			
s_2			
\vdots			

Each entry is a single state (or stack).

NFA

	C_1	C_2	\dots
s_1	$\{s_2, s_3\}$	$\{s_1\}$	
s_2		$\{s_3, s_4, s_5\}$	
\vdots			

Each entry is a set of state(s).

For DFA M , $x_1 x_2 \dots x_n \in L(M)$ iff.

$$\delta(\dots \delta(\delta(q, x_1), x_2) \dots, x_n) \in F$$

where the output of δ is zero/one state.

For NFA M , $x_1 x_2 \dots x_n \in L(M)$ iff

$$\delta(\dots \delta(\delta(q, x_1), x_2) \dots, x_n) \cap F \neq \emptyset$$

where the output of δ is a (possibly empty) set of states.

NFA and DFA are equivalent.

Two automata M and M' are equivalent iff. $L(M) = L(M')$

Lemma 1: DFA \rightarrow NFA: given a DFA, you can construct an equivalent NFA.
Trivially correct.

Lemma 2: NFA \rightarrow DFA

Converting NFA into DFA: ① normal-edge. The new DFA tracks set of states.
② remove epsilon.

Algorithm to remove ϵ -transition:

$\text{elclose}(s)$ = set of all states reachable from s in ≥ 0 epsilon transitions.

Rule #1: To construct ϵ -free, equivalent FSM M' from M , s is an accepting state of M' iff. $\text{elclose}(s)$ contains an accepting state.

Rule #2: $s \xrightarrow{\epsilon} t$ is a transition in M' iff $q \xrightarrow{\epsilon} t$ from some $q \in \text{elclose}(s)$

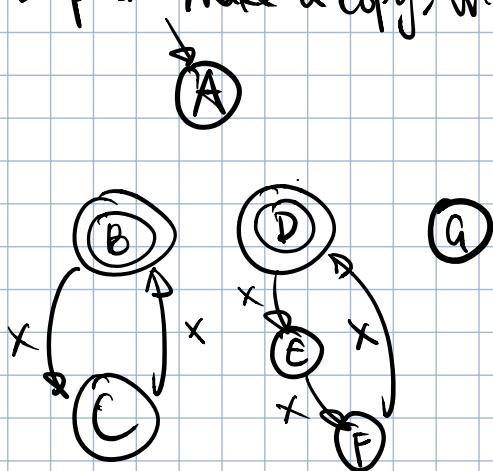
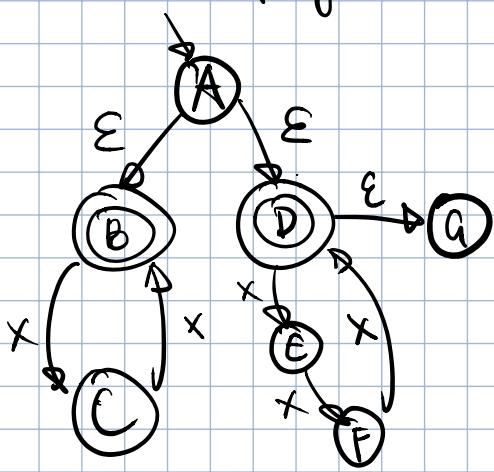
Graphically,

M (NFA): $Q_s \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} Q_q \xrightarrow{\epsilon} t, q \in \text{elclose}(s)$

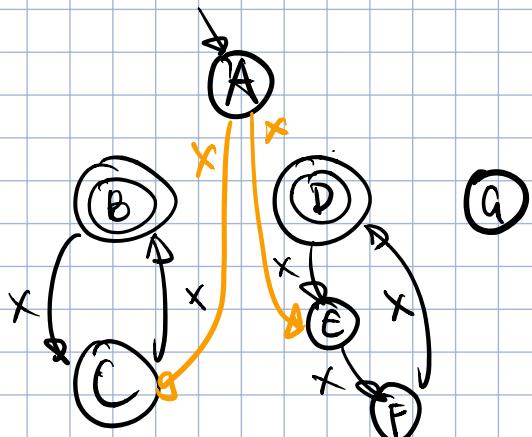
M' (DFA): $Q_s \xrightarrow{\epsilon} Q_t$

M : accepting X^n where n is even or divisible by 3.

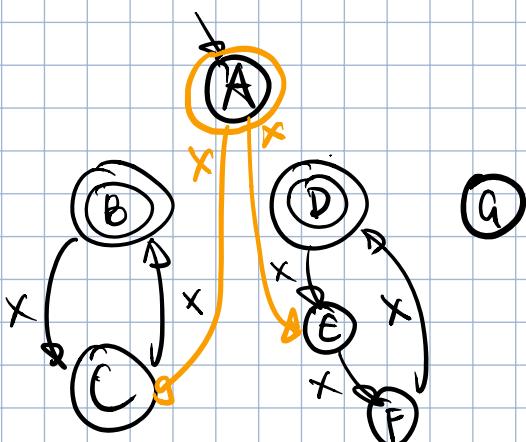
Step 1: Make a copy, with ϵ -transition removed



Step 2: Follow Rule 1



Step 3: Follow Rule 2



Conclusion: ① DFAs and NFAs are equivalent

② ϵ -transition don't add expressiveness to NFA.

Lec 4 RegExs & DFAs

Writing regexps differently doesn't affect performance: all compiled into DFAs.

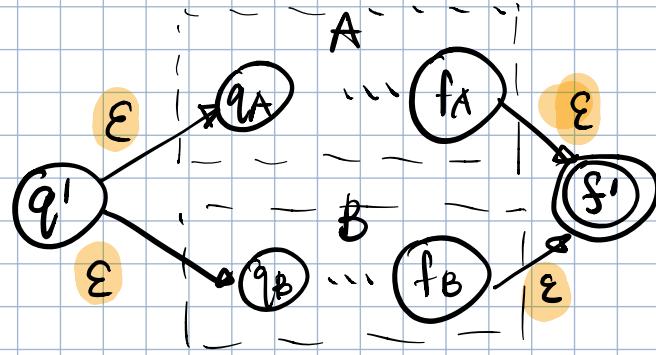
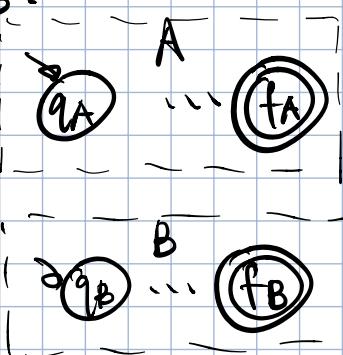
RegExps to NFA

Idea = literals/epsilon → simple, small DFAs

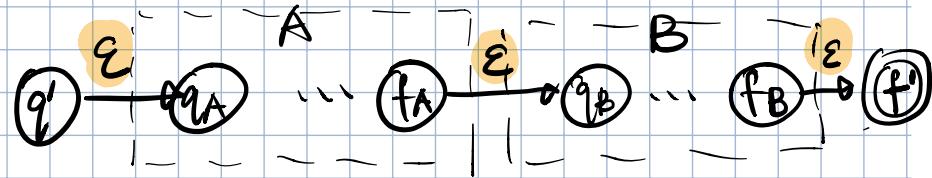
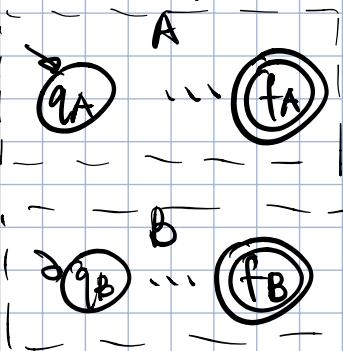
Operators → methods of joining DFAs.

Rule: each expression has a unique start state and a unique accepting state.

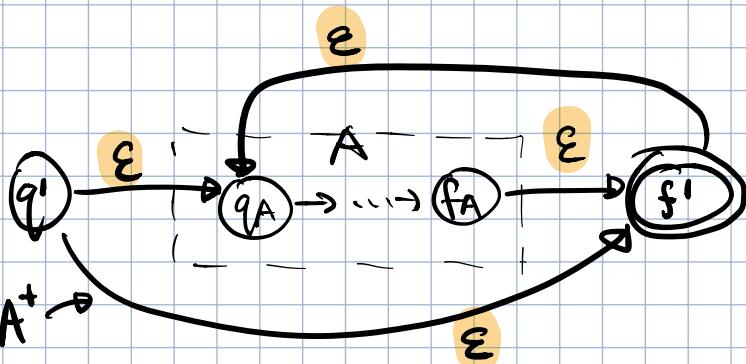
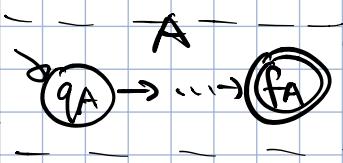
$A \mid B$:



$A \cdot B$ (concatenation):



A^* (iteration):



Remove this and you get A^+

RegExp operator precedence:

$*$ > $.$ > $|$

Regexp → tree representation → NFA w/ ϵ → NFA w/o ϵ → DFA

Essentially: $\text{RegExp} == \text{NFA} == \text{DFA}$

Why we want DFA ?

The transitions S can be expressed as a table and an efficient array representation.

Use FSM for tokenization ?

Up to now, FSM only checks language membership of a string. We need to relate the FSM with some actions .

Solution : action + ability to put char back.

Ready Construct DFA and action table

In this setting, at each run the DFA would return the longest token. Thus, the scanner would keeps running the FSM until it consumes the entire string / gets stuck.

From here (Lee 2 - Lee 4), we finished the discussion on scanner . Next topic : Parser.

Lec 5 CFG

RegExp is great for scanner(tokenization), but insufficient for parsing.

Limitations of RegExp:

① cannot specify all structs. E.g "matching": balanced parenthesis.

② doesn't enforce structure on the tokens.

$$CFG \equiv (N, \Sigma, P, S)$$

N : set of non-terminals

Σ : set of terminals

P : set of production rules

$S \in N$, the initial non-terminal symbol

Derivation syntax: \rightarrow , \rightarrow^+ , \rightarrow^*

Lec 6 Define Syntax Using CFGs w/o ambiguity

Imposing derivation order doesn't remove ambiguity.

Grammar G is ambiguous if it has > 1 parse tree for a string s .

Remove ambiguity: precedence + associativity

- Fix precedence : It's a common practice to put literals and parenthesized terms at the same precedence level.
- Fix associativity : {left recursion \rightarrow left associativity
right recursion \rightarrow right associativity}

Lec 7 SDT

Syntax Directed Translation: augment CFG rules w/ translation rules.

Suppose we already have the parse tree, then apply SDT rules bottom-up.

Lec 8 Java API

Lec9 Bottom-up parsing algorithm: CYK

Chomsky normal form (CNF)

$$X \rightarrow t$$

$$X \rightarrow A B$$

The only rule allowed to derive epsilon is the start symbol S.

CYK = build the parse tree in $O(n^3)$, bottom-up

Useless Grammar: ① A non-terminal that cannot derive terminals.

② A non-terminal that can't be derived from the start symbol

Convert grammar into CNF: 4 steps

① eliminate ϵ rules

$$F \rightarrow id(A) \quad F \rightarrow id()$$

$$\begin{array}{l} A \rightarrow \epsilon \\ A \rightarrow N \end{array} \Rightarrow \begin{array}{l} | id(A) \\ | id(N) \end{array}$$

$$\begin{array}{l} X \rightarrow A \times A \\ | \\ A \rightarrow \epsilon \\ | \\ A \rightarrow \epsilon \end{array} \Rightarrow \begin{array}{l} X \rightarrow A \times A \\ X \rightarrow \times A \\ X \rightarrow A \times \\ X \rightarrow \times \\ A \rightarrow \epsilon \end{array}$$

② eliminate unit rules

$$\begin{array}{l} F \rightarrow id() \\ | id(A) \\ A \rightarrow N \end{array} \Rightarrow \begin{array}{l} F \rightarrow id() \\ | id(N) \end{array}$$

③ Fix productions with terminals (and other things) on RHS

$$\begin{array}{l} F \rightarrow id() \\ | id(N) \end{array} \Rightarrow \begin{array}{l} F \rightarrow I L R \\ | I L N R \\ I \rightarrow id \\ L \rightarrow (\\ R \rightarrow) \end{array}$$

Rename each terminal
as an non-terminal.

④ Fix production with > 2 nonterminals on RHS.

$$A \rightarrow B C D \Rightarrow A \rightarrow B E \\ \qquad \qquad \qquad E \rightarrow C D$$

Lec 10 Top-down parsing for LL(1) (Recursive Descent)

Restricting the grammar to be LL(1), O(n) top-down parser.

Algorithm:

```
stack.push(eof)
stack.push(Start non-term)
t = scanner.getToken()
Repeat
    if stack.top is a terminal y
        match y with t
        pop y from the stack
        t = scanner.next_token()
    if stack.top is a nonterminal X
        get table[X,t]
        pop X from the stack
        push production's RHS in reverse order
Until one of the following:
    stack is empty accept
    stack.top = terminal that doesn't match t
    stack.top = non-term and parse table entry is empty
reject
```

Necessary (not sufficient) conditions for LL(1):

- ① No left recursion
- ② Left factored

If left-recursion exists, no prediction can be guaranteed to be correct.
If not left factored, obvious ambiguity.

• Remove left-recursion:

$$A \rightarrow A\alpha | \beta \Rightarrow A \rightarrow \beta A' \\ A' \rightarrow \alpha A' | \epsilon$$

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n \quad | \quad \beta_1 | \beta_2 | \dots | \beta_m \Rightarrow A \rightarrow \beta_1 | \beta_2 | \dots | \beta_m A' \\ A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon$$

• Left factoring

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 \Rightarrow A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 | \beta_2$$

$$A \rightarrow \alpha\beta_1 | \dots | \alpha\beta_m | y_1 | \dots | y_n \Rightarrow A \rightarrow \alpha A' | y_1 | \dots | y_n \\ A' \rightarrow \beta_1 | \dots | \beta_m$$

Lee 11 Build the Parse Table

$FIRST(\alpha)$: the set of terminals that begin the strings derivable from α .
 If α can derive ϵ , then ϵ is also in $FIRST(\alpha)$.

$$FIRST(\alpha) = \{ t \mid (t \in \Sigma \wedge \alpha \xrightarrow{*} t\beta) \vee (t = \epsilon \wedge \alpha \xrightarrow{*} \epsilon) \}$$

If $\alpha \in \Sigma$, $FIRST(\alpha) = \{ \alpha \}$

If $\alpha \in N$, Let $\alpha = Y_1 Y_2 \dots Y_k$

- Put $FIRST(Y_1) - \{\epsilon\}$ in $FIRST(\alpha)$
- If ϵ is in $FIRST(Y_1)$: add $FIRST(Y_2) - \{\epsilon\}$ to $FIRST(\alpha)$
- If ϵ is in $FIRST(Y_2)$: add $FIRST(Y_3) - \{\epsilon\}$ to $FIRST(\alpha)$
- ...
- If ϵ is in $FIRST$ of all Y_i symbols, put ϵ into $FIRST(\alpha)$

$FOLLOW(\alpha)$: For non-terminal α , $FOLLOW(\alpha)$ is the set of terminals that can appear immediately to the right of α .

$$FOLLOW(A) = \{ t \mid (t \in \Sigma \wedge S \xrightarrow{*} \alpha A \beta) \vee (t = \text{eof} \wedge S \xrightarrow{*} \alpha A) \}$$

$FOLLOW(A)$ for $X \rightarrow \alpha A \beta$

If A is the start, add **eof**

Add $FIRST(\beta) - \{\epsilon\}$

Add $FOLLOW(X)$ if ϵ in $FIRST(\beta)$ or β is empty

When building $FOLLOW(A)$, look at every production rule whose RHS contains A .

Build the parse table

```

for each production  $X \rightarrow \alpha$  {
    for each terminal  $t$  in  $FIRST(\alpha)$  {
        put  $\alpha$  in Table[X][ $t$ ]
    }
    if  $\epsilon$  is in  $FIRST(\alpha)$  {
        for each terminal  $t$  in  $FOLLOW(X)$  {
            put  $\alpha$  in Table[X][ $t$ ]
        }
    }
}
  
```

Iterating over each production rule

Whole process:

- ① Build $FIRST$ for all nonterminals
- ② Build $FIRST$ for all RHTs
- ③ Build $FOLLOW$ for all nonterminals
- ④ Fill the parse table

FIRST(α) for $\alpha = Y_1 Y_2 \dots Y_k$

Add FIRST(Y_1) - { ϵ }

If ϵ is in FIRST($Y_{1 \text{ to } i-1}$): add FIRST(Y_i) - { ϵ }

If ϵ is in all RHS symbols, add ϵ

FOLLOW(A) for $X \rightarrow \alpha A \beta$

If A is the start, add eof

Add FIRST(β) - { ϵ }

Add FOLLOW(X) if ϵ in FIRST(β) or β empty

Table[X][t] "S → Bc | DB" should be considered as $\begin{cases} S \rightarrow Bc \\ S \rightarrow DB \end{cases}$

for each production $X \rightarrow \alpha$

for each terminal t in FIRST(α)

put α in Table[X][t]

if ϵ is in FIRST(α) {

for each terminal t in FOLLOW(X) {

put α in Table[X][t]

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

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

$$\text{FIRST}(D) = \{d, \epsilon\}$$

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

$$\text{FIRST}(DB) = \{d, a, c\}$$

$$\text{FIRST}(ab) = \{a\}$$

$$\text{FIRST}(cS) = \{c\}$$

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

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

$$\text{FOLLOW}(S) = \{\text{eof}, c\}$$

$$\text{FOLLOW}(B) = \{c, \text{eof}\}$$

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

CFG

$$S \rightarrow Bc | DB$$

$$B \rightarrow ab | cS$$

$$D \rightarrow d | \epsilon$$

Not LL(1)

	a	b	c	d	eof
S	Bc DB		Bc DB		
B	ab		cS		
D	ϵ		ϵ	d	

Table collision \Leftrightarrow Not LL(1)

Now we can use recursive descent, instead of CYK, to build the parse tree for LL(1) grammar.

Lee 12 SDT for Top-Down Parsing

- Last time: remove left-recursion, perform left factoring to get LL(1). However, grammar transformation affect the structure of the parse tree.
- So far, SDT shown as a procedure performed on parse-tree. But for LL(1), the parser never needs to explicitly build the parse tree. Instead, the LL(1) parser uses a **semantic stack** to implicitly track the state. Also, SDT rules become actions (popping, pushing) on the semantic stack.
- Number actions and put them on the symbol stack by adding action number symbols at the end of the production.
- Action numbers go **AFTER** the corresponding nonterminals, **BEFORE** their corresponding terminal. (We put action number before the corresponding terminal so as to consume its value.)

<u>CFG</u>	<u>SDT Actions</u>
$\text{Expr} \rightarrow \text{Expr} + \text{Term} \ #1$	#1 tTrans = pop ; eTrans = pop ; push(tTrans + eTrans)
 Term	
$\text{Term} \rightarrow \text{Term} * \text{Factor} \ #2$	#2 fTrans = pop; tTrans = pop ; push(fTrans * tTrans)
 Factor	
$\text{Factor} \rightarrow \#3 \text{ intlit}$	#3 push(intlit.value) ↑

The order matters! In reverse order.

Benefits of using Action Numbers: grammar transformation doesn't affect parse tree
Action numbers can go into the parse table.

Now we finished the topic of parser (of LL(1)).

Next topic: Semantic analyzer

Lee 13 Semantic Analysis with Emphasis on Name Analysis.

Static semantic analysis:

① Name analysis

process declarations, add to symbol table

For each scope, of process statements, update IDs to point to entries

② Type analysis

process statements: use symbol table info to determine type of each exp

Practically feasible to catch errors:

undeclared identifiers;

multiple declared identifiers,

ill-typed terms.

Static scope v.s. Dynamic scope

compile time

runtime