

# Compiler construction

Mr.	Raphael	Finkel
Dr.	Rafi	Goldstein
Prof.	דודו	
—		

---

user (client)

specification.

implementation

---

programmer

compiler

[syntax  
semantics]

---

how compiler fits into a chain of operations

compiler outputs

pure machine code (Linux Kernel)

→ augmented machine code

libraries

operating system calls. (common)

→ virtual machine code (Java) (CSX-go)

To represent the output

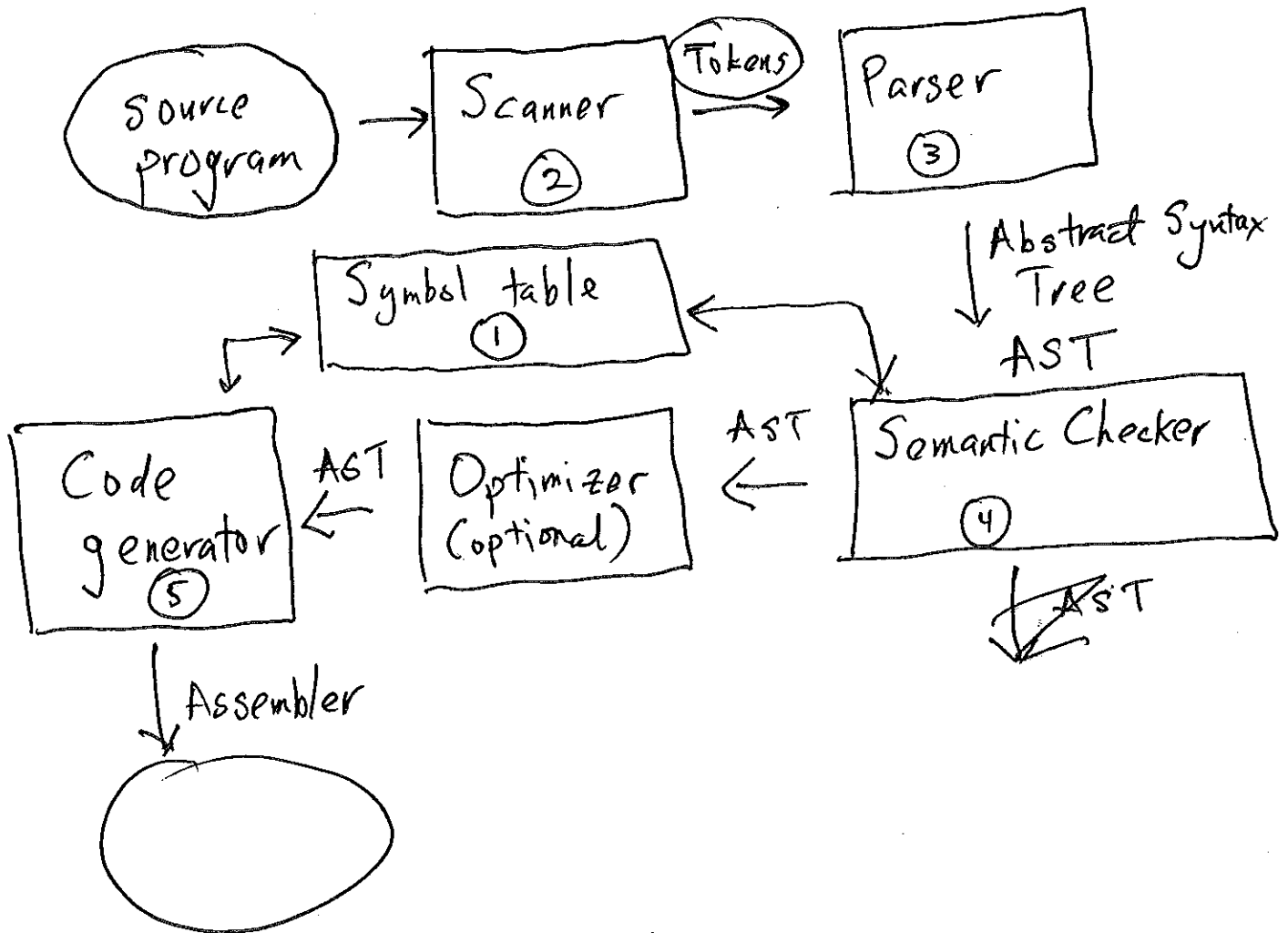
→ assembler (compiler does not need to resolve all references)

(modular compilation)

relocatable binary (Unix: .o file)

(defers resolving external references)  
absolute binary

### Structure of the compiler



# Scanner

builds a stream of tokens

reads chars

Example: <sup>input</sup> if ( a < 39 ) {

Output:

if	(ident, reserved)
"("	(punctuation)
a	(ident)
<	(punctuation)
39	(integer literal)
)	(punctuation)
{	(punctuation)

Method : regular expressions.

Scanner generator :

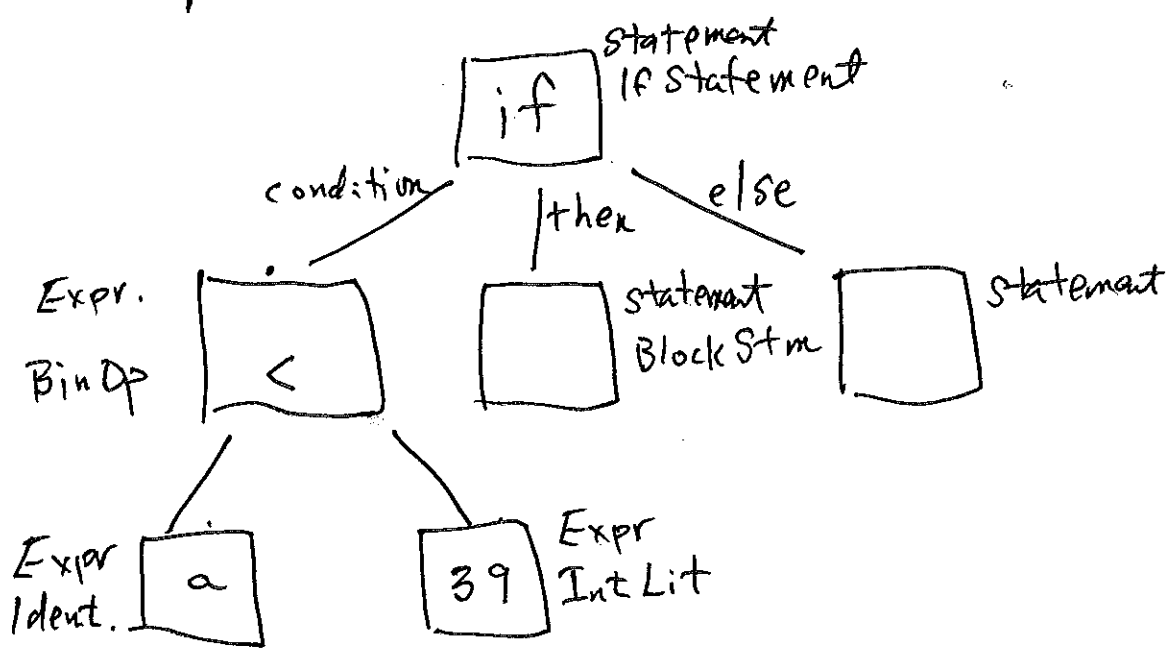
lex, flex, jflex

# Parser

reads token stream

builds AST

example: if (a < 39) {



Method: Context-free grammar.

Parser generator:

yacc, bison, javaCup

# Semantics Checker

Navigates through entire AST (DFS)

verifying declarations.

binding identifiers to ST (Symbol Table)

verify types of expressions,

adds info to AST. (type)

generate "error type" info.

add nodes to AST. (type conversion)

verify exception handling.

verify reachability.

# Optimizer

Analyzes AST or generated code to

improve it (time, space)

many techniques.

simplify expressions.

moving code

eliminating trivial arithmetic.

# Code generator

DFS of AST

generates output

assembler

# Programming Lg Considerations

(6)

## Special features

passing parameter "by name"

dynamic-sized arrays

nested name scopes. (static chain)

anonymous functions, first-class functions.

multiple-yield iterators.

automatic reclamation of object store  
(garbage collection)

## Computer architecture considerations.

limited registers.

cost of operations.

effect of memory hierarchy.

## Specialty compilers

Debugging support

Participation in an IDE

~~Retargeting~~

Retargetable compilers.

---

①  $x_2 \oplus y_1 = x_1 \wedge y_1$

②  $y_2 = x_2 \wedge y_1 = x_1 \wedge \overbrace{y_1 \wedge y_1}^0 = x_1$

③  $x_3 = x_2 \wedge y_2 = x_1 \wedge y_1 \wedge x_1 = x_1 \wedge x_1 \wedge y_1 = y_1$

The AC (adding calculator) language

types: int, float

Keywords: f i p

variables: lowercase a...z  
except f, i, p

Program:

f b i a a = 5 b = a + 3.2 p b \$

Syntax in: Context-free grammar.

formalism: Backus-Naur Form (BNF)

Prog  $\rightarrow$  Decls Stms \$

Decls  $\rightarrow$  Decl Decls

Decls  $\rightarrow$   $\lambda$

Decl  $\rightarrow$  float decl id | int decl id

Stms  $\rightarrow$  Stm Stms |  $\lambda$

Stm  $\rightarrow$  id assign Val Expr

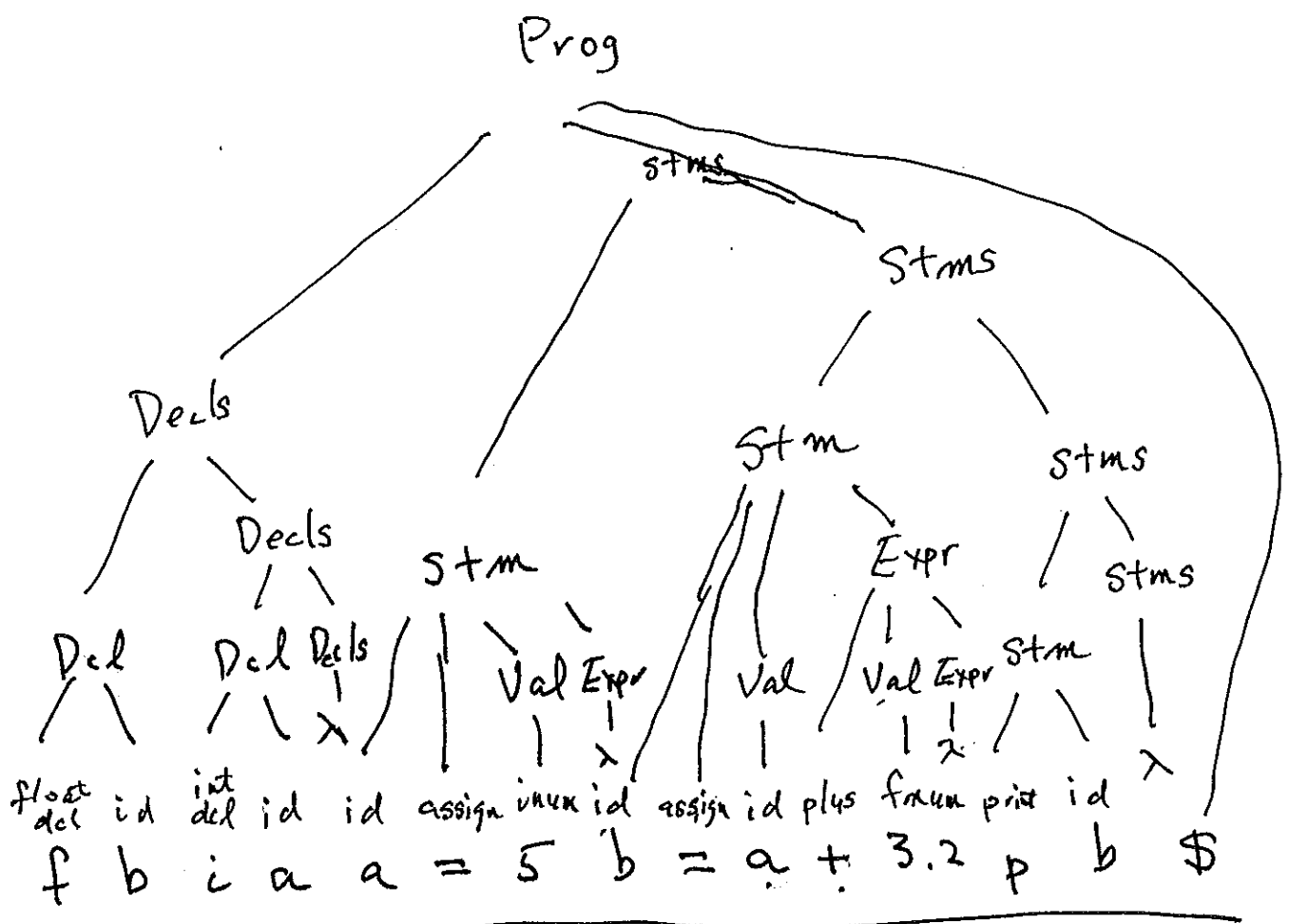
~~id~~ | print id

Val  $\rightarrow$  id | inum | fnum

Expr  $\rightarrow$  plus Val Expr | minus Val Expr |  $\lambda$

rule, production  
LHS: a nonterminal  
RHS: nts + ts  
any order  
\$ is terminal

# Parse Tree



Scanner: translates a stream of chars into a stream of tokens

type(~~operator~~, inum...)   
 plus   
 value ("3.2")

design choices:   
 are all operators separate token types? (yes)   
 are all reserved words " (yes)

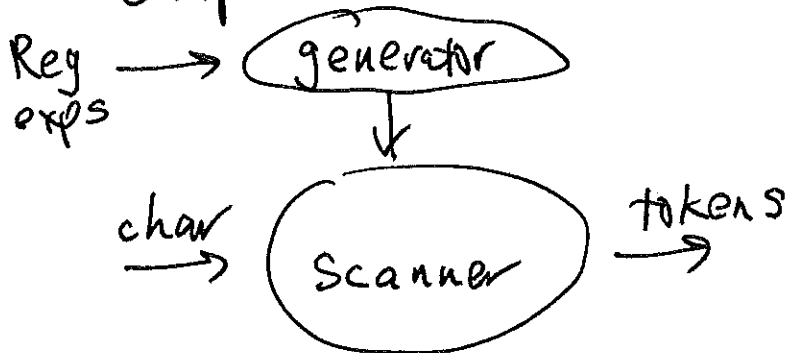


To build a scanner

- 1) Hard-coded (specific to lg)
- 2) Scanner generator (lex, flex, jflex)

Input: regular expressions.

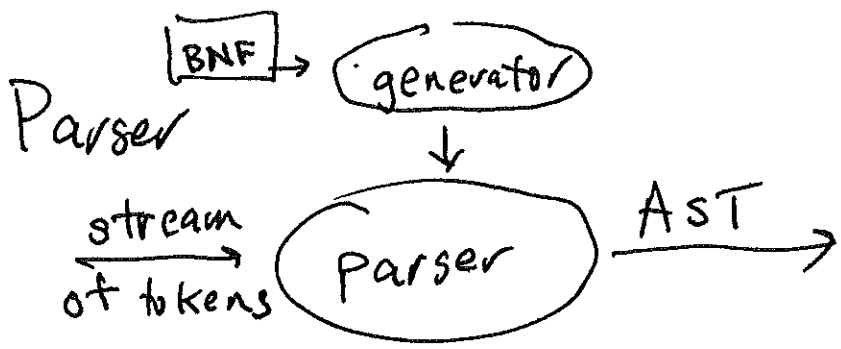
Output: Scanner



### Formal lg hierarchy

Language type	Formalism	Automaton
* Regular	Regular expression	FSA (finite-state automaton)
* Context-free	BNF CFG   context-free grammar	PDA (push-down automaton)
Context-sensitive	CSG	LBA (linear-bounded automaton)
Type 0	various	Turing machine

$S \rightarrow aSb$   
 $S \rightarrow \lambda$   
 $(S)$   
 $(\lambda)$   
 aabb



methods:

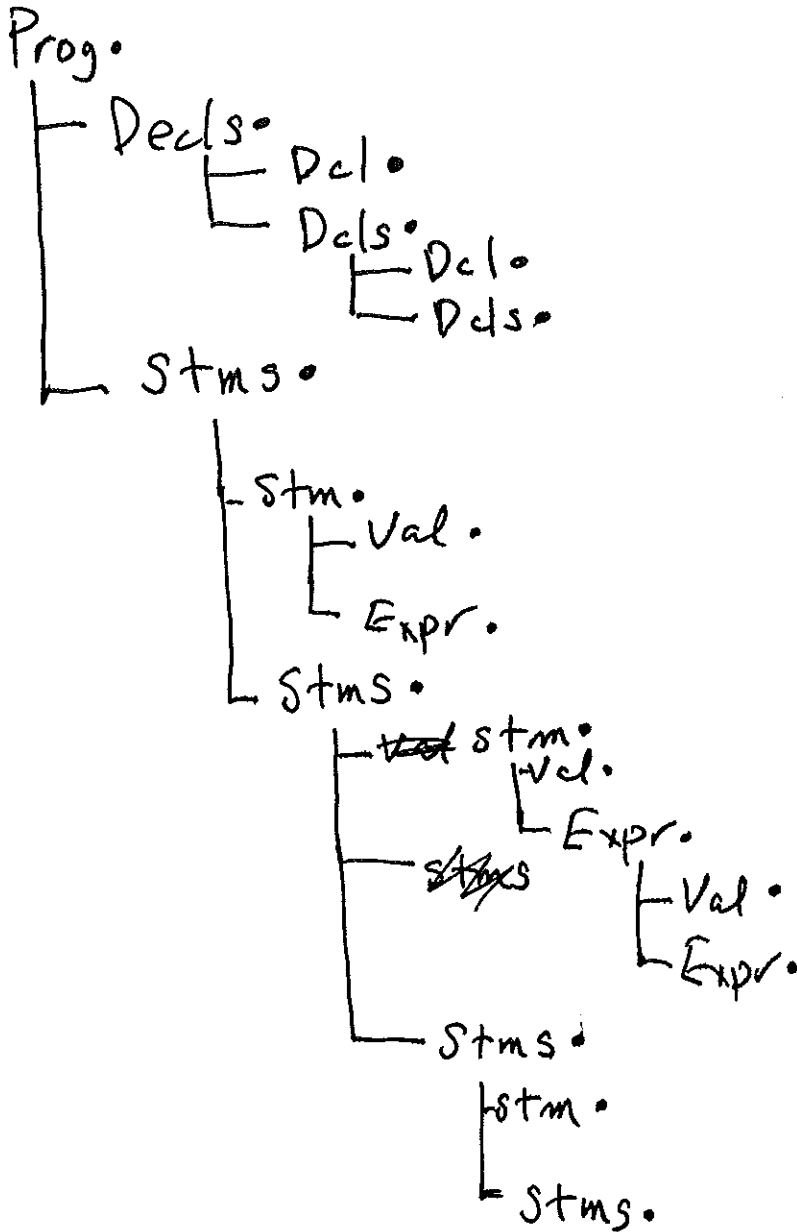
- 1) recursive descent (top-down)
  - every non-terminal (in BNF) is represented by a procedure.
  - case for each RHS for that non-terminal
  - match() for each terminal on RHS.
  - (recursive) call to a procedure for each non-terminal on RHS.

to pick which case:  
 "predict sets"

- 2) LR(0) LR(1) SALR(1) (bottom-up)
- 3) Parser generator (yacc, bison, javaCup, ...)

f b i a a = 5 b = a + 3.2 p b \$

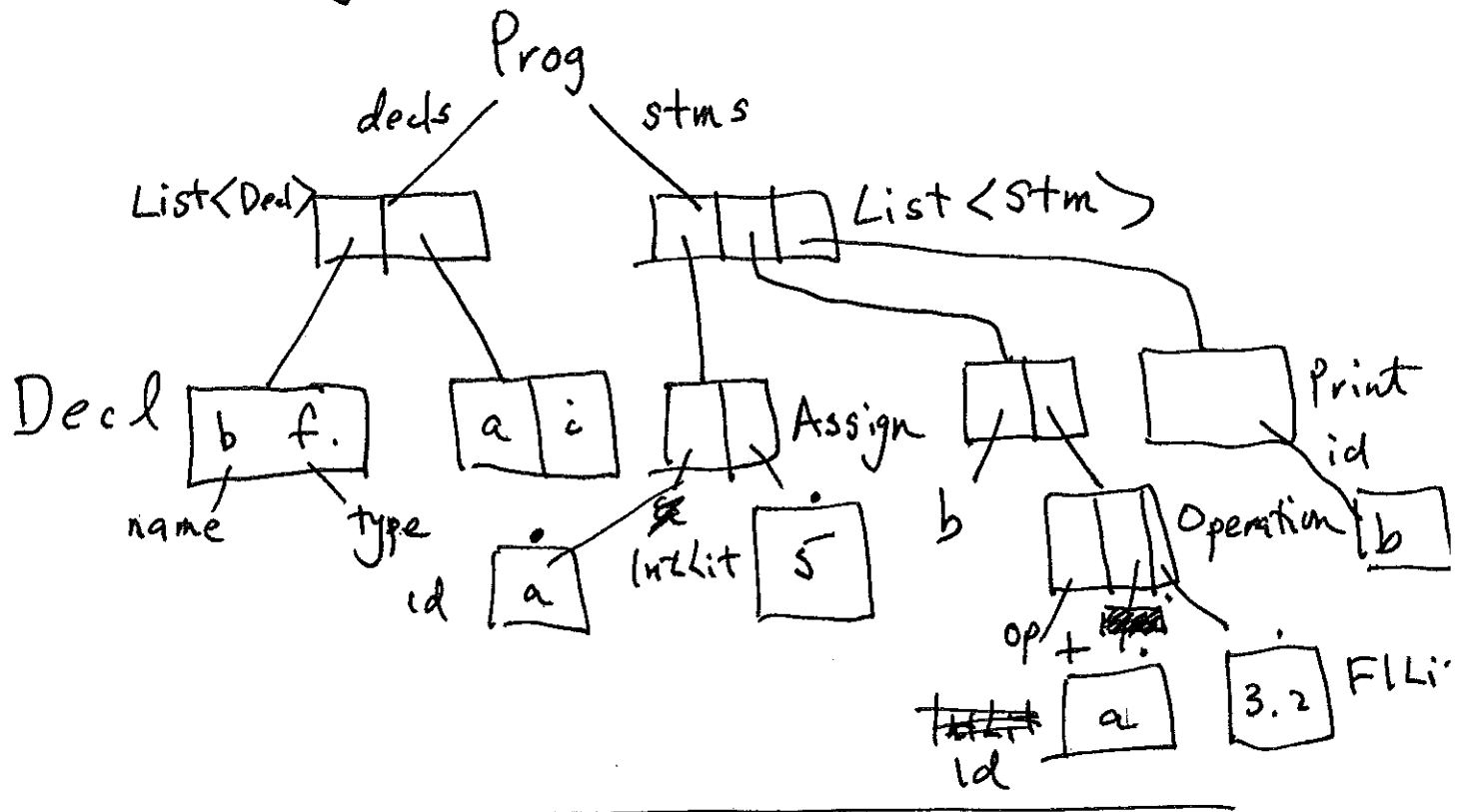
(11)



Example:

Free Prog(~~ts~~ Stream <Token> ts)

# Abstract Syntax Tree



Semantic checking: recursively visits AST  
 Uses symbol table

at declaration: insert: id, type  
 verifying uniqueness.

at use:  
 verify existence  
 mark type

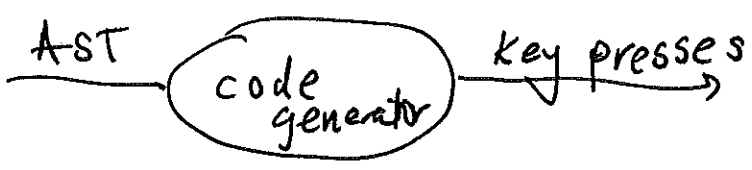
why: there are errors that scanning,  
 parsing cannot catch.

- a) strong typing.
- b) disambiguate some meanings

Java: x.y.z

- c) resolve overloaded operators.

# Generating code



## Calculator:

has registers, named by letters.

load (L) store (S)

precision: K button.

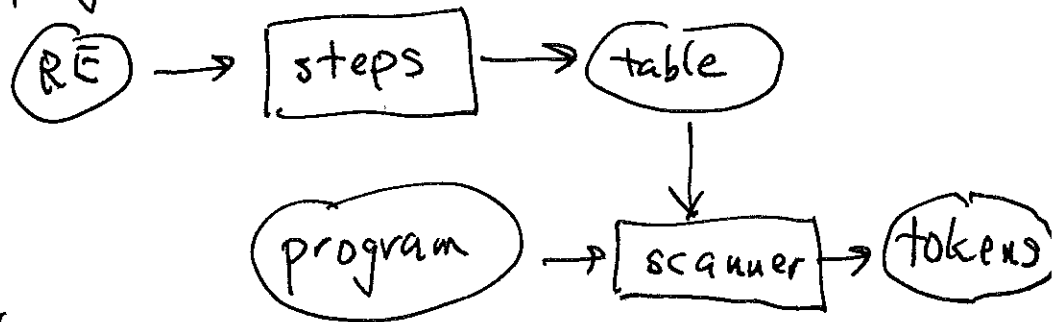
print: P button.

Method: Visit AST depth-first  
calling codeGen()

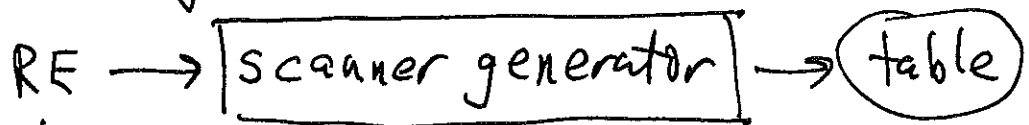
Scanner: full story  
formal, systematic

short story:

- 1) tokens are defined by regular expressions RE
- 2) which are encoded as non-deterministic finite-state automata (NFA)
- 3) which <sup>are</sup> converted to deterministic finite-state automata (DFA)
- 4) which are described by tables  
state x input  $\Rightarrow$  action x state
- 5) which are executed by a simple program.



shorter story:



Complexities

- 1) string literals: can contain escaped "
- 2) over-eagerness  
3..4 (Pascal)  
'a (Ada)  
DO 200 I = 1 . 10
- 3) speed.

# Regular expressions

RE defines a language, which is the set of valid strings over an alphabet  $\Sigma$   
↓  
valid input characters.

- 1)  $\emptyset$  [no strings at all]
- 2)  $\lambda$  [the empty string]
- 3) any letter in  $\Sigma$  [the string containing just that letter]
- 4) concatenation of REs. [all 2-part strings]
- 5) alternation of REs [all of each]  
 (notation:  $|$  (bar symbol) )  
 (example:  $a | b a \lambda$  )
- 6) closure operations [ \* : 0 or more concatenation  
 + : 1 or more ]  
 (notation:  $*$  + )  
 (example:  $a b^* c$  )
- 7) parentheses for grouping
- 8) escaped metacharacters  
 (example:  $a b \backslash * c$  )

### Examples of REs.

// Java comment : // (not(eol))\*eol

Decimal literal :

$(0|1|2|3|4|5|6|7|8|9)^+ \cdot (0|1|2|3|4|5|6|7|8|9)^+$

abbrev:  $D = (0|1|\dots|9)$

$D^+ \cdot D^+$

Integer literal :

~~$-^* D^*$~~

$(-|\lambda) D^+$

$(-|+|\lambda) D^+$

Comment delimited by ##

~~$## (\Sigma)^* ##$~~

~~$## (not(##))^* ##$~~

$## ((\#|\lambda) not(\#))^* ##$

question: ## a ### not in set

Balanced braces

example:  $\{ \{ \} \} \{ \} \}$  not a regular set



Fortran real literal:

need at least one digit on at least one side of the decimal point.

$$D^+ \cdot D^* \mid \cdot D^+$$

Identifier, including underscores (  ) but:  
never adjacent, frontal, or terminal.  
not starting with digit.

$$L (LID)^* (\_ (LID)^+)^*$$

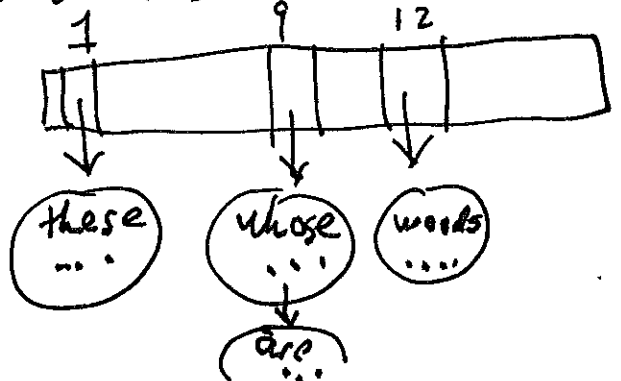
$$L ((\_ | \_ ) (LID))^*$$

Hashing: Key  $\rightarrow$  value (lookup function)  
string  $\rightarrow$  ~~int~~ def.  $\Theta(1)$

method: introduce a function from keys to integer indices.

Build an array indexed by those indices containing the "values"

- $h(\text{"whose"}) \rightarrow 9$
- $h(\text{"woods"}) \rightarrow 12$
- $h(\text{"these"}) \rightarrow 1$
- $h(\text{"are"}) \rightarrow 9$



Difference between Map and HashMap? (18)

implements  
 ↙ ↘  
 Map      and      HashMap?  
 ↑                      ↑  
 Interface              class

~~Map <String, Symb> hm~~

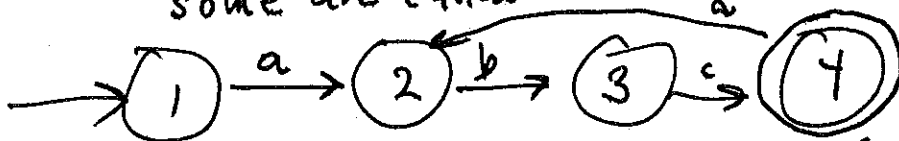
Map <String, Symb> hm =  
new HashMap <String, Symb>;

List <int> li =  
new LinkedList <int>;

RE → automata

FSA: finite-state automaton  
 model of a simple computer.

Set of states (depicted by circles)  
 One is called "start state" (arrow from nowhere)  
 some are called "final states" (double circle)



transitions between states (labelled arrow)

abc | abc abc | ...  
 (abc)<sup>+</sup>                      member of Σ

## Equivalent concepts

(1) Regular expression (RE)

Regular language

(3) Finite-state automaton (deterministic)

(2) Non-deterministic finite-state automaton

(4) table: rows are states  
columns are inputs  
cells are state numbers

	a	b	c	
$\emptyset$ 1	2	5	5	
2	5	3	5	
3	5	5	4	
4	2	5	5	final
5	5	5	5	

generate : produce sentences in lg.

recognize : accept only sentences in lg.

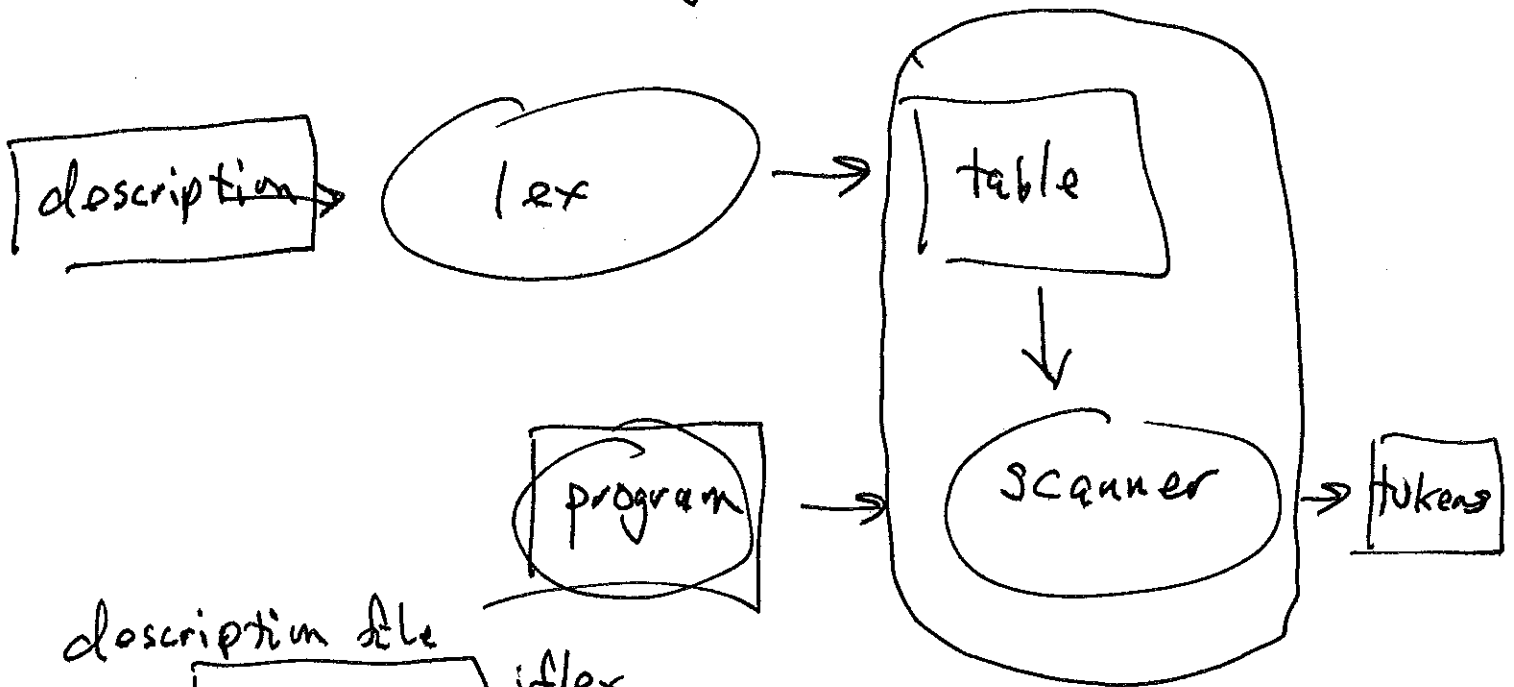
transduce : take sentences in lg,  
create new sentences in  
other lg.

FSA is deterministic if:

- 1) no conflicting transitions
- and
- 2) no  $\lambda$  transitions

else: non-deterministic

lex flex jflex .....



description file

Java code  
 % %  
 shortcuts  
 % %  
 rules.

jflex

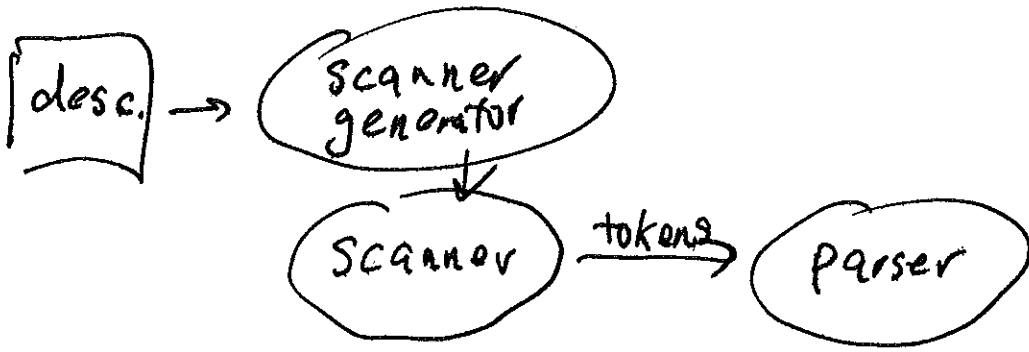
% ignorecase  
character classes

$D = [0-9]$

$L = [a-z]$  or  $[a-zA-Z]$

Letter =  $[a-e] | [g-l]$

~~E~~ (f|F)(o|O)(r|R) ~~[j-o]~~ [j-o] | [g-z]



shared token codes: # small integers  
 usually: table built by parser generator  
 yacc: y.tab.h  
 javaCup: sym.java

### Regular expressions for lex

bracket syntax for a range of values

[a-i] "character classes"

escape syntax for metacharacter

$\backslash$  [      $\backslash$  +      $\backslash \backslash$   
 "[     "+"

no need to quote alphanumeric characters

case is significant

[pP][rR][iI][mM][tT]

% ignore case

metacharacters: \* + | ( ) . ?

^ : start of an input line

\$ : end of an input line

processing code: in braces; opening brace  
on same line as its RE.

Where is the text that was matched?

lex: string variable: yytext  
jflex: method: yytext()

evanescent!

If two REs overlap:

longest match wins.

matching same string: earlier rule wins.

Catch-all RE:

- error rule to match any character  
(at end)

End-of-input: token with code 0.

~~Weird~~

Weird situations

Pascal 31..42

Back up to last accepting state

Ada foo'bar

Right context:  $[0-9]^+ / \dots$

Alternatives to lex (flex, jflex, ...)

GLA, re2c: generate direct code

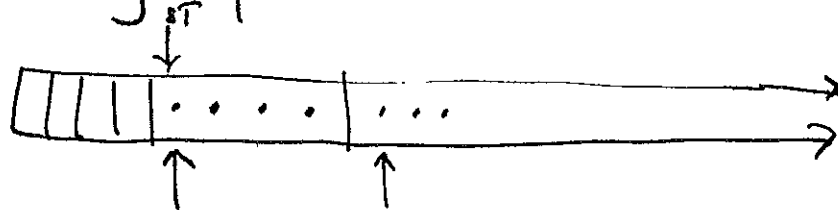
# Practical considerations

23

Identifiers: put in ST or not?

LA lg not block structured, yes.

String space



next string  
begins here

If case is significant:  
don't convert case.

If case not significant  
convert all identifiers to lower.

## Literals

Integers: convert to internal (binary) representation.  
out of range?

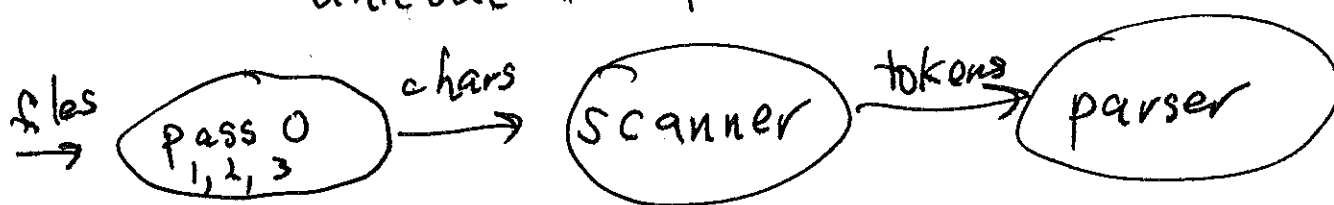
String: expand escaped characters  
also store original form.

Reserved words (for, while, ...)   
1) each one associated with a different token type.

2) treat like identifiers, but look them up in a special table.

## Compiler directives

- 1) file inclusion: stack of open files (recursion)
- 2) conditional compilation: (for debugging, to customize code for specific environments) previous step before the scanner (pass 0)
- 3) Unicode escapes



End of file: generate an EOF token.

Multi-character lookahead

when scanner enters an error state, back up, returning characters until it enters an accepting state. If none, announce "error", skip one character.

Speed: use scanner generator  
buffer reads  
use a profiling tool



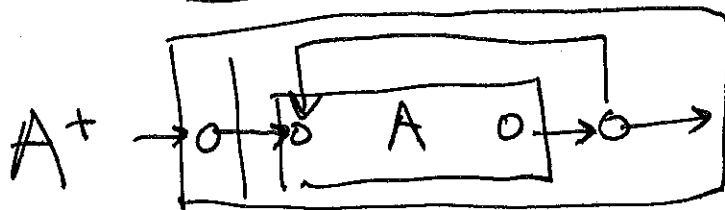
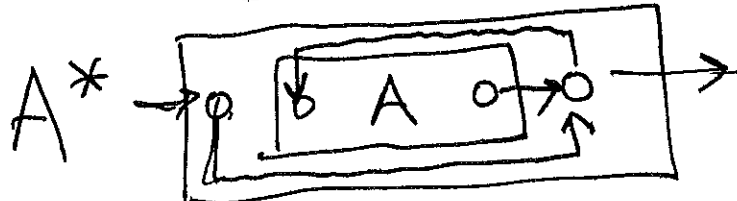
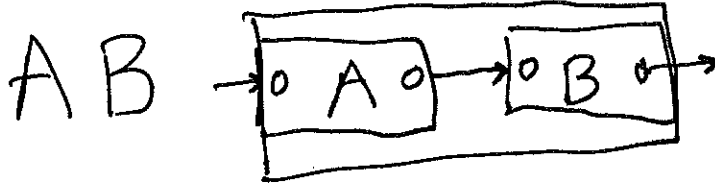
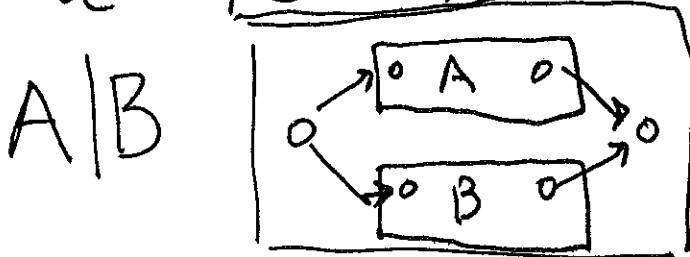
# Error recovery

generate an "error" token.

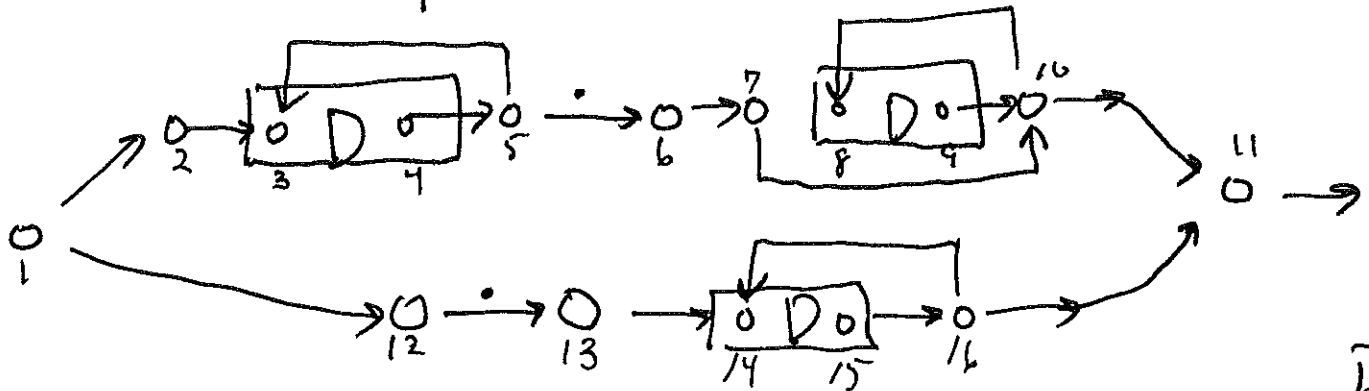
runaway strings and comments.

might disallow crossing line boundary.

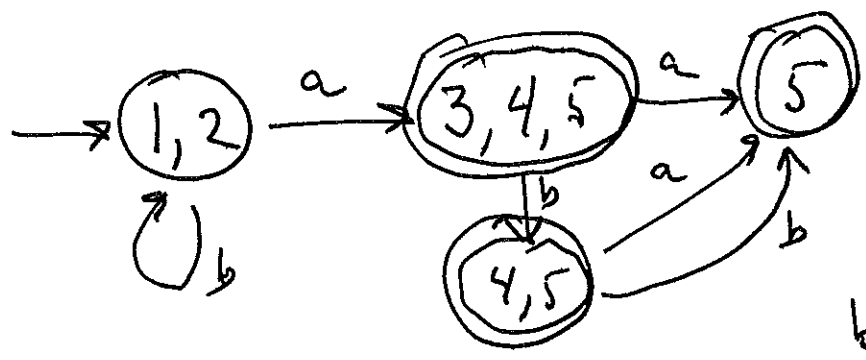
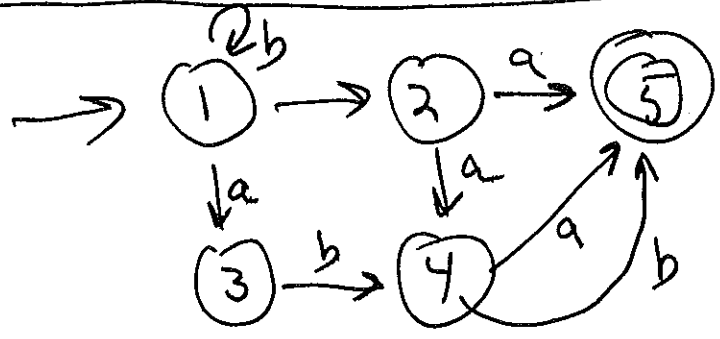
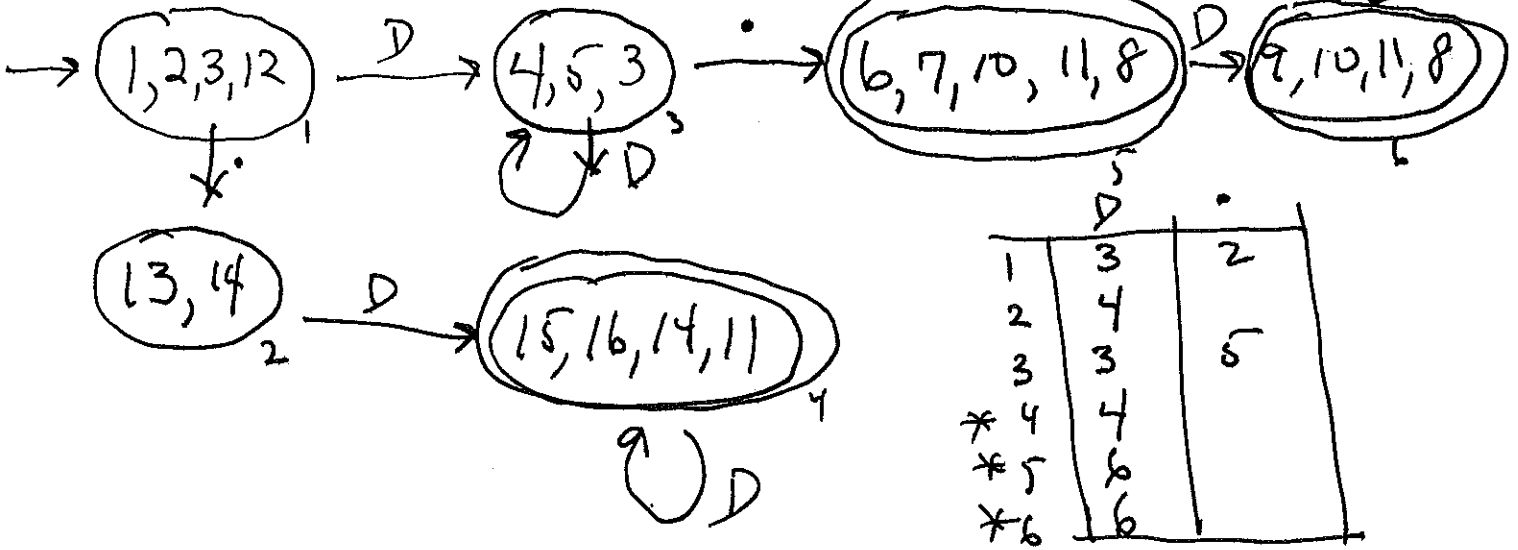
RE  $\Rightarrow$  NDFA  $\Rightarrow$  FSA



$D^+ \cdot D^* \mid \cdot D^+$

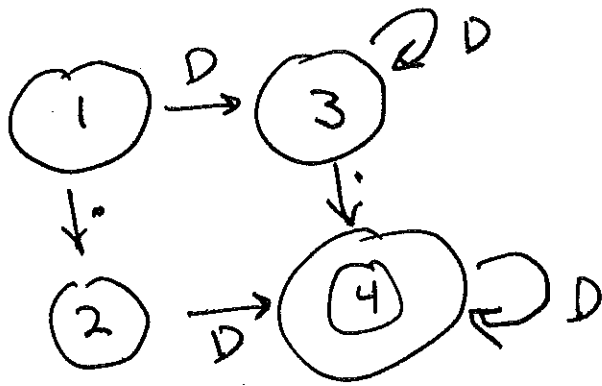


subset construction



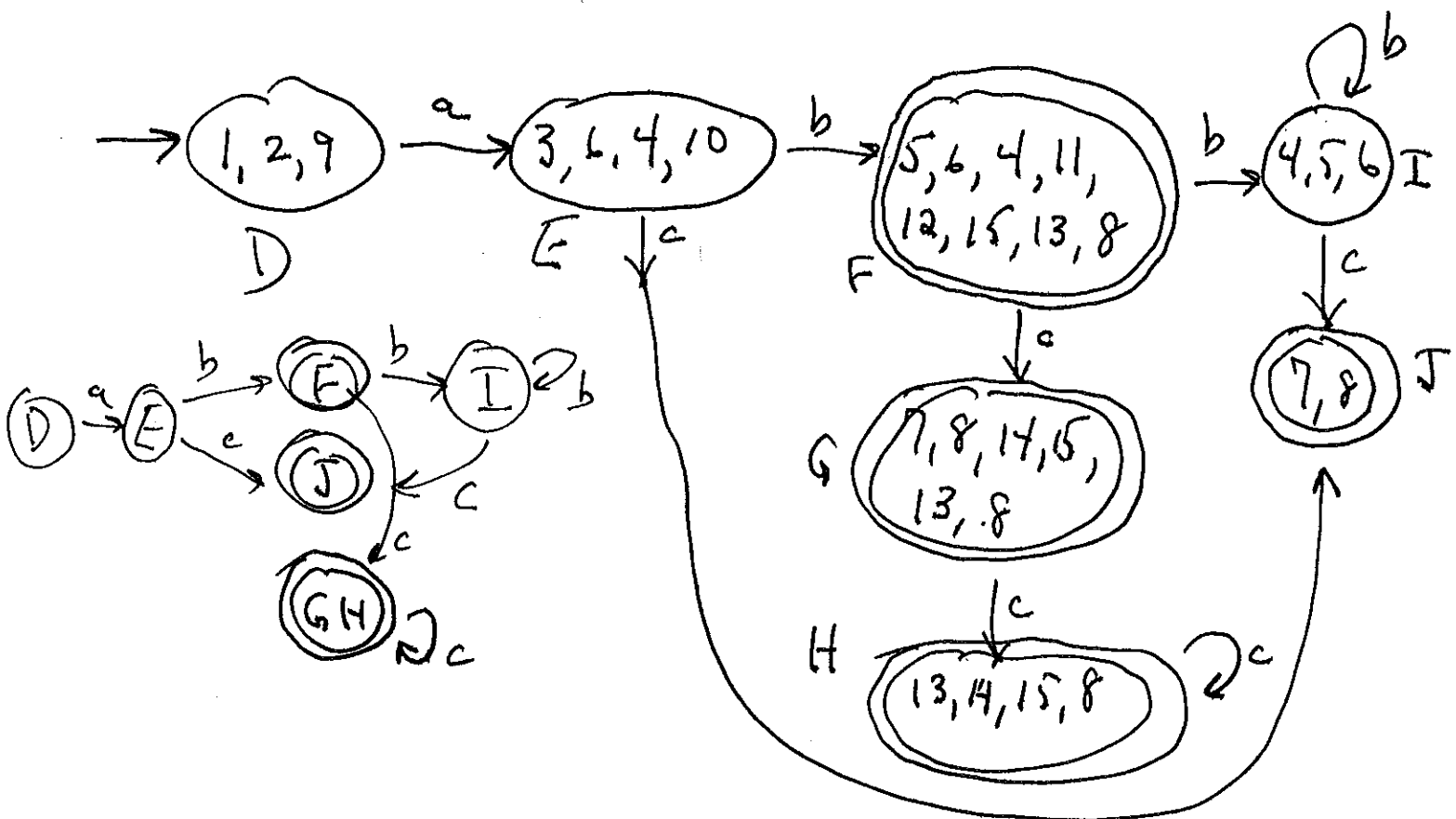
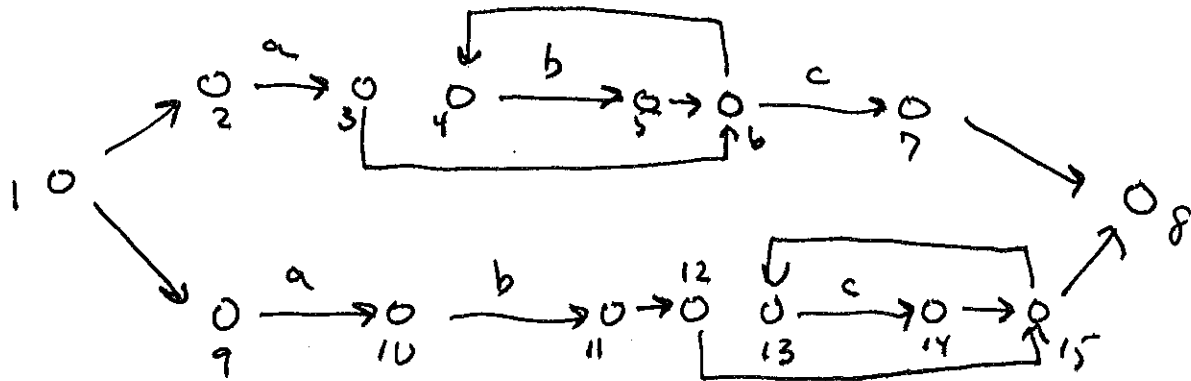
$\frac{1}{b} a (a/b)$

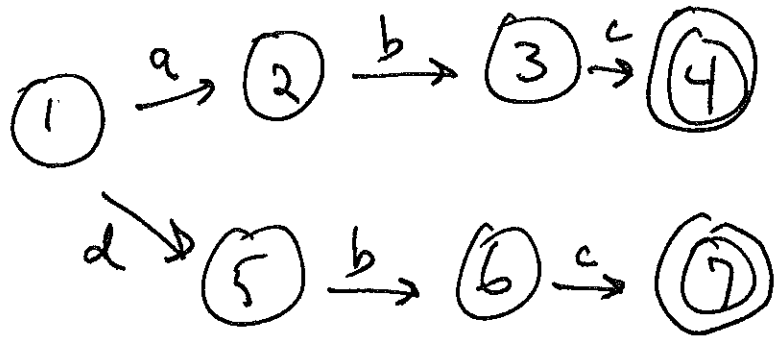
$b^* a (\lambda | a | b (\lambda | a | b))$



	D	0
1	3	2
2	4	
3	3	4
*4	4	

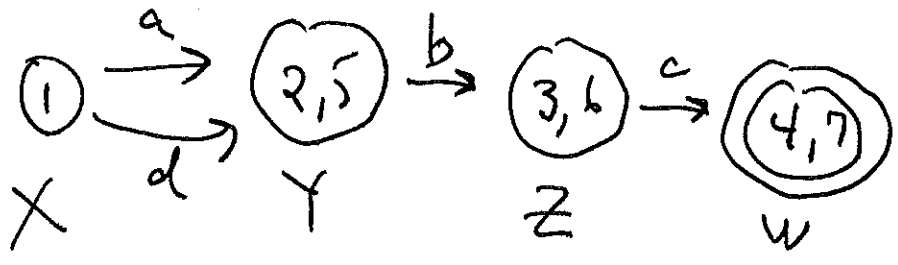
RE:  $(ab^*c) | (abc^*)$





p. 98

	a	b	c	d
1	2	.	.	5
2	.	3	.	.
3	.	.	4	.
*4	.	.	.	.
5	.	6	.	.
6	.	.	7	.
*7	.	.	.	.

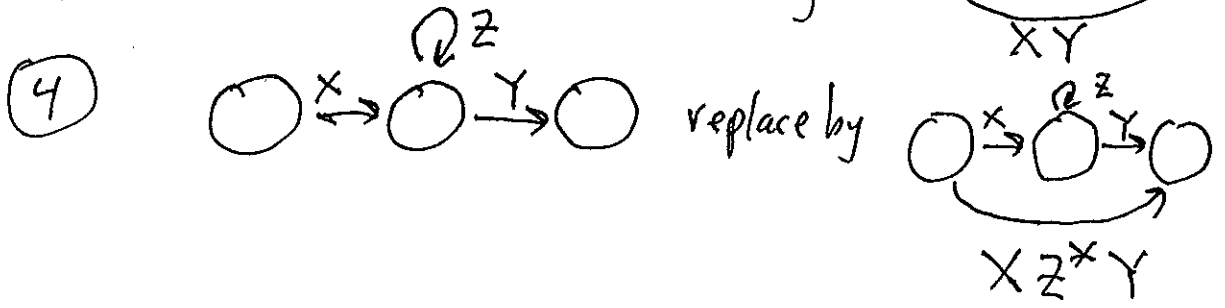
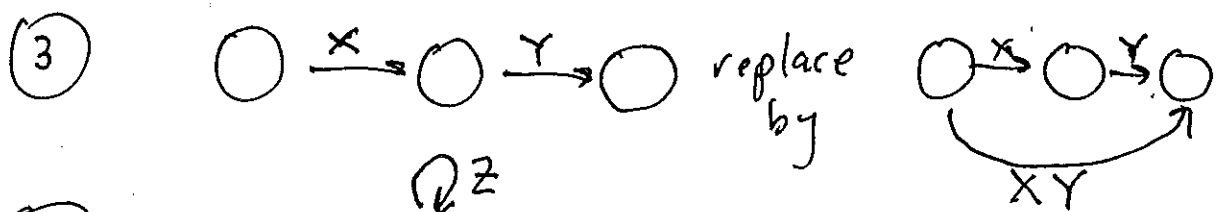


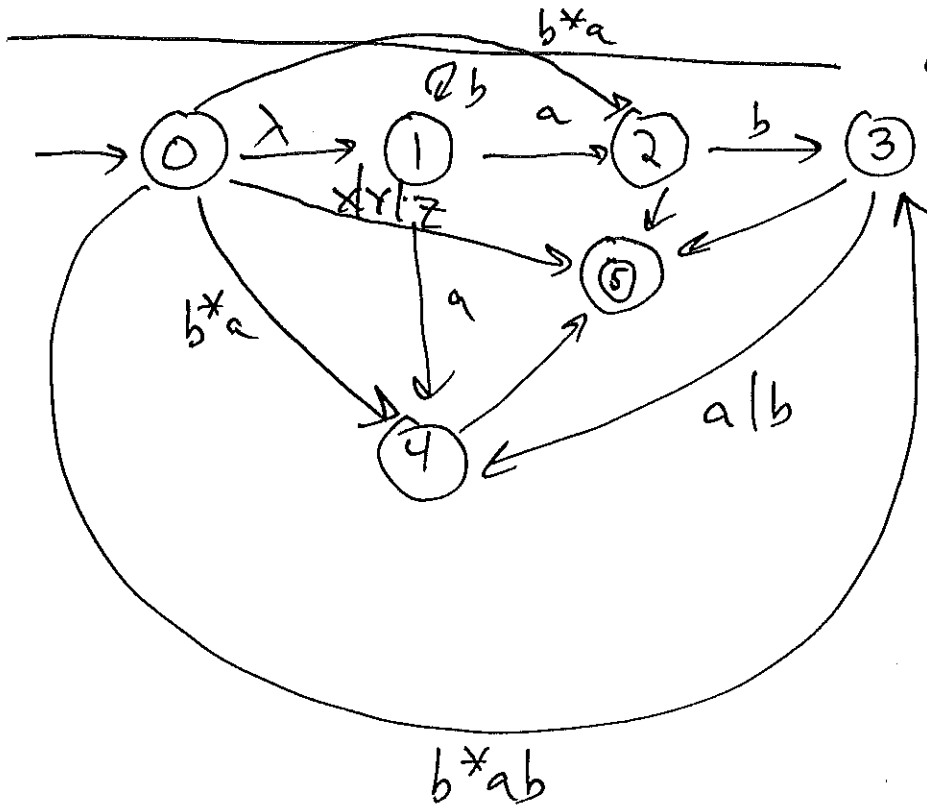
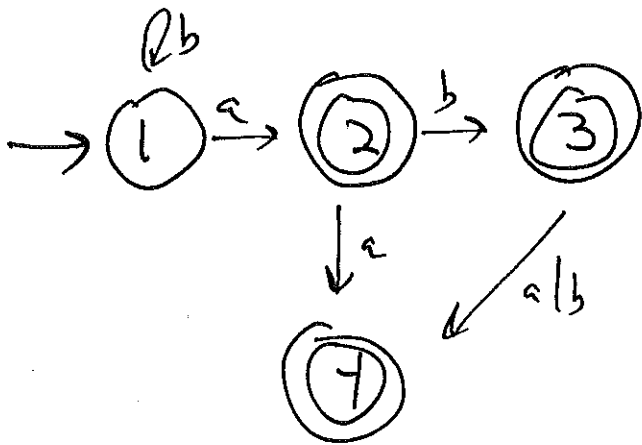
	a	b	c	d
X	Y	.	.	Y
Y	.	Z	.	.
Z	.	.	W	.
*W	.	.	.	.

RE → NDFA → DFA

by transforming the DFA

1) introduce a single accepting state





This figure is mis-copied!

- X  $b^*a$
- Y  $b^*ab$
- Z  $b^*a$
- W  $b^*ab(ab)$

---

Midterm: delayed by 1 week!

---

Chapter 4.

Context-free [ grammars CFGs  
languages — set of strings

Strict superset of regular languages.

instead of a DFA, need a push-down automaton (PDA)

We won't use PDAs

we will use CFGs

alphabet of terminal symbols. (tokens)

set of non-terminal symbols

convention: CamelCaps

start symbol (non-terminal)

usually P

rules (productions)

$$A \rightarrow X_1 \dots X_m$$

↑  
nonterminal

↑  
terminal or nonterminal

$$A \rightarrow \epsilon$$

↑  
epsilon (empty RHS)

language of rules:

BNF

Backus

Naur

Form

sentential form

(31)

$P \Rightarrow A B \Rightarrow a a B$

derivation

Conventions:

leftmost derivation:

~~at~~ always expand the leftmost non-terminal

rightmost derivation

always expand the rightmost non-terminal

read off the derivation backwards

Parse tree

describes a derivation

start symbol is the root.

interior nodes are nonterminals

leaves are terminals.

Regular grammar: like CFG, but RHS

only: a single element of  $\Sigma$

followed optionally by a single non-terminal.

CFG: represented by BNF, can be parsed in  $\Theta(n^3)$  time.

Useful CFGs can be parsed in  $\Theta(n)$  time.

### Context-sensitive grammar

like CFG, but LHS can ~~be~~ have extra "context" around the non-terminal.

can cover the concepts of "declared" but not worth it.

### Type-0 grammar

rules let you rewrite arbitrary patterns.

what kind of CFG?

Not ambiguous. (undecidable)

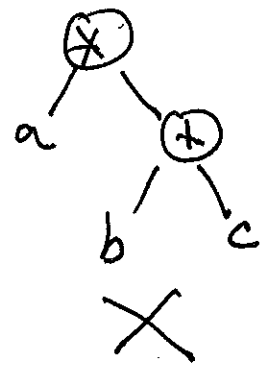
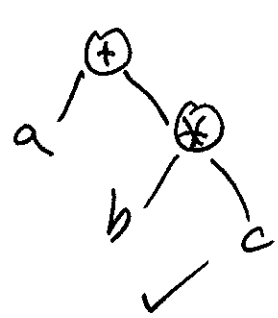
- |   |                    |   |
|---|--------------------|---|
| 1 | $S \rightarrow Aa$ | $S \xrightarrow{1} Aa \xrightarrow{3} aa$ |
| 2 | $\rightarrow aA$   | $S \xrightarrow{2} aA \xrightarrow{3} aa$ |
| 3 | $A \rightarrow a$  |   |
| 4 | $B \rightarrow a$  |   |

Not unreduced

~~not~~ nonterminals that cannot be reached.

Not wrong language.

$a+b*c$





# Extended BNF

to have optional part of RHS

[ optional part ]

to have 0 or more repetitions

{ this stuff }

$$A \rightarrow a [b c] d$$

rewrite as

$$A \rightarrow a N d$$

$$N \rightarrow b c$$

$$\rightarrow \lambda$$

$$A \rightarrow a \{ b c \} d$$

rewrite as

$$A \rightarrow a N d$$

$$N \rightarrow \lambda$$

$$\rightarrow b c N$$

## Notation

recognizer : string  $\rightarrow$  bool

parser : string  $\rightarrow$  parse tree

top-down parser : leftmost derivation  $\downarrow$  LR(1)

bottom-up parser : rightmost derivation examine  $\leftarrow$  LR(1)

1 symbol lookahead

(leftmost)



How to represent a BNF grammar.

BNF  
~~Rules~~: Hash table mapping  
Non-terminals → List of Rules.

Rule: List of RHS elements

Element: either a non-terminal or a terminal

Javai

~~List~~ RHS

```

decl { List <Element> RHS = . . .
      for (e: E
code { for (Element e: RHS) {
      // use e
      }
    }

```

advice: alternative is to use  
 i = RHS.iterator()  
 i.hasNext()  
 i.next()

} if navigating the structure requires modifying it.

LL(1) : top down  
recursive descent

magic to be explained: how to choose among competing rules.

Algorithms.

- 1) Rule Derives Empty (rule)  $\rightarrow$  bool
- 2) Symbol Derives Empty (nonterminal)  $\rightarrow$  bool
- 3) First(nonterminal) = set of terminal
- 4) Follow(nonterminal) = set of terminal
- 5) Predict(rule) = set of terminal

1, 2) List<sup>L</sup> all nonterminals that directly yield  $\lambda$ .

for each  $N \in L$ , find all rules where  $N$  is in the RHS. Pretend to remove it. (by subtracting from count of RHS symbols)

If any RHS count becomes 0, add its LHS to  $L$ .

$$\text{First}(\alpha) = \{ b \in \Sigma \mid \alpha \Rightarrow^* b\beta \}$$

↑ string of symbols      ← Terminals

### Algorithm

Hint: if BNF is written top-down, then compute First() from the end to the beginning.

Consider first symbol of  $\alpha$ .

Easy:  $\alpha$  is empty. ~~answer~~ answer = { }

~~if~~ first symbol is terminal  
 answer = { that symbol }

Hard: first symbol is a nonterminal, N.

for each rule with LHS = N,  
 compute First(RHS) [avoid computing First(N)]  
 answer = Union of all those computations.

if N can derive  $\lambda$ , then  
 union in ~~the~~ First(rest of  $\alpha$ )

Follow(N) = { terminals that can come after N }

$$= \{ b \in \Sigma \mid S \xRightarrow{*} \alpha N b \beta \}$$

↑  
start symbol

### Algorithm

Hint: if BNF is top-down, start at top and work down.

for each place that N appears in a RHS of some ~~rule~~ rule r, union in First( symbols following N )  
"tail"

if tail can derive  $\lambda$  (or is empty),

union in Follow(LHS(~~r~~))  
r

Predict(r) = First(RHS(r))

if r derives  $\lambda$  then Union Follow(LHS(r))

Recursive-descent parser (LL(1))  
top-down  
1 token lookahead  
read from start of program

One procedure for each non-terminal.

A case for each rule, determined by the predict set.

"Absorb" the RHS by match (terminals)  
recursive call (for non-terminals)

```

proc P() {
  switch (peekc) {
    case f, case i, case id, case p, case $ :
      Ds();
      S0();
      match('$');
      break;
    default:
      error()
  } // switch
} // proc P()

```

```

proc Ds() {
  switch (peekc) {
    case f, case i : D(); Ds(); break;
    case id, case p, case $ : break;
    default : error();
  } // switch
} // proc Ds()

```

p. 148 example

	First(RHS)	Follow(T)	(39) Prad. 2
$S \rightarrow A C \$$	$\{a, b, g, c, \$\}$	$\emptyset$	$\{a, b, g, c, \$\}$
* $C \rightarrow c$	$\{c\}$		$\{c\}$
$\rightarrow \lambda$	$\emptyset$	$\{\$, d\}$	$\{\$, d\}$
* $A \rightarrow a B C d$	$\{a\}$	$\{c, \$\}$	$\{a\}$
$\rightarrow B Q$	$\{b, g\}$		$\{b, g, c, \$\}$
* $B \rightarrow b B$	$\{b\}$	$\{c, d, g, \$\}$	$\{b\}$
$\rightarrow \lambda$	$\emptyset$		$\{c, d, g, \$\}$
* $Q \rightarrow g$	$\{g\}$	$\{c, \$\}$	$\{g\}$
$\rightarrow \lambda$	$\emptyset$		$\{c, \$\}$

proc A() {

switch (peekc) {

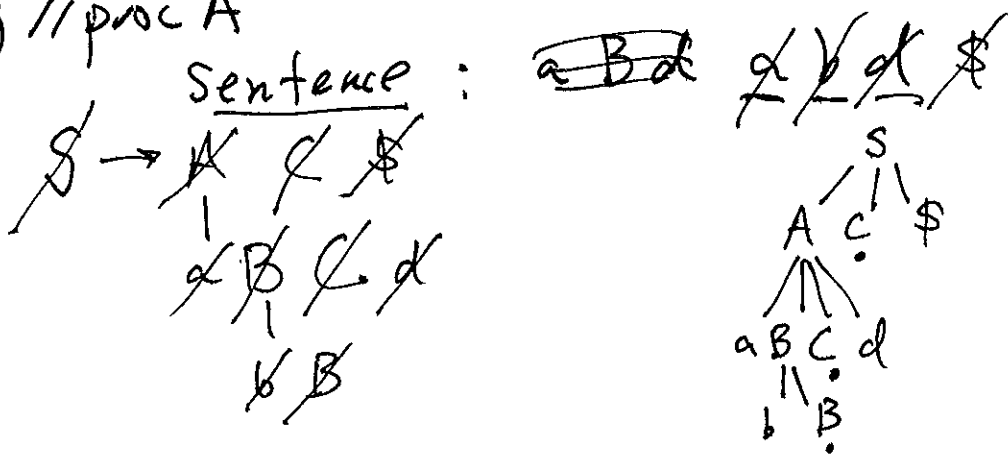
case a : match(a); procB(); procC(); match(d);  
          break;

case {b, g, c, \$}: procB(); procQ(); break;

default: error()

  } // switch

} // proc A



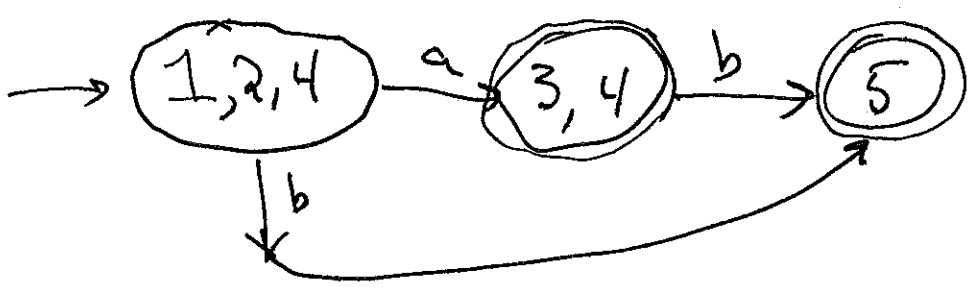
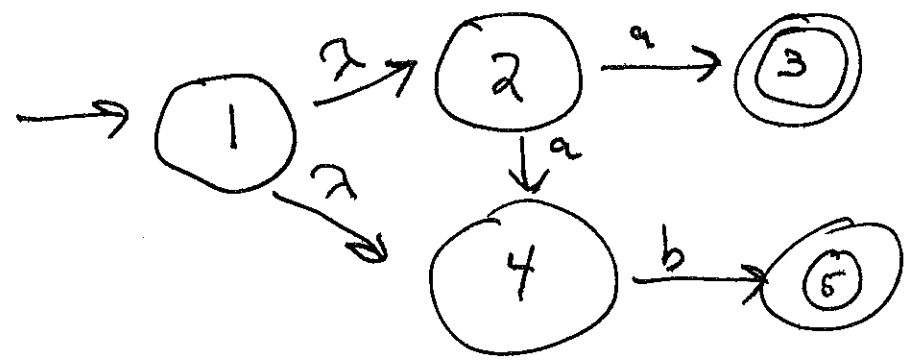
$S \rightarrow a b$	Predict
$\rightarrow a c$	$\{a\}$
	$\{a\}$

conflict in predict sets  $\Rightarrow$  grammar not LL(1)

fix: (factor out common prefix)

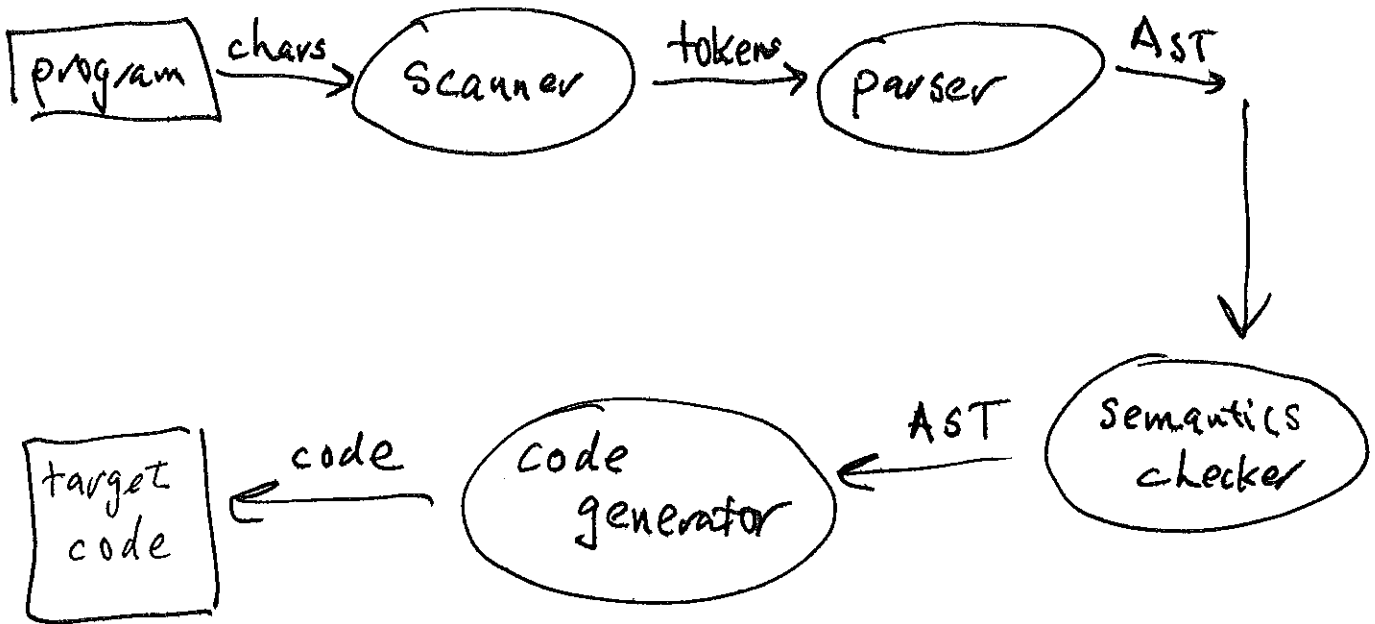
~~$S \rightarrow A b$~~   
 $S \rightarrow a B$   
 $B \rightarrow b$   
 $\rightarrow c$

### Subset Construction





1a)

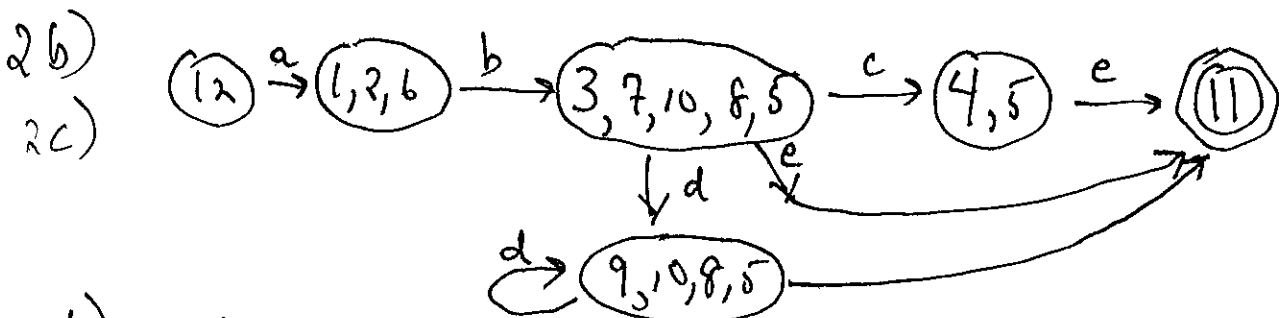
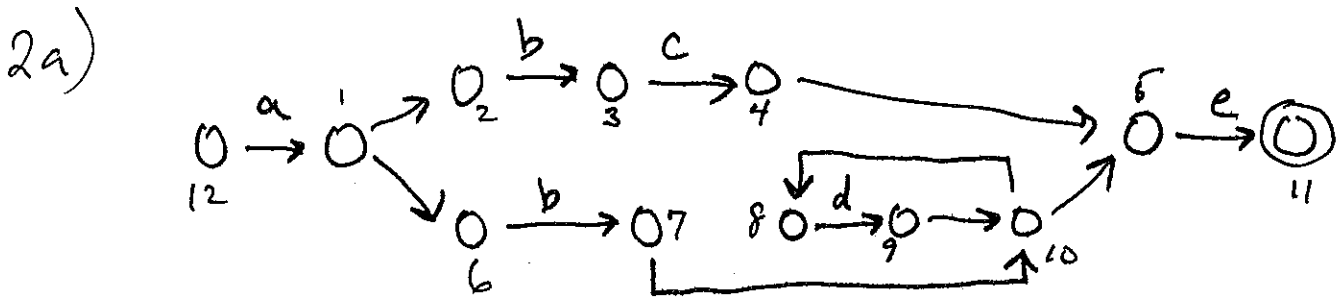


1b) jflex

1c) no

1d)  $\Theta(1)$

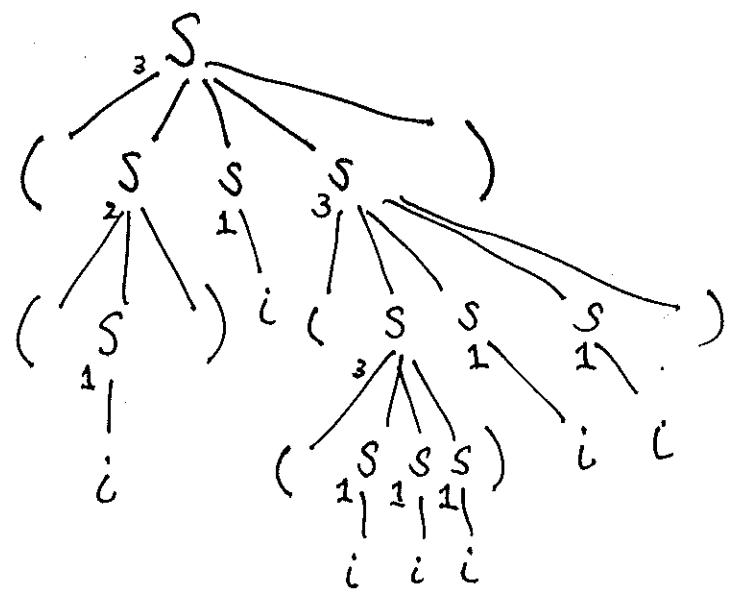
$a(bc|bd^*)e$



2d) abdde

- 1  $S \rightarrow i$
- 2  $S \rightarrow (S)$
- 3  $S \rightarrow (S S S)$

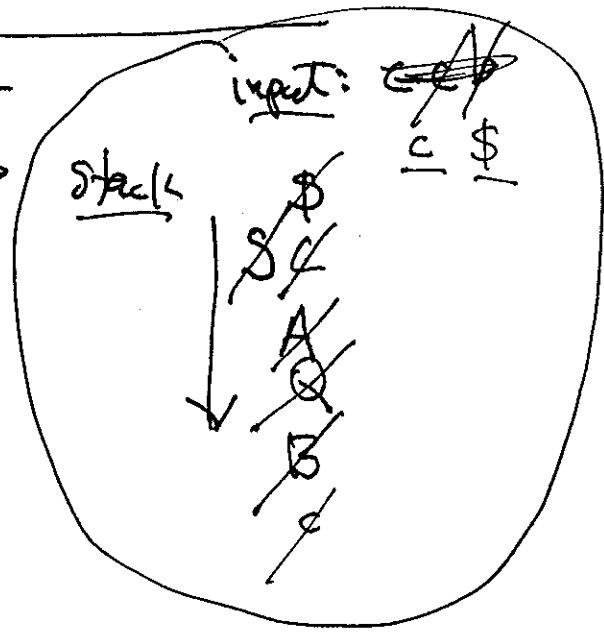
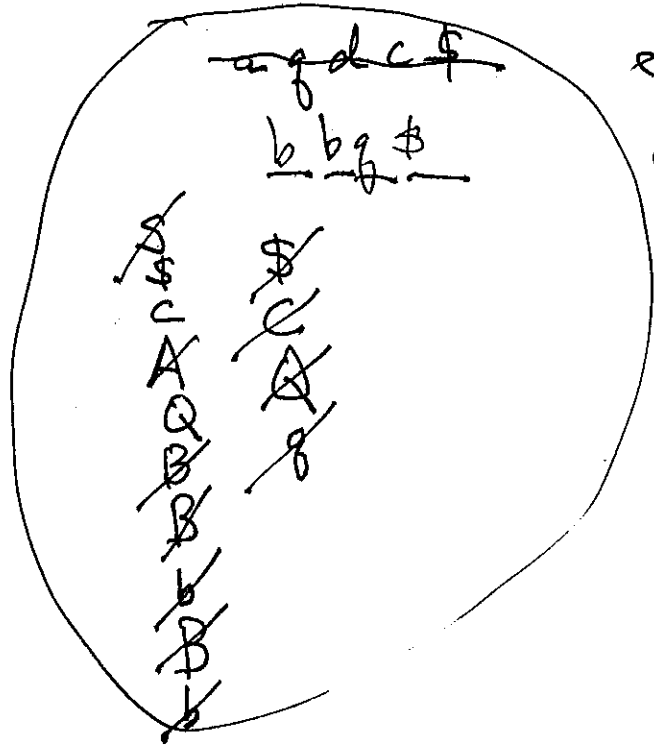
((i)i ((i i L) i i))



Not  $LL(1)$

Table-driven  $LL(1)$  parser

examples



Advantage:  
 "generic" parser; the parser generator just builds a table.

Problem: The parse table can be sparse.

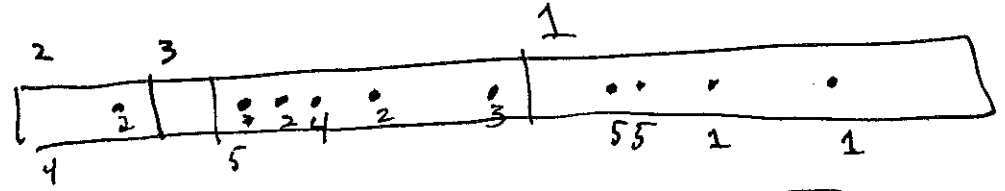
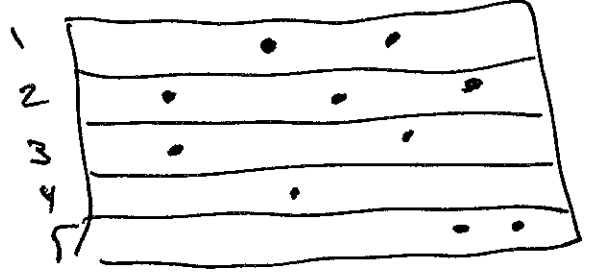
- Solution:
- 1) Live with it.
  - 2) Compress the table.

How?

- a) default value for each row.
- b) hash the table:

$T[i, j]$ : search at  $h(i, j)$

- c) double-offset indexing



To make a grammar LL(1)

- |   |                        |   |                      |
|---|------------------------|---|----------------------|
| 1 | $S \rightarrow id$     | 1 | $S \rightarrow id$   |
| 2 | $S \rightarrow (S$     | 2 | $S \rightarrow (S R$ |
| 3 | $S \rightarrow (S S S$ | 3 | $R \rightarrow )$    |
|   |                        | 4 | $\rightarrow S S )$  |

factoring

$S \rightarrow \text{if } E \text{ then } S_s \text{ end}$   
 $\rightarrow \text{if } E \text{ then } S_s \text{ else } S_s \text{ end}$

factor:

$S \rightarrow \text{if } E \text{ then } S_s R$   
 $R \rightarrow \text{end}$   
 $\quad \text{else } S_s \text{ end}$

fixing left recursion

$S_s \rightarrow S_s \boxed{\begin{matrix} ; \\ S \end{matrix}} \alpha$   
 $\rightarrow S$

$S_s \rightarrow X Y$   
 $X \rightarrow S$   
 $Y \rightarrow ; S Y$   
 $Y \rightarrow \lambda$

$E \rightarrow \text{var} + E$   
 $\rightarrow \text{var}$

$E \rightarrow \text{var } R$   
 $R \rightarrow + E$   
 $\rightarrow \text{var } \lambda$

General case:

$A \rightarrow A \alpha$

change to:

$A \rightarrow X Y$   
 $\quad \uparrow \quad \uparrow$   
 $\quad \text{new nonterminals}$

for every other production

$A \rightarrow \beta$

change to:  $X \rightarrow \beta$

add production

$Y \rightarrow \alpha Y$

add production

$Y \rightarrow \lambda$

if  $B_1$   
     if  $B_2$   
          $S_1$   
     else  
          $S_2$  ;

Dangling "else" problem

2  $S \rightarrow$  if  $E$  then  $S$   $V$   
 3        $\rightarrow$  other  
 4  $V \rightarrow$  else  $S$   
 5        $\rightarrow \lambda$

<u>table</u>	<u>if</u>	<u><del>E</del></u>	<u>then</u>	<u>else</u>	<u>other</u>
$S$	2				3
$V$				4,5	

Properties of LL(k) parsers

- 1) compute correct leftmost parse
- 2) grammar is unambiguous
- 3) table driven parser uses  $\Theta(n)$  time,  
 $\Theta(n)$  space. ( $n = \#$  of tokens)

(4b)

Recovering from syntax errors.

- 1) print line, col, indicate what terminal was expected (predict set or match() param), what was seen.
- 2) try to enter a state you can continue from.
  - bad recovery causes a cascade of errors.
  - "panic mode": skip input until find a frequent delimiter, like semicolon.
  - Wirth: pass extra parameter to match(): set of possible lookaheads.

Repairing syntax errors

modify input to obtain acceptable parse.

situation

program  
 $\underline{\alpha} X \beta$   
↑

erroneous

- 1) modify  $\alpha$ : prefer not, because  $\alpha$  is fine as it is.
- 2) insert  $\delta$  to get  $\alpha \underline{\delta} X \beta$ . "insert-correctable language"
- 3) delete  $X$  and keep going.  
makes progress.

# Chapter 6: Bottom-up parsing.

good news: handles largest class of grammars that can be parsed deterministically.

bad news: less intuitive than top-down.

general idea: start with leaves of the parse tree: terminals.

occasionally reduce a string of ~~terminals~~ symbols to a non-terminal.

parsers are called shift-reduce parsers.

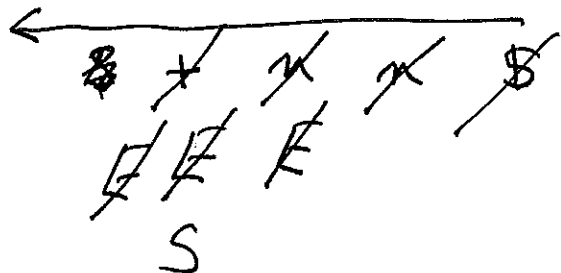
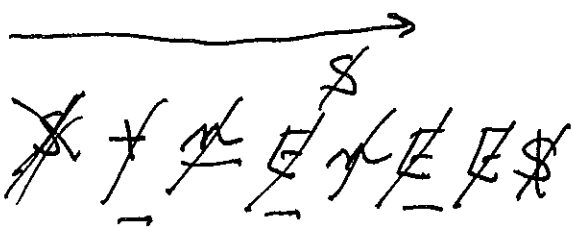
LR(K) ← number of tokens of lookahead  
↑ rightmost derivation  
read from start to end of program

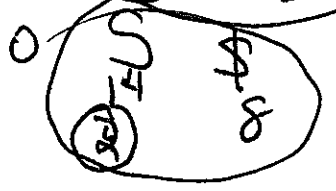
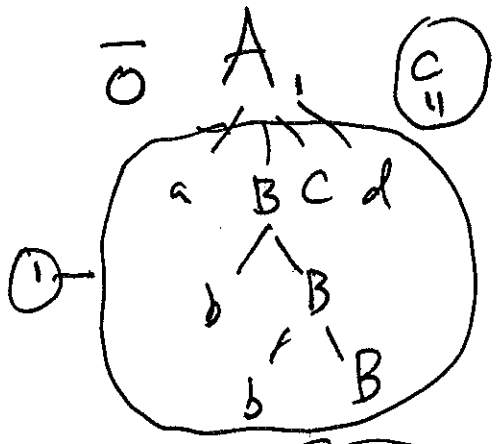
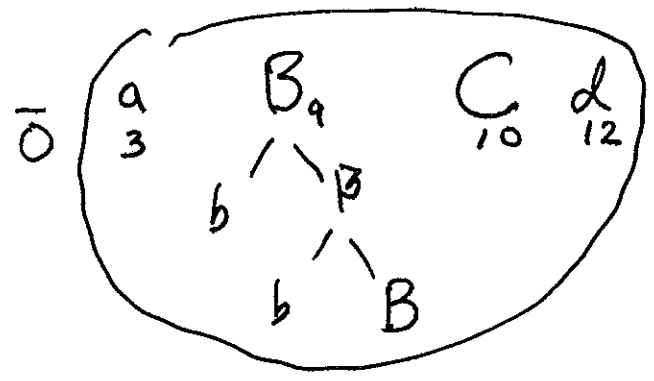
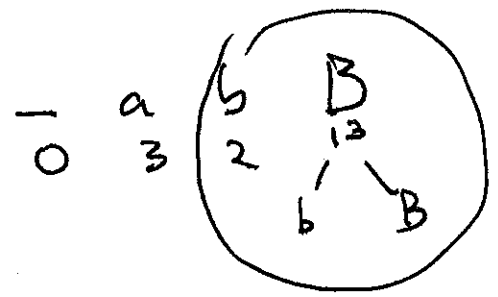
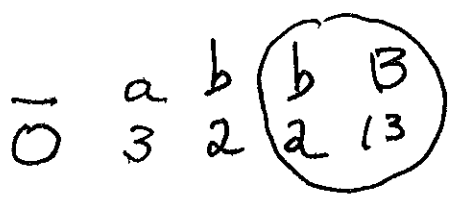
1  $S \rightarrow E \$$

2  $E \rightarrow + E E$

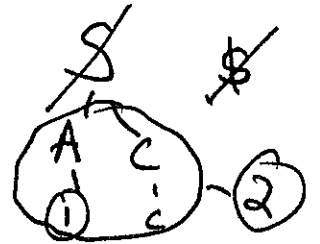
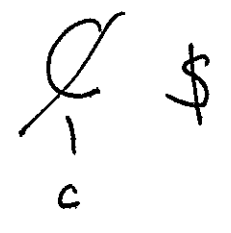
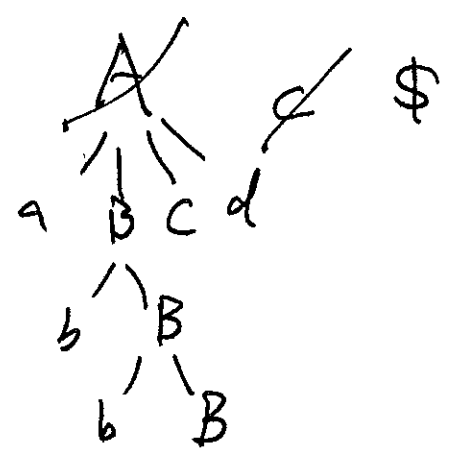
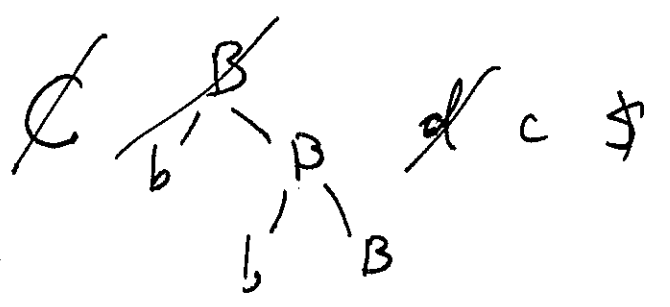
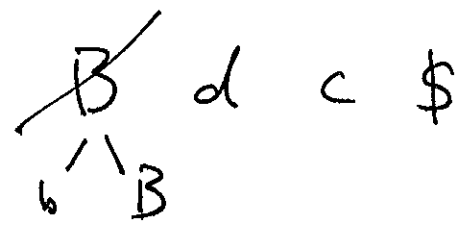
3  $\rightarrow n$

parse: + n n \$





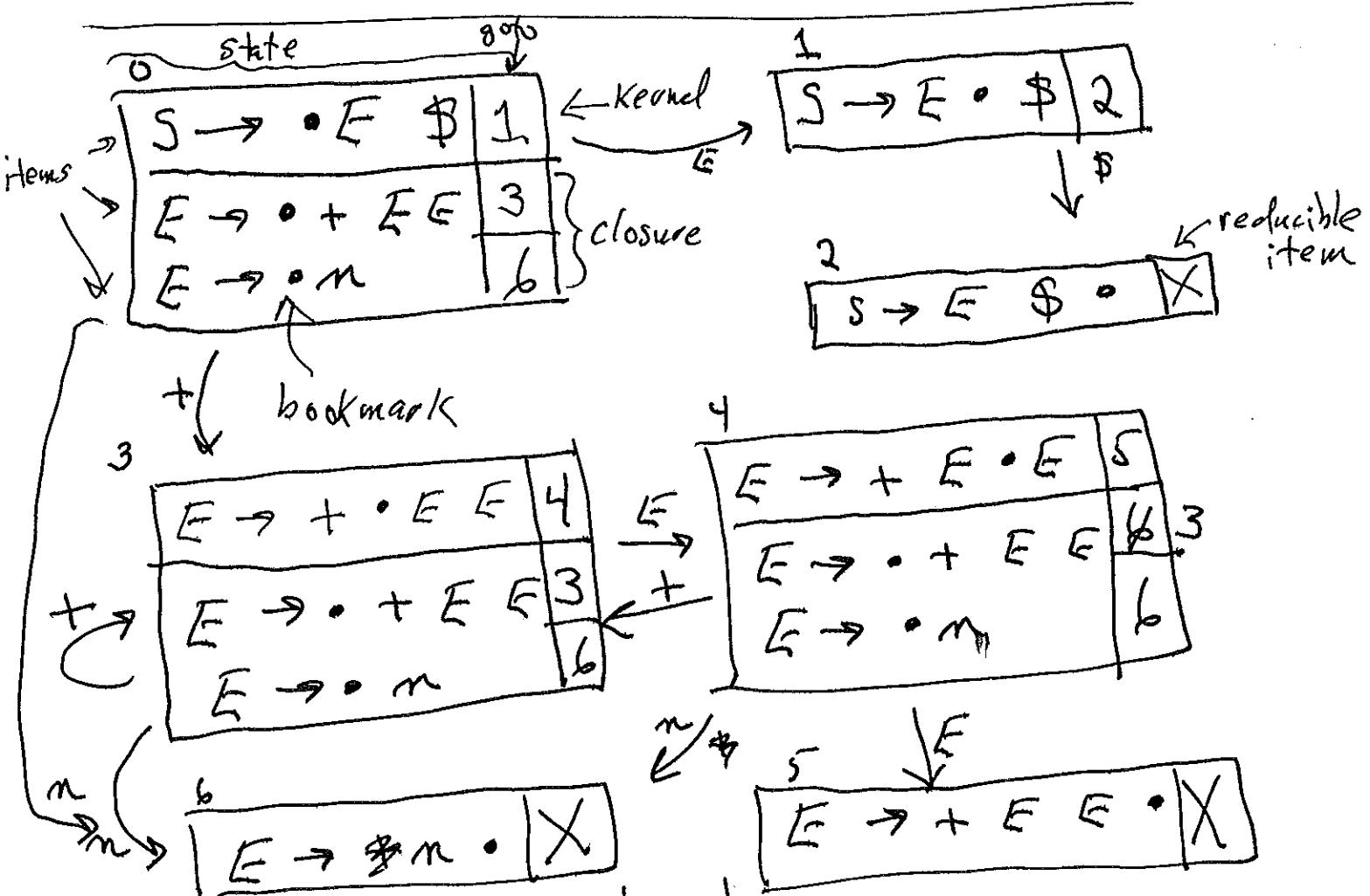
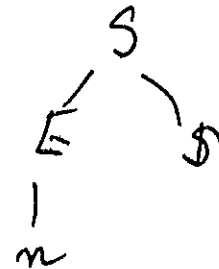
~~B~~ a b b d c \$



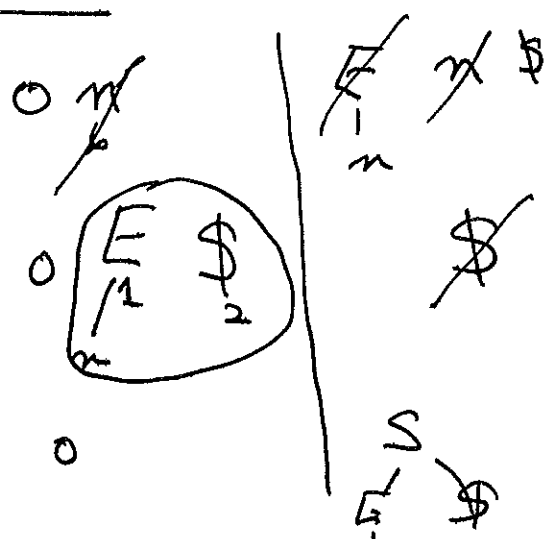


grammar (6.2 p 184)

- 1  $S \rightarrow E \$$
- 2  $E \rightarrow + E E$
- 3  $E \rightarrow n$



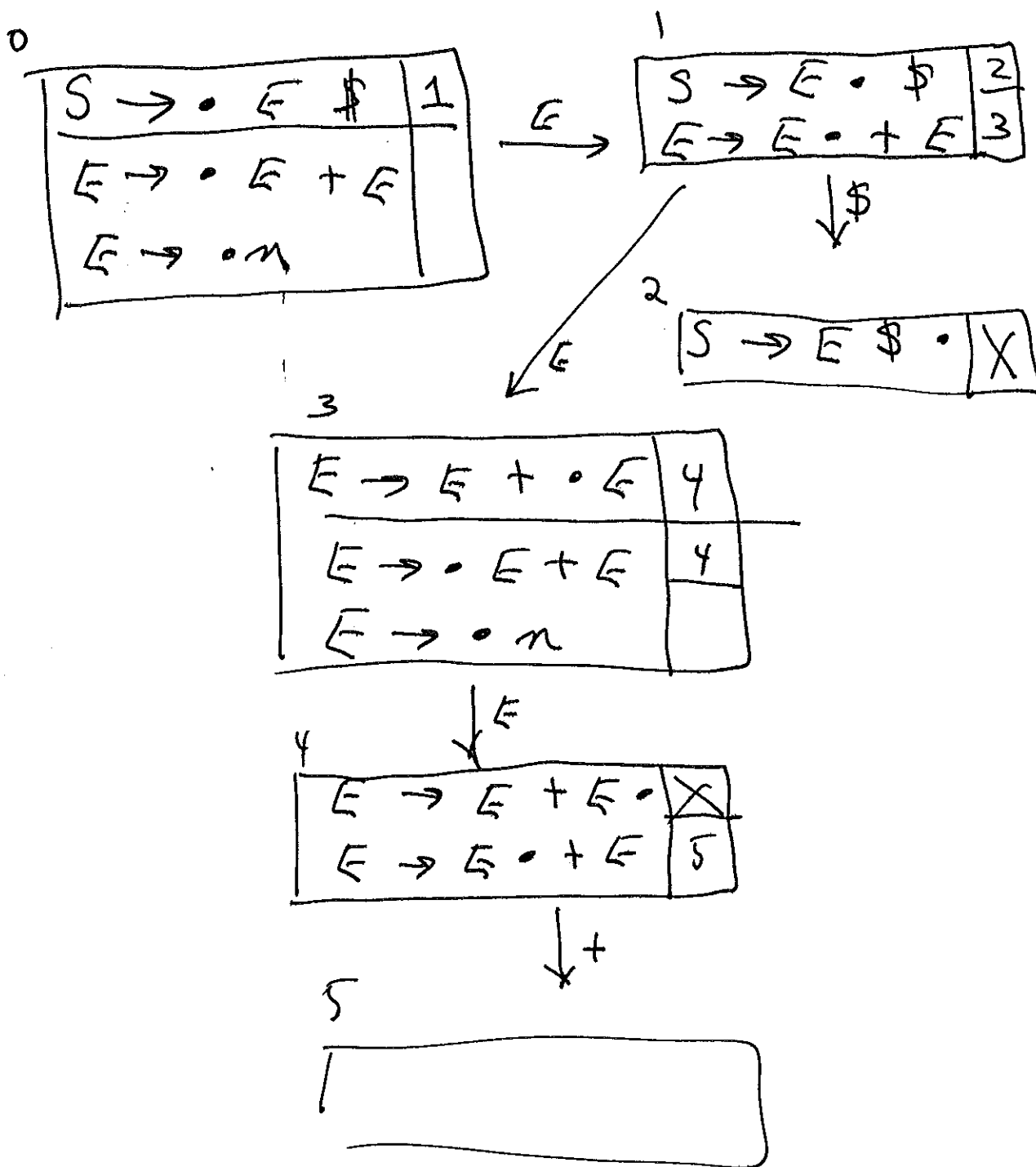
	+	n	\$	S	E
0	[3]	[6]		accept	[1]
1	X	X	[2]	X	X
2	1	1	1	1	1
3	[3]	[6]			[4]
4	[3]	[6]			[5]
5	2	2	2	2	2
6	3	3	3	3	3



$$S \rightarrow E \$$$

$$E \rightarrow E + E$$

$$E \rightarrow n$$



Example 10g p. 227

(51)

- 1  $St \rightarrow S \$$
- 2  $S \rightarrow i = A ;$
- 3  $A \rightarrow i = A$
- 4  $A \rightarrow E$
- 5  $E \rightarrow E + P$
- 6  $E \rightarrow P$
- 7  $P \rightarrow i$
- 8  $P \rightarrow ( i ; i )$
- 9  $P \rightarrow ( A )$

$St \rightarrow \cdot S \$$	1
$S \rightarrow \cdot i = A ;$	3

$St \rightarrow S \cdot \$$	2
-----------------------------	---

$St \rightarrow S \$ \cdot X \$$	
----------------------------------	--

$S \rightarrow i \cdot = A ;$	4
-------------------------------	---

$S \rightarrow i = \cdot A ;$	5
$A \rightarrow \cdot i = A$	7
$A \rightarrow \cdot E$	10
$E \rightarrow \cdot E + P$	10
$E \rightarrow \cdot P$	
$P \rightarrow \cdot ( i ; i )$	
$P \rightarrow \cdot ( A )$	7

$S \rightarrow i = A \cdot ;$	6
-------------------------------	---

$S \rightarrow i = A ; \cdot X \$$	
------------------------------------	--

$A \rightarrow i \cdot = A$	8
$P \rightarrow i \cdot$	+);

$A \rightarrow i = A \cdot X ;$	
---------------------------------	--

$A \rightarrow E \cdot$	X ; )
$E \rightarrow E \cdot + P$	11 +

$A \rightarrow i = \cdot A$	9
$A \rightarrow \cdot i = A$	7
$A \rightarrow \cdot E$	10
$E \rightarrow \cdot E + P$	10
$E \rightarrow \cdot P$	
$P \rightarrow \cdot i$	
$P \rightarrow \cdot ( A )$	7
$P \rightarrow \cdot ( i ; i )$	

$E \rightarrow E + \cdot P$	12
$P \rightarrow \cdot i$	
$P \rightarrow \cdot ( A )$	
$P \rightarrow \cdot ( i ; i )$	

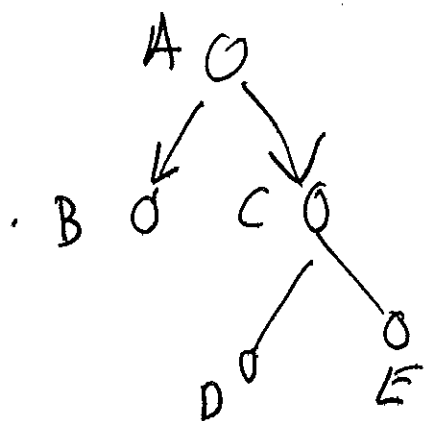
$E \rightarrow E + P \cdot X + ;$	
-----------------------------------	--

# Syntax-directed translation

taking actions during parsing.

↑ syntactic actions

these actions may manipulate semantic values of symbols.



parse tree

information flowing downward: inherited.

flowing upward: synthesized.

Top-down parsing can pass information both ways.

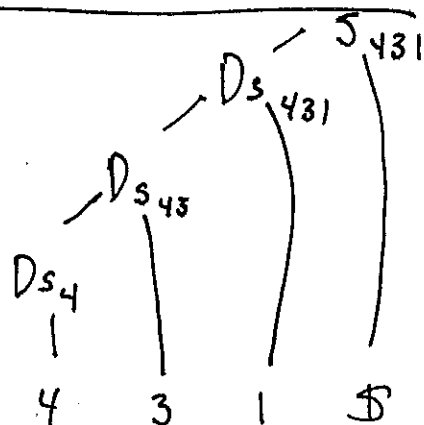
Bottom-up parsing: synthesized only.

Example 1: decimal digits

$$S_e \rightarrow D_s_f \quad \{e = f\}$$

$$D_s_a \rightarrow D_s_b \quad d_c \quad \{a = b \cdot 10 + c\}$$

$$D_s_a \rightarrow d_c \quad \{a = c\}$$



### Example 2: decimal and octal digits

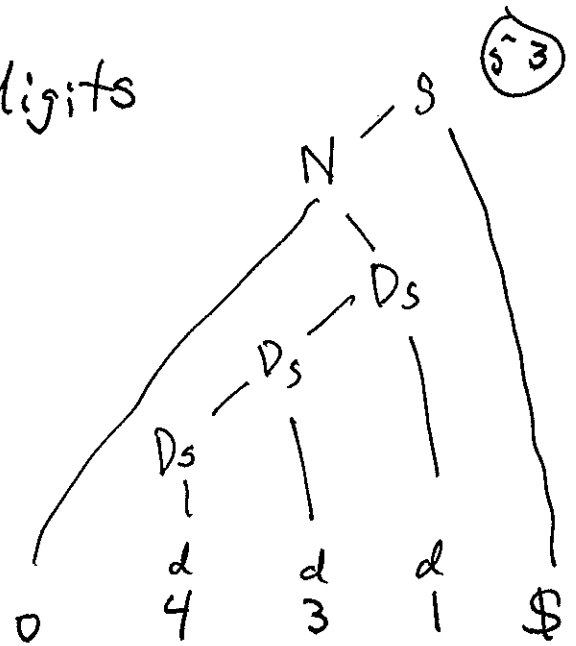
$$S \rightarrow N \$$$

$$N \rightarrow 0 D_s$$

$$N \rightarrow D_s$$

$$D_s \rightarrow D_s d$$

$$D_s \rightarrow d$$



problem: cannot inform the semantic actions whether octal or decimal.

### Example 3

rule  $S_a \rightarrow N_b \$$   
 $a = b$

cloning  $N_a \rightarrow 0 O_{s_b}$   
 $a = b$

$$N_a \rightarrow D_{s_b}$$

$$a = b$$

$$D_{s_a} \rightarrow D_{s_b} d_c$$

$$a = b \cdot 10 + c$$

$$D_{s_a} \rightarrow d_c$$

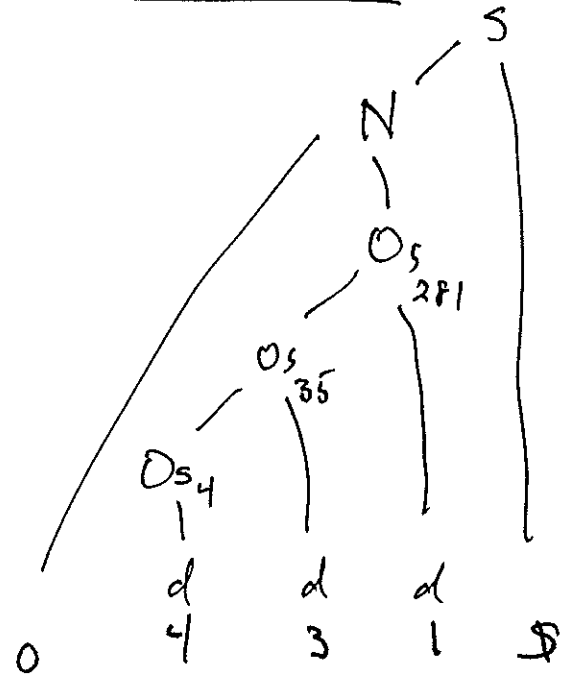
$$a = c$$

$$O_{s_a} \rightarrow O_{s_b} d_c$$

$$a = b \cdot 8 + c$$

$$O_{s_a} \rightarrow d_c$$

$$a = c$$



# Example 4: forcing semantic actions

6-4

$$S_a \rightarrow N_b$$

$a = b$

$$N_a \rightarrow \text{Sig } D \quad Ds_b$$

$a = b$

$$N_a \rightarrow \text{Sig } 0 \quad Ds_b$$

$a > b$

$$\text{Sig } D \rightarrow \lambda$$

$\text{base} = 10$

$$\text{Sig } 0 \rightarrow 0$$

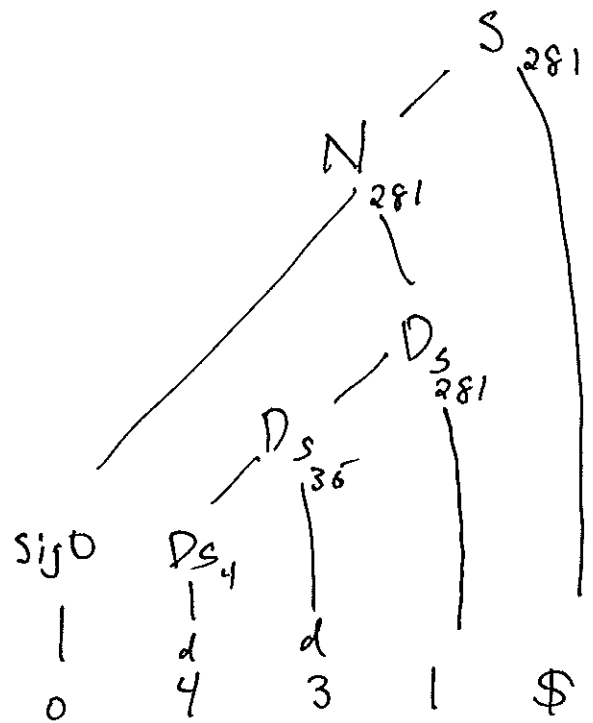
$\text{base} = 8$

$$Ds_a \rightarrow Ds_b \quad d_c$$

$a = b \cdot \text{base} + c$

$$Ds_a \rightarrow d_c$$

$a = c$



Generalize the lg: x5 4 3 1 \$

$$S_{n'} \rightarrow Ds_{n'} \quad \$$$

$n' = n$

$$Ds_{n'} \rightarrow Ds_n \quad d_c$$

$n' = n$   
 $b' = b$   
 $a' = a \cdot n + c$

$$Ds_{n'} \rightarrow \text{SetBase } n$$

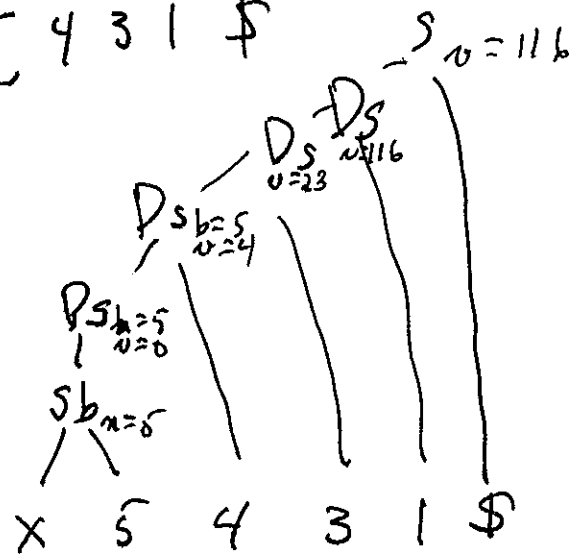
$n < 0$   
 $b \leftarrow n$

$$\text{SetBase } n \rightarrow \lambda$$

$n < 10$

$$\rightarrow X \quad d_c$$

~~$n \leftarrow a$~~   
 $n \leftarrow c$



Enhance the lg to have any base

$$x[37] 4 2 1 \$$$

also allow:

$$x[x[8]24] 4 2 1 \$$$

left as an exercise.

### Top-down syntax-directed translation



### Expression grammar (like Lisp)

$$S_b \rightarrow V_b \$ \quad b' = b$$

$$V_b \rightarrow n \quad b' = n$$

~~$$\rightarrow (+ V V)$$~~

$$\rightarrow (E_b) \quad b' = b$$

$$E_b \rightarrow + V_a V_b \quad b' = a + b$$

$$\rightarrow * V_s V_b \quad b' = b$$

$$V_s \rightarrow V_a V_b \quad b' = a * b$$

$$\rightarrow \lambda \quad b' = 1$$

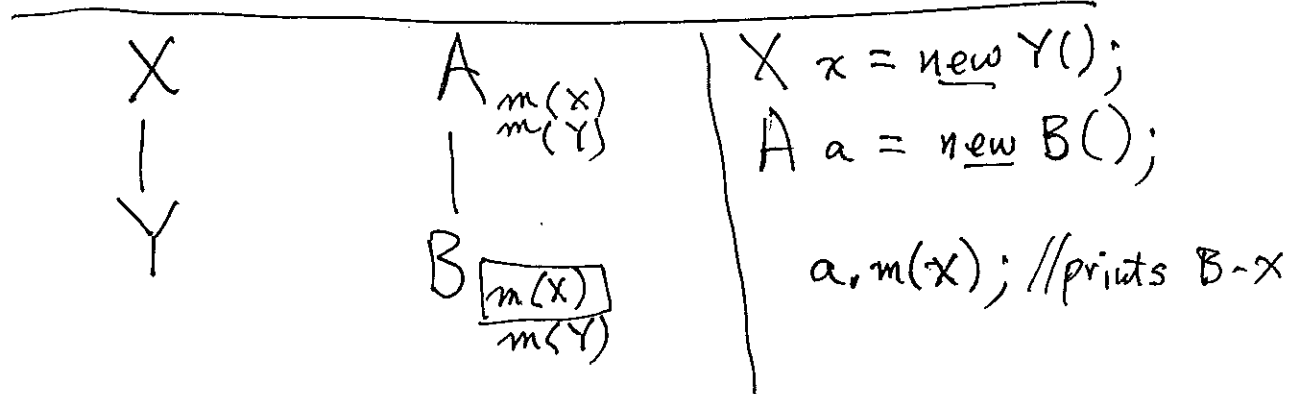
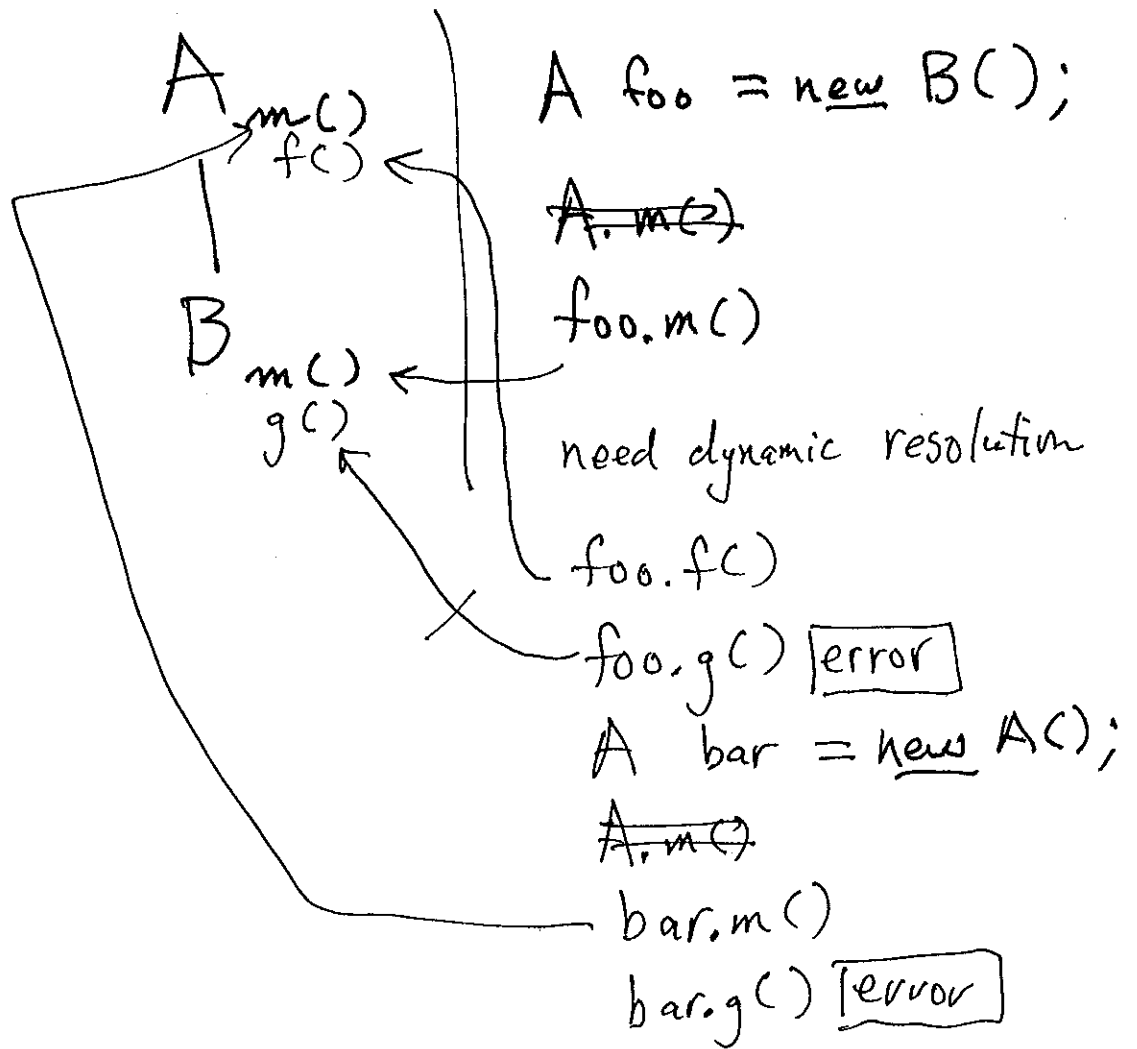
# Method resolution in Java

(56)

dynamic (at run time): based on value of the object

static (at compile time): based on class of the object

foo.m()

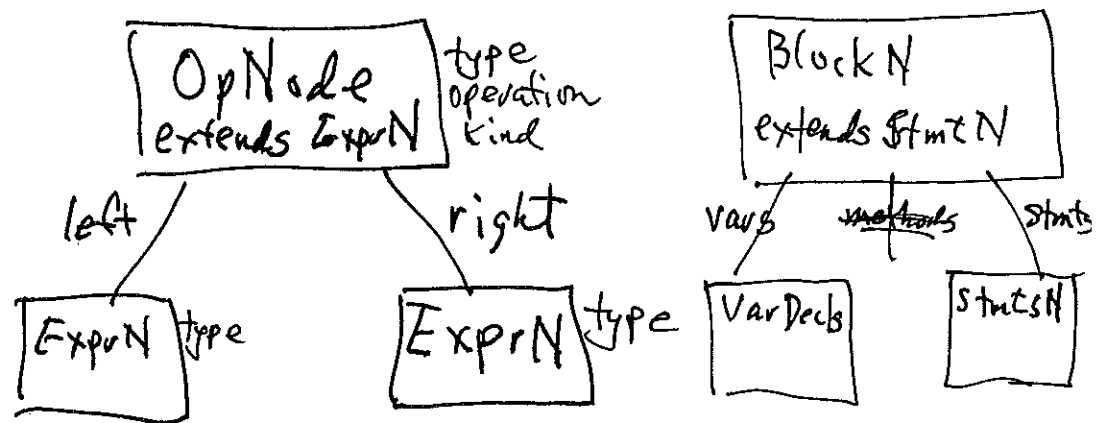
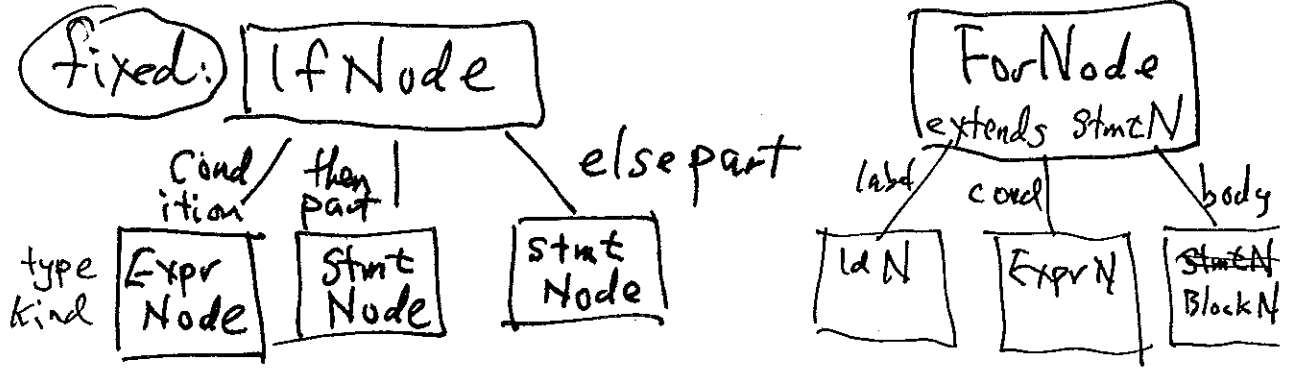




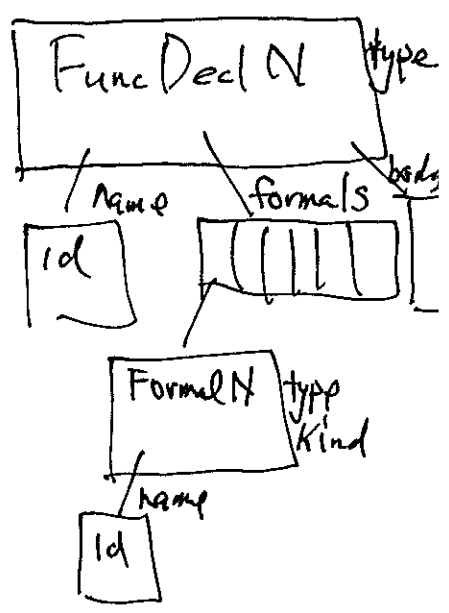
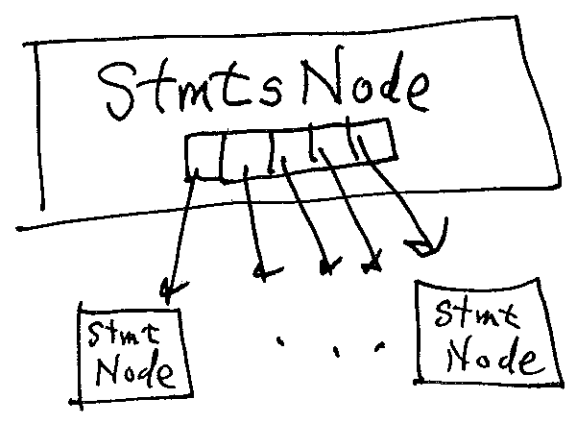
# ASTs

distinguish: parse trees from ASTs.  
|  
concrete syntax trees

Number of children

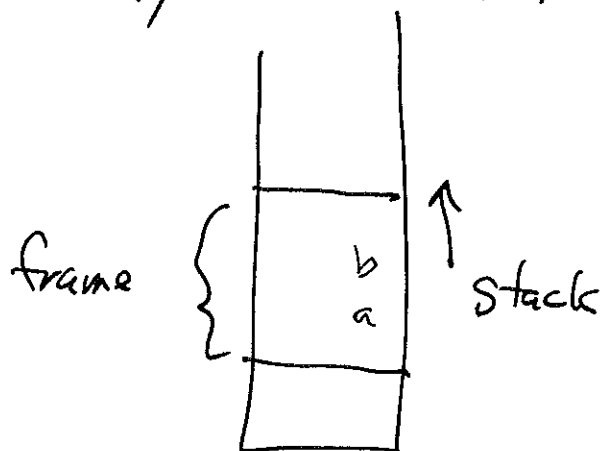


**Variable**



$a := b + 5;$   
 ↑                    ↑                    ↙  
 L-value            R-value (current arithmetic value)  
 (location where variable is stored at runtime)

example: a, b local ints



Facts:

constants (literals) have only R value.

this has only R value

in C, get an L value from an R value  
with & (referencing)

in C, given an L ~~R~~ value, use

\* (dereferencing) to get its R value.

# The Java Virtual Machine (JVM)

executes bytecode packaged in class files  
("test.class")

bytecode designed to be compact.

0-address: implicit stack

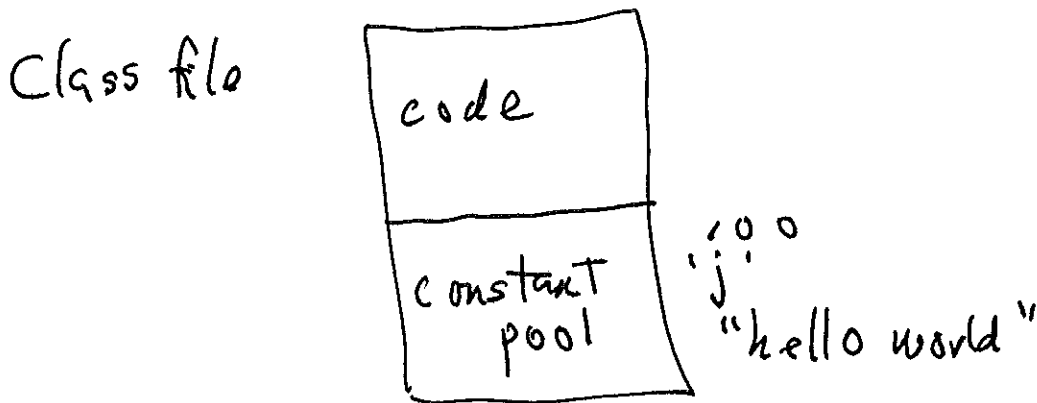
bytecode is safe: a verifier checks for  
type errors and stack errors.

you can disassemble a class file.

```
javap -c test.class
```

Types are described by short strings. p 400

Constants referred to in class files by an <sup>16-bit</sup> index  
into "constant pool". Entries are typed.



## Operations

# Instructions in bytecode

arithmetic: pops 2 operands, does arithmetic,  
~~places~~ pushes result.

iadd (integer addition)

fadd (float addition)

d mul (double multiplication)

Registers to hold local variables and parameters.

Numbered: 0, 1, 2 ... for parameters,  
then local variables.

parameter 0 in dynamic methods:  
this.

registers are untyped.

## Stack operations

register → stack

iload 39

fload

aload 14

reg 14 holds a reference.  
place the reference on  
stack.

shorthands

iload - 3

stack → register

istore - 4

↑  
type

↑  
register #

stack elements, registers are 32-bit values  
long, double use 2 registers, 2 stack elements  
↓  
even, odd

bool, char, byte, short: use 32 bits.

ldc: load a constant literal.

ldc "hello"  
ldc 'x'  
ldc false  
ldc 14  
} translated by Jasmin  
to reference elements  
in constant pool.

type conversion

i2f

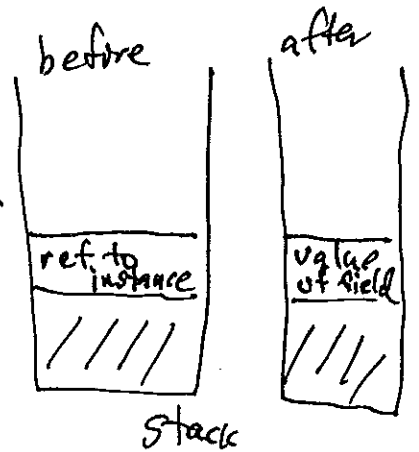
get/put a field of a class (~~inst~~ class → stack)  
get static name type class

↳ Ljava/io/PrintStream

↓  
java/lang/System/out

put static name type  
(stack → class)

get field  
put field



# Calling a method

invokestatic name



java/lang/Math/pow(DD)D

parameters on stack:  
left-most deepest

invokevirtual name

push ref to instance  
push parameters  
call

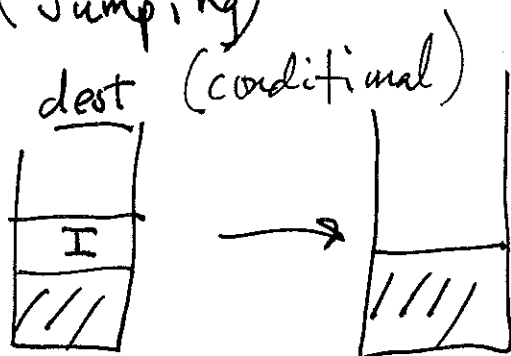
invokespecial name (for constructors)

push ref to uninitialized object  
(created by new)

push parameters

## Branching (Jumping)

ifgt



code:

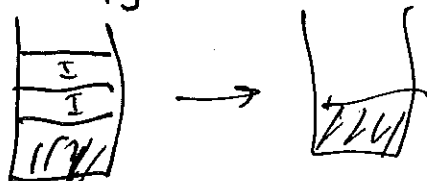
```

inst
inst
inst
L: inst
inst

```

goto dest (unconditional)

if\_icmpgt dest (conditional)



# Stack operations

pop

swap : interchanges top 2 stack elements.

dup-x1 : duplicates the element at top of stack to a position 3 below the top.

```

while Bool {
  body
}

```

```

L2: code for Bool
    conditional branch if false to L1
    code for Body
    go to L2
L1:

```

better code

```

jump L3
L4: code for Body
L3: code for Bool
    conditional branch if true to L3 L4

```

(64)

```

if B {
  S1
} else {
  S2
}

```

code for Bool  
 conditional branch to L1 (if false)  
 code for S1  
 unconditional branch to L2  
 L1: code for S2  
 L2:

```

func foo() {

```

```

  var a; (reg 0)
  var b; (reg 1)
  if B {
    var c; (reg 2)
  } else {
    var d; (reg 3)
  }
}

```

Keep track of reg.  
 Sequence # as ~~an~~ a  
 hidden ST entry.



# Chapter 8: Symbol tables

First used in semantics checker.

Can check for reserved words using ST during scanning.

C: typedef.

Handling nested scopes.

Stack of hash tables

push : open a scope	} semantics checker code generator
pop : close a scope	

possible to build a single unified hash table.

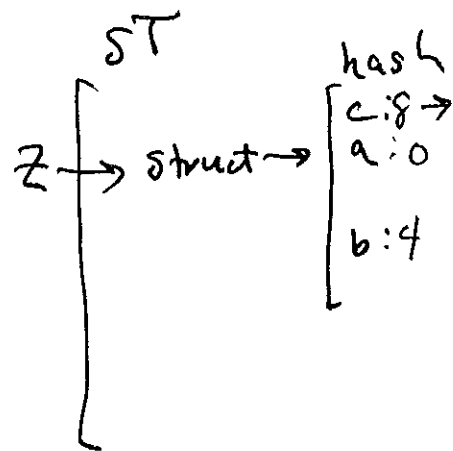
How to treat structs.

represent a struct type in ST as a hash table.

```

struct {
  a: int;
  b: float;
  c: struct {
    d: int;
  }
} z

```



} z

## Overloading

- field (instance variable in class)  
no overloading, only hiding.

- methods

foo(a: int)

foo(b: bool)

- 1) could store in ST: Key = "foo, int"  
"foo, bool"

but to search at the invocation point  
is very expensive.

- 2) store with key "foo"  
data: list of all methods with  
that name.

resolution

to search at invocation:  
get the list of possibilities from  
the closest scope where the  
list is not empty.  
select the "best" of those  
possibilities.

- operators

put definitions in ST.

- enumeration constants (WEO): list in ST

```

var
  a, b record
    x, y: int
  end;
begin
  a.x := 7;
  with a do
    x := 1;
    y := y + 1;
  end;

```

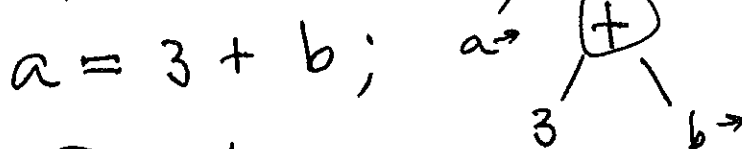
Pascal with  
for efficiency  
and clarity

To handle:  
Open a scope when  
encounter with.

### Processing declarations (semantics checker)

Semantics checker: ST should point the identifier's entry to  
a type descriptor: type(int/char) kind(array?)  
Every access in the AST to that identifier  
points to that same descriptor.

Therefore, code generation need not use ST.



Particular AST nodes.

IfNode

ForNode

Labelled node (ForNode):  
 sem.check: unify, put in  
 code gen: save 2 labels in  
Jasmin

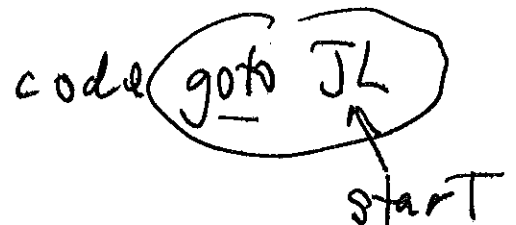
# Labelled for loop.

## Semantics checking

- 1) ensure label is not in ST globally.
- 2) place label in ST. (avoid dup)
- 3) open new scope, put label in new scope, too. (to store code locations)
- 4) every break, continue: check ST for existence of label. point the AST to the label information.

## Code generation

- 1) generate Jasmin Labels (JLs) for the loop (2 or 3)
- 2) store "start", "end" labels in the label information.
- 3) generate code
- 4) every break: code goto JL   
 ↑  
end
- 5) every continue:



# Booleans and generating code for JVM

Represent as 0 (false) 1 (true)

```
ldc 1
iconst_1
```

Handle  $a < b$

```

code to evaluate a (maybe simple: i load r)
code to evaluate b
if_icmplt L1
  integer <
  iconst_0
  goto L2
L1: iconst_1
L2:

```

```

Handle b1 and b2
code for b1
code for b2
iand

```

```

Short-circuit?
code for b1
ifeq L1
code for b2
L1: goto L2 iconst_0
L2:

```

Handle  $b1 = b2$

```

code for b1
code for b2
if_icmpeq

```

similar

```

class foo {
  public static int a, b;
  void static bar() {
    a = b;
  }
}

```

```

.class foo
  .field static a I
  .field static b I
  .method static bar() V
  .local
  .limit locals 1
  .limit stack 10
  getstatic foo/b I
  putstatic foo/a I
  .end method

  .method static main
    ([Ljava/lang/String;) V
    (putstatic for each const)
    invokestatic foo/main() V
    return
  .end method

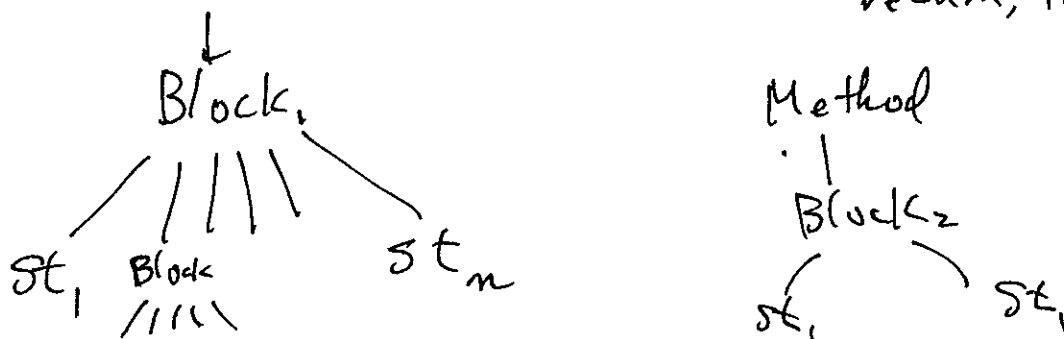
```

# Other semantic checks

- ① type checking
- ② reachability checking
- ③ exception checking

② invalid to have an unreachable statement.  
 each statement node in AST has

bool { is Reachable  
 terminates Normally (not break, continue, return, throw)



st<sub>1</sub> inherits is Reachable, Block<sub>1</sub> synthesizes terminates from st<sub>n</sub>.  
 Block<sub>2</sub> is marked is Reachable.

declarations terminate normally.



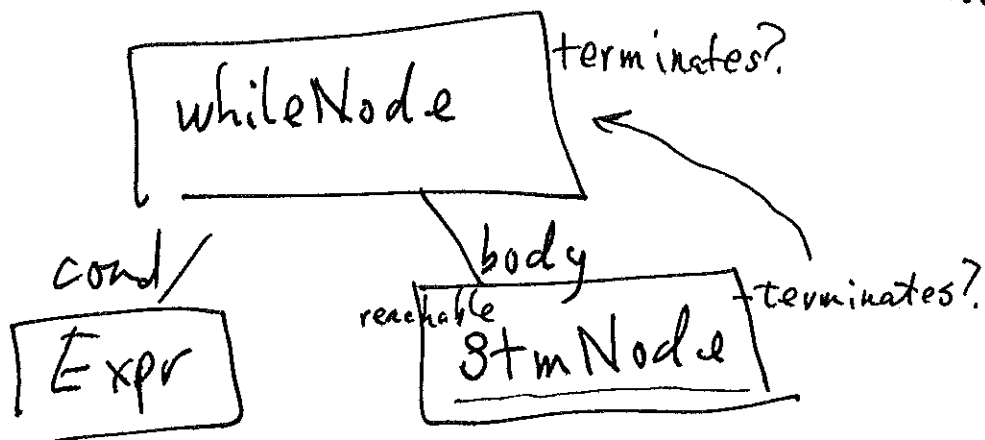
B's is Reachable is derived from A's terminates

return, break, continue, throw do not terminate normally

resulting EN is OR of EN of the children. part's is Reachable.

while B do S end

assume S is reachable  
if S terminates normally, so does while  
otherwise while does not terminate normally.



switch:  $\epsilon N$  is the OR of the  $\epsilon N$ s of all cases.

③ Semantic check in Java for exceptions.   
 ~~checked~~ checked

rules:

any method that can exit with exception must say so.

a catch may only refer to an exception that the try can throw.

multiple catch clauses on a try each may only catch an exception still possible.



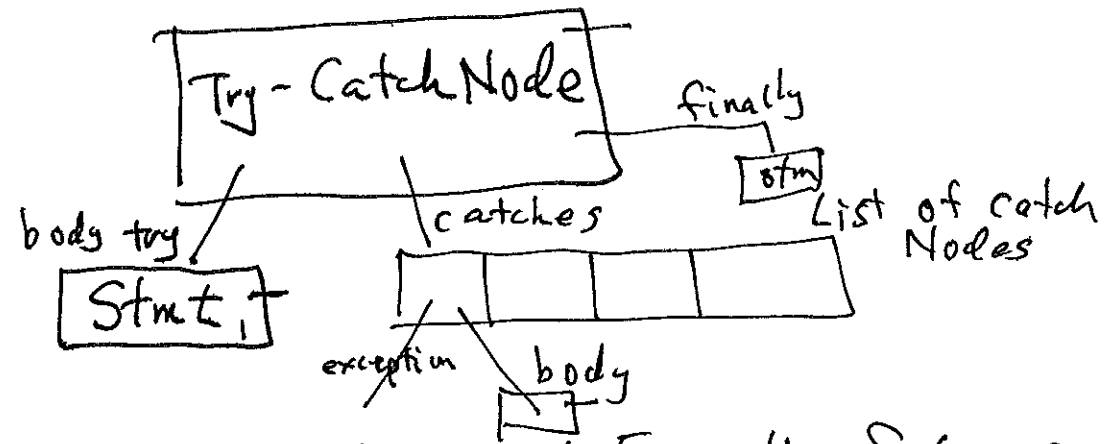
To implement ③, checker for stmts must generate a list of exceptions the Stm can throw. <sup>set</sup>

Block Node can throw any exception that its internal statements can throw.

If Node: inherits all exceptions that condition, then-part, else-part can throw.

Fcn Call Node: can throw all exceptions listed in the Fcn Def Node.

### Try-catch Node



Keep track of possible Exception Set = p  
 initially,  $p = \text{Stmt}. \text{exceptionSet};$   
 for each catch Node  
   verify the exception is in p  
   remove that exception from p  
   add into p all exceptions from body  
 add in the exceptionSet of the finally clause  
 place result as the exceptionSet of TryCatch No

DeclList  $\rightarrow$  DeclList ; Decl ] before  
 $\rightarrow$  Decl

---

DeclList  $\rightarrow$  X Y  
 X  $\rightarrow$  Decl  
 Y  $\rightarrow$  ; Decl Y  
 Y  $\rightarrow$   $\lambda$  ] after

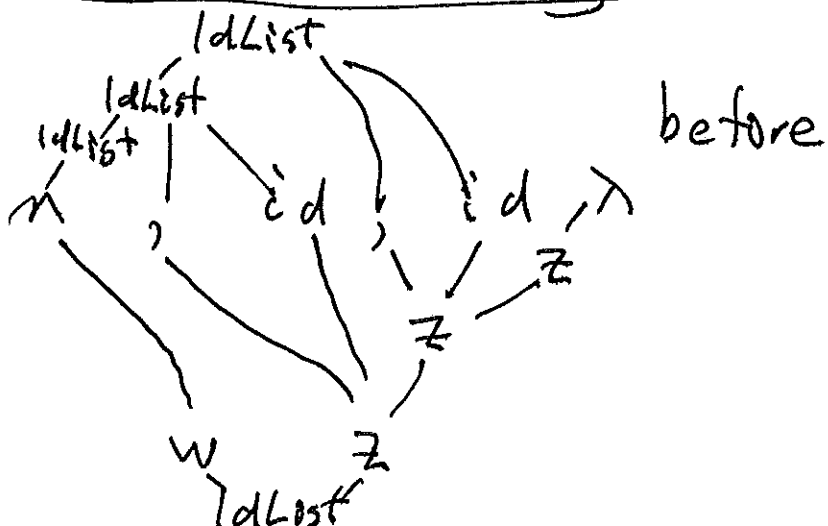
---

IdList  $\rightarrow$  IdList , id ] before  
 $\rightarrow$  id  
 $\rightarrow$   $\lambda$

---

IdList  $\rightarrow$  W Z  
 W  $\rightarrow$  id  
 W  $\rightarrow$   $\lambda$   
 Z  $\rightarrow$  , id Z  
 Z  $\rightarrow$   $\lambda$  ] after

---



State 0

Start $\rightarrow \cdot S \$$	1
$S \rightarrow \cdot id = E ;$	3

S

Start $\rightarrow S \cdot \$$	2
--------------------------------	---

\$

Start $\rightarrow S \$ \cdot$	X
--------------------------------	---

id  
X 3

$S \rightarrow id \cdot = E ;$	4
--------------------------------	---

=>

$S \rightarrow id = \cdot E ;$	5
$E \rightarrow \cdot E + P$	5
$E \rightarrow \cdot P$	
$P \rightarrow \cdot id$	
$P \rightarrow \cdot ( E )$	
$P \rightarrow \cdot id = E$	

E

$S \rightarrow id = E \cdot ;$	6
$E \rightarrow E \cdot + P$	7

;

$S \rightarrow id = E ; \cdot$	X
--------------------------------	---

if B<sub>1</sub> then if B<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>

X