

Compiler course

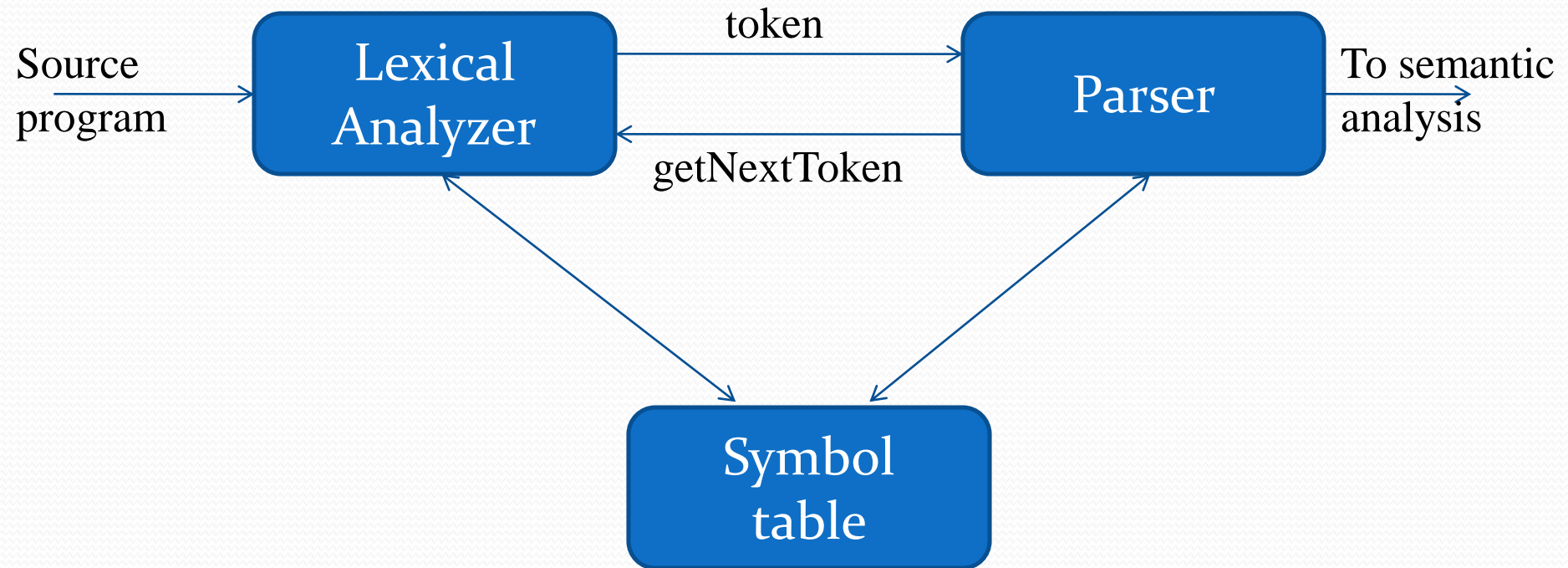
Chapter 1

Lexical Analysis

Outline

- Role of lexical analyzer
- Specification of tokens
- Recognition of tokens
- Lexical analyzer generator
- Finite automata
- Design of lexical analyzer generator

The role of lexical analyzer



Why to separate Lexical analysis and parsing

1. Simplicity of design
2. Improving compiler efficiency
3. Enhancing compiler portability

Tokens, Patterns and Lexemes

- A token is a pair a token name and an optional token value
- A pattern is a description of the form that the lexemes of a token may take
- A lexeme is a sequence of characters in the source program that matches the pattern for a token

Example

Token	Informal description	Sample lexemes
if	Characters i, f	if
else	Characters e, l, s, e	else
comparison	< or > or <= or >= or == or !=	<=, !=
id	Letter followed by letter and digits	pi, score, D2
number	Any numeric constant	3.14159, 0, 6.02e23
literal	Anything but “ sorrounded by “	“core dumped”

```
printf(“total = %d\n”, score);
```

Attributes for tokens

- $E = M * C ** 2$
 - <id, pointer to symbol table entry for E>
 - <assign-op>
 - <id, pointer to symbol table entry for M>
 - <mult-op>
 - <id, pointer to symbol table entry for C>
 - <exp-op>
 - <number, integer value 2>

Lexical errors

- Some errors are out of power of lexical analyzer to recognize:
 - $fi(a == f(x)) \dots$
- However it may be able to recognize errors like:
 - $d = 2r$
- Such errors are recognized when no pattern for tokens matches a character sequence

Error recovery

- Panic mode: successive characters are ignored until we reach to a well formed token
- Delete one character from the remaining input
- Insert a missing character into the remaining input
- Replace a character by another character
- Transpose two adjacent characters

Sentinels



```
Switch (*forward++) {
    case eof:
        if (forward is at end of first buffer) {
            reload second buffer;
            forward = beginning of second buffer;
        }
        else if {forward is at end of second buffer) {
            reload first buffer;\
            forward = beginning of first buffer;
        }
        else /* eof within a buffer marks the end of input */
            terminate lexical analysis;
        break;
    cases for the other characters;
}
```

Specification of tokens

- In theory of compilation regular expressions are used to formalize the specification of tokens
- Regular expressions are means for specifying regular languages
- Example:
 - Letter_(letter_ | digit)*
- Each regular expression is a pattern specifying the form of strings

Regular expressions

- ϵ is a regular expression, $L(\epsilon) = \{\epsilon\}$
- If a is a symbol in Σ then a is a regular expression, $L(a) = \{a\}$
- $(r) \mid (s)$ is a regular expression denoting the language $L(r) \cup L(s)$
- $(r)(s)$ is a regular expression denoting the language $L(r)L(s)$
- $(r)^*$ is a regular expression denoting $(L(r))^*$
- (r) is a regular expression denoting $L(r)$

Regular definitions

$d_1 \rightarrow r_1$

$d_2 \rightarrow r_2$

...

$d_n \rightarrow r_n$

- Example:

$\text{letter_} \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid Z \mid _$

$\text{digit} \rightarrow 0 \mid 1 \mid \dots \mid 9$

$\text{id} \rightarrow \text{letter_} (\text{letter_} \mid \text{digit})^*$

Extensions

- One or more instances: $(r)^+$
- Zero of one instances: $r^?$
- Character classes: $[abc]$

- Example:
 - `letter_` -> $[A-Za-z_]$
 - `digit` -> $[0-9]$
 - `id` -> $\text{letter_}(\text{letter}|\text{digit})^*$

Recognition of tokens

- Starting point is the language grammar to understand the tokens:

stmt -> **if** expr **then** stmt

| **if** expr **then** stmt **else** stmt

| ϵ

expr -> term **relop** term

| term

term -> **id**

| **number**

Recognition of tokens (cont.)

- The next step is to formalize the patterns:

digit -> [0-9]

Digits -> digit+

number -> digit(.digits)? (E[+-]? Digit)?

letter -> [A-Za-z_]

id -> letter (letter|digit)*

If -> if

Then -> then

Else -> else

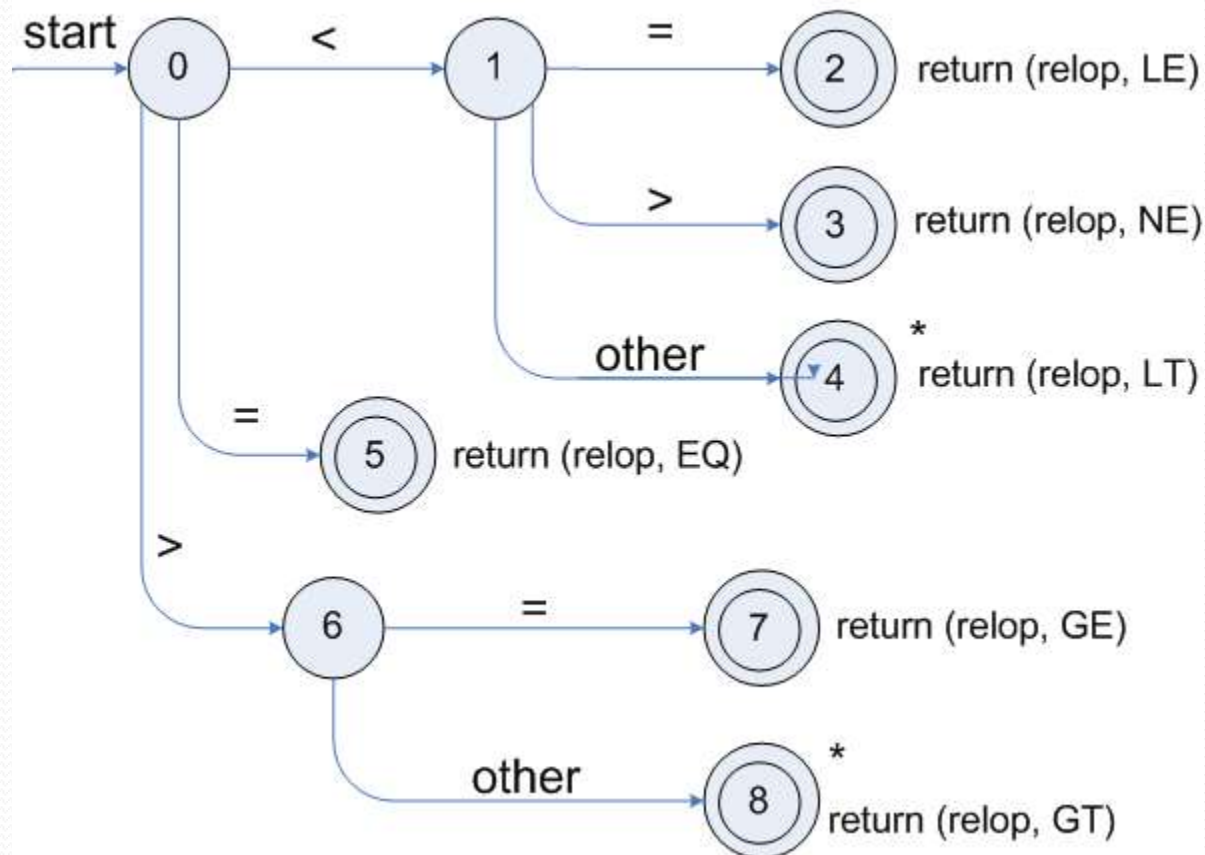
Relop -> < | > | <= | >= | = | <>

- We also need to handle whitespaces:

ws -> (blank | tab | newline)+

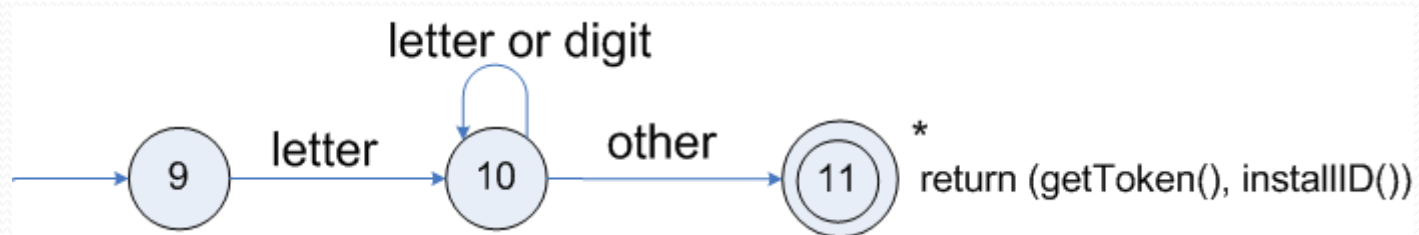
Transition diagrams

- Transition diagram for relop



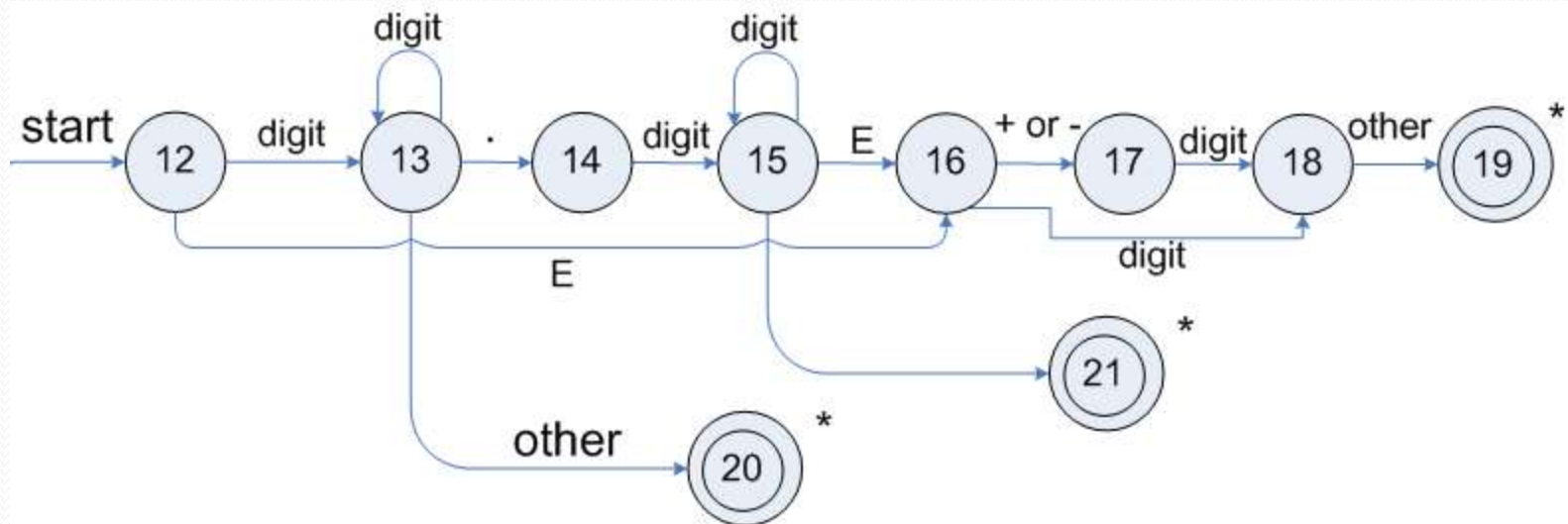
Transition diagrams (cont.)

- Transition diagram for reserved words and identifiers



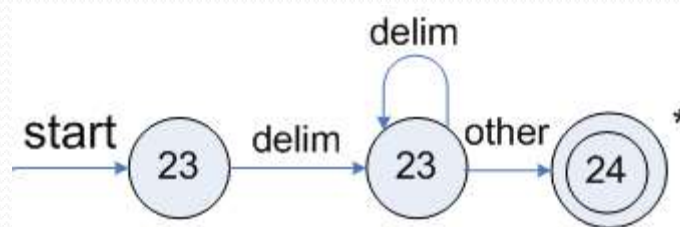
Transition diagrams (cont.)

- Transition diagram for unsigned numbers



Transition diagrams (cont.)

- Transition diagram for whitespace



Architecture of a transition-diagram-based lexical analyzer

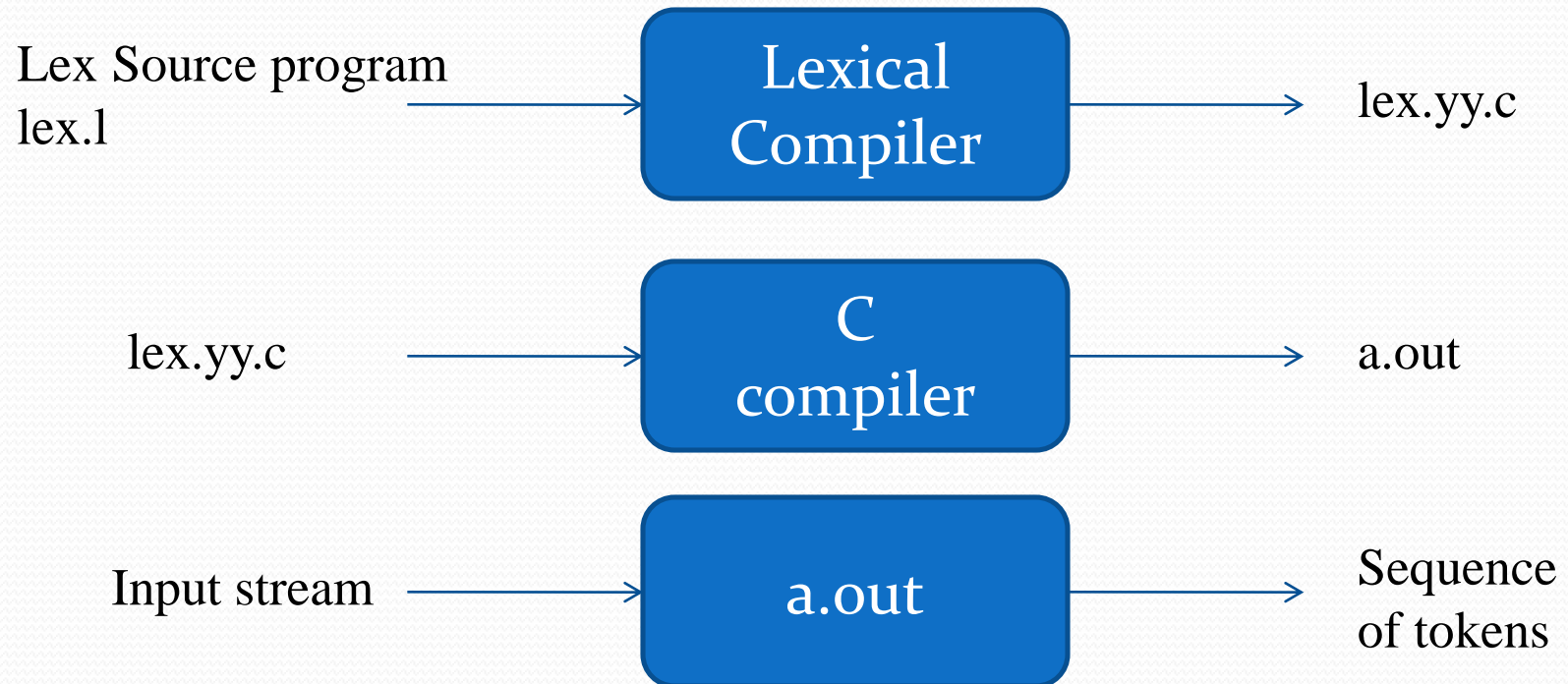
```
TOKEN getRelop()
{
    TOKEN retToken = new (RELOP)
    while (1) {          /* repeat character processing until a
                        return or failure occurs */
        switch(state) {
            case 0: c= nextchar();
                    if (c == '<') state = 1;
                    else if (c == '=') state = 5;
                    else if (c == '>') state = 6;
                    else fail(); /* lexeme is not a relop */
                    break;

            case 1: ...

            ...

            case 8: retract();
                    retToken.attribute = GT;
                    return(retToken);
        }
    }
```

Lexical Analyzer Generator - Lex



Structure of Lex programs

declarations

%%

translation rules



Pattern {Action}

%%

auxiliary functions

Example

```
%{
    /* definitions of manifest constants
    LT, LE, EQ, NE, GT, GE,
    IF, THEN, ELSE, ID, NUMBER, RELOP */
}%

/* regular definitions
delim      [ \t\n]
ws         {delim}+
letter     [A-Za-z]
digit      [0-9]
id         {letter}({letter}|{digit})*
number     {digit}+(\.{digit}+)?(E[+-]?{digit}+)?

%%
{ws}       { /* no action and no return */}
if         {return(IF);}
then       {return(THEN);}
else       {return(ELSE);}
{id}       {yylval = (int) installID(); return(ID); }
{number}   {yylval = (int) installNum(); return(NUMBER);}
...
```

```
Int installID() { /* funtion to install the
lexeme, whose first character is
pointed to by yytext, and whose
length is yyleng, into the symbol
table and return a pointer thereto
*/
```

```
}
```

```
Int installNum() { /* similar to
installID, but puts numerical
constants into a separate table */
```

```
}
```

Finite Automata

- Regular expressions = specification
- Finite automata = implementation

- A finite automaton consists of
 - An input alphabet Σ
 - A set of states S
 - A start state n
 - A set of accepting states $F \subseteq S$
 - A set of transitions $\text{state} \rightarrow^{\text{input}} \text{state}$

Finite Automata

- Transition

$$s_1 \xrightarrow{a} s_2$$

- Is read

In state s_1 on input “a” go to state s_2

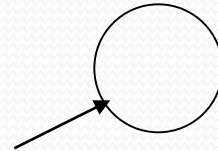
- If end of input
 - If in accepting state => accept, otherwise => reject
- If no transition possible => reject

Finite Automata State Graphs

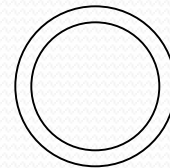
- A state



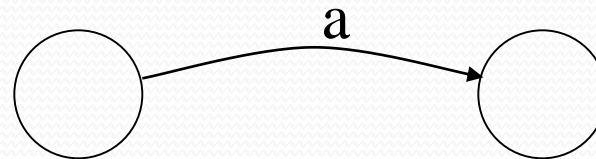
- The start state



- An accepting state

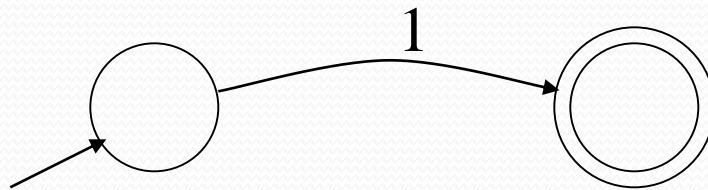


- A transition



A Simple Example

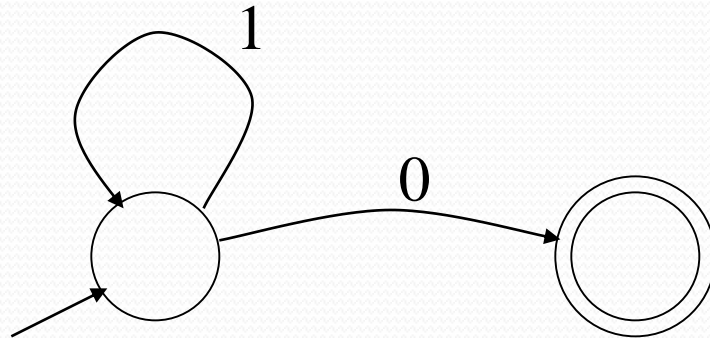
- A finite automaton that accepts only “1”



- A finite automaton accepts a string if we can follow transitions labeled with the characters in the string from the start to some accepting state

Another Simple Example

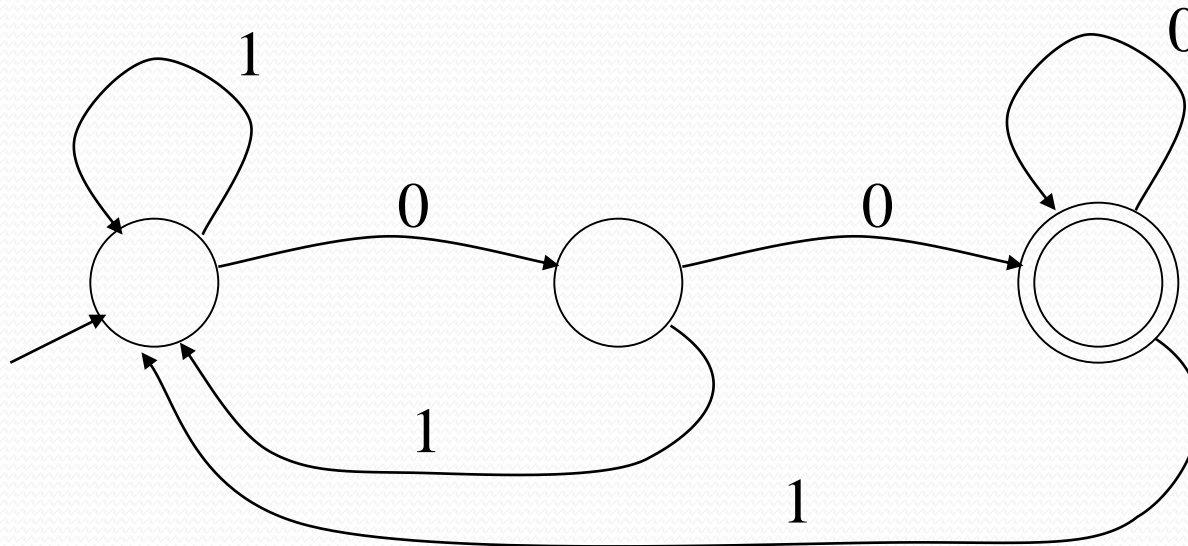
- A finite automaton accepting any number of 1's followed by a single 0
- Alphabet: $\{0,1\}$



- Check that “1110” is accepted but “110...” is not

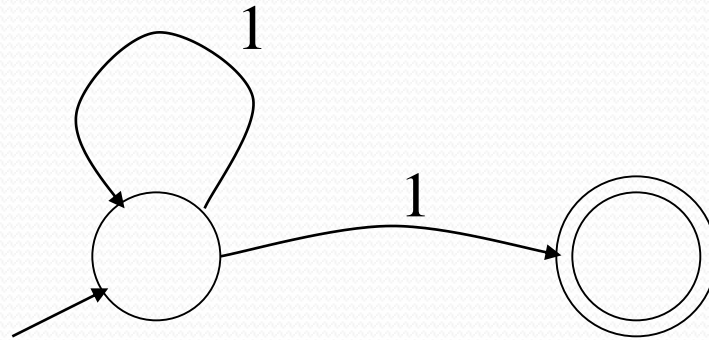
And Another Example

- Alphabet $\{0,1\}$
- What language does this recognize?



And Another Example

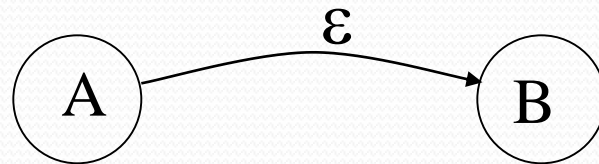
- Alphabet still $\{ 0, 1 \}$



- The operation of the automaton is not completely defined by the input
 - On input “11” the automaton could be in either state

Epsilon Moves

- Another kind of transition: ϵ -moves



- Machine can move from state A to state B without reading input

Deterministic and Nondeterministic Automata

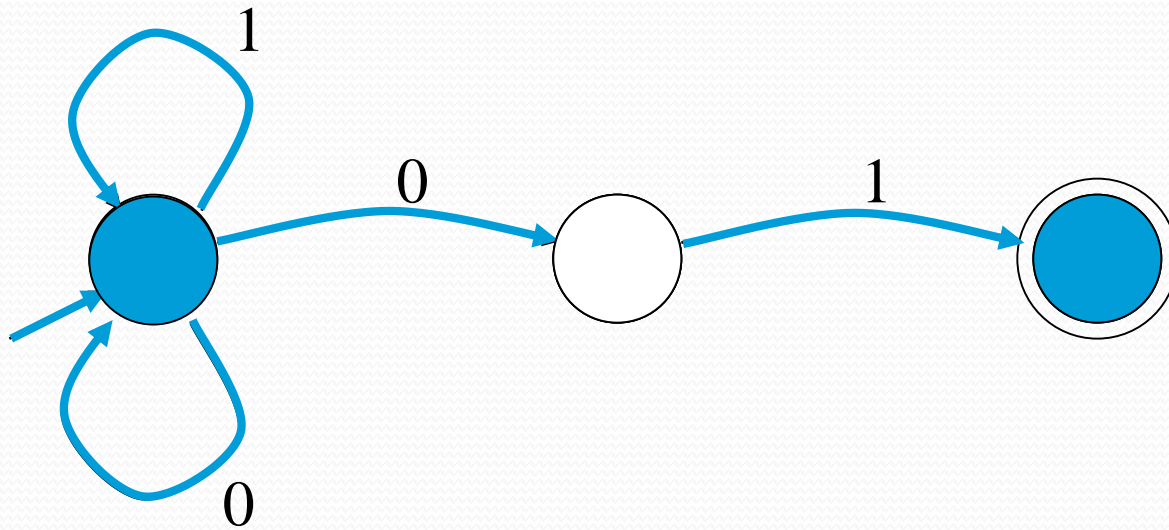
- Deterministic Finite Automata (DFA)
 - One transition per input per state
 - No ϵ -moves
- Nondeterministic Finite Automata (NFA)
 - Can have multiple transitions for one input in a given state
 - Can have ϵ -moves
- *Finite* automata have *finite* memory
 - Need only to encode the current state

Execution of Finite Automata

- A DFA can take only one path through the state graph
 - Completely determined by input
- NFAs can choose
 - Whether to make ϵ -moves
 - Which of multiple transitions for a single input to take

Acceptance of NFAs

- An NFA can get into multiple states



- Input: 1 0 1
- Rule: NFA accepts if it can get in a final state

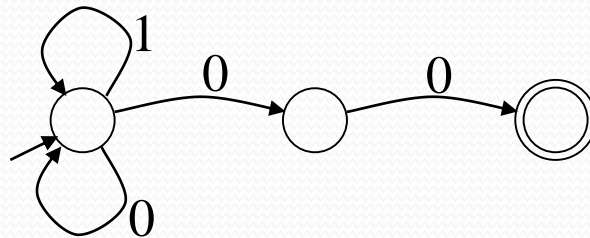
NFA vs. DFA (1)

- NFAs and DFAs recognize the same set of languages (regular languages)
- DFAs are easier to implement
 - There are no choices to consider

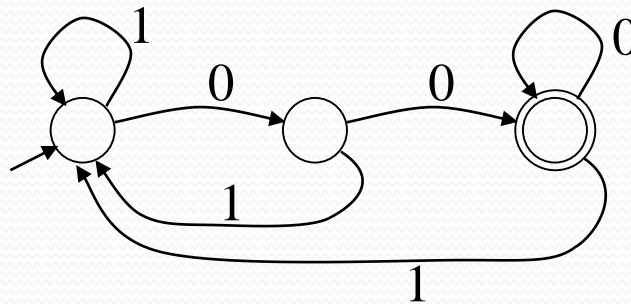
NFA vs. DFA (2)

- For a given language the NFA can be simpler than the DFA

NFA



DFA

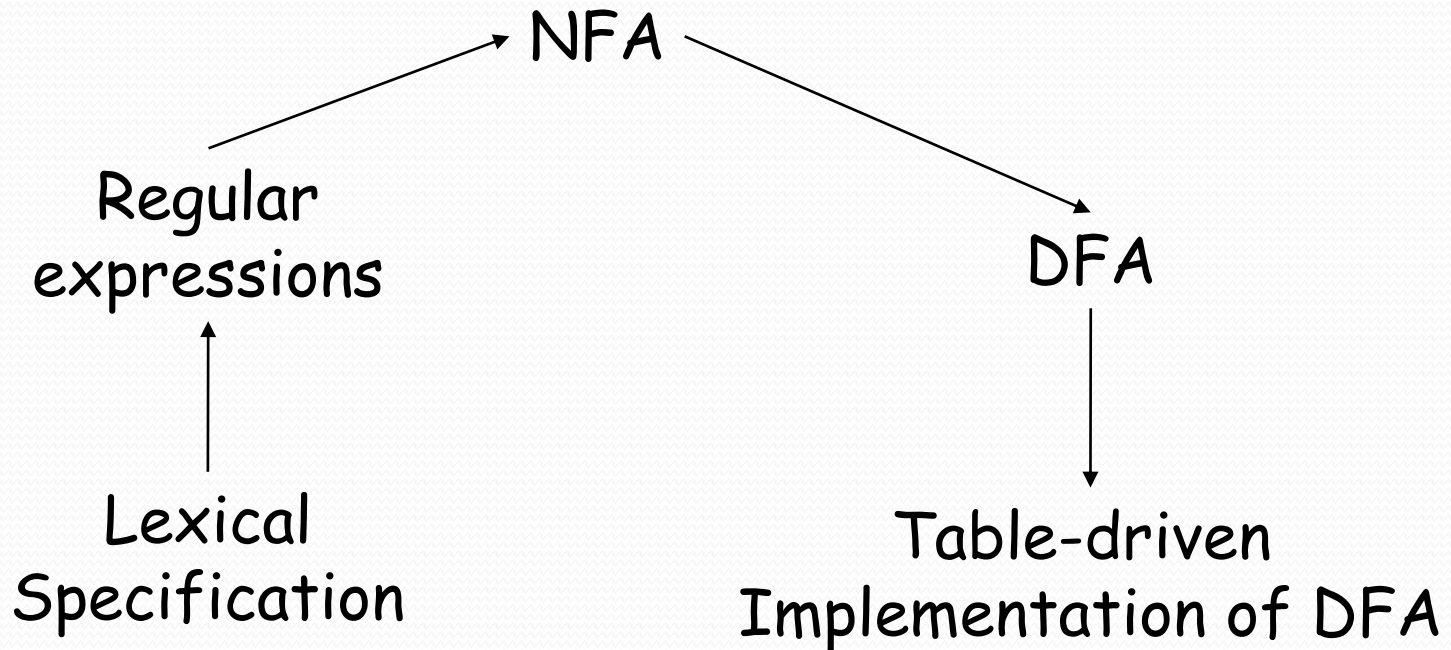


- DFA can be exponentially larger than NFA

Regular Expressions to Finite

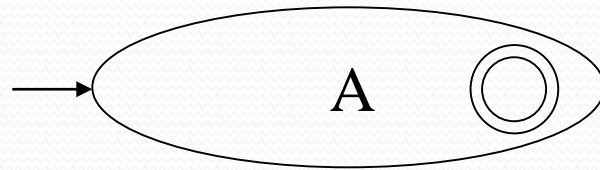
Automata

- High-level sketch

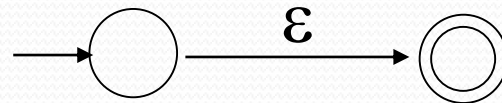


Regular Expressions to NFA (1)

- For each kind of rexp, define an NFA
 - Notation: NFA for rexp A



- For ε

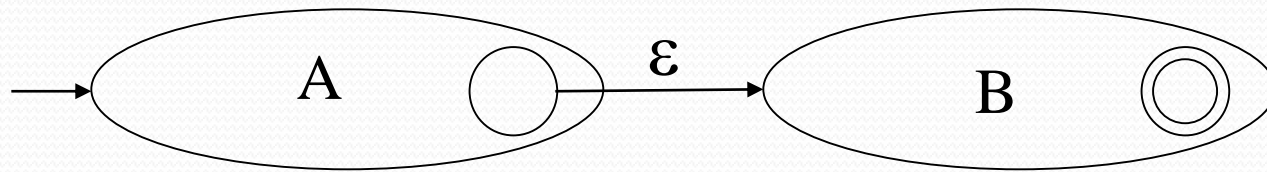


- For input a

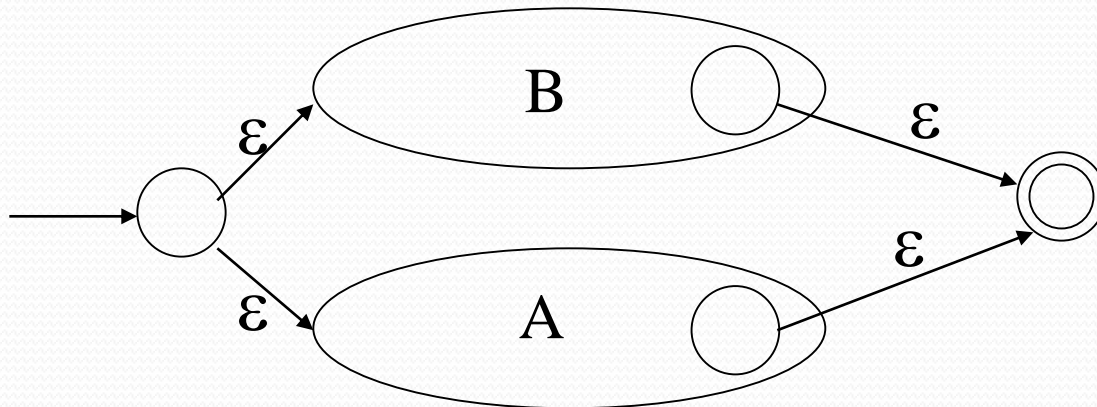


Regular Expressions to NFA (2)

- For AB

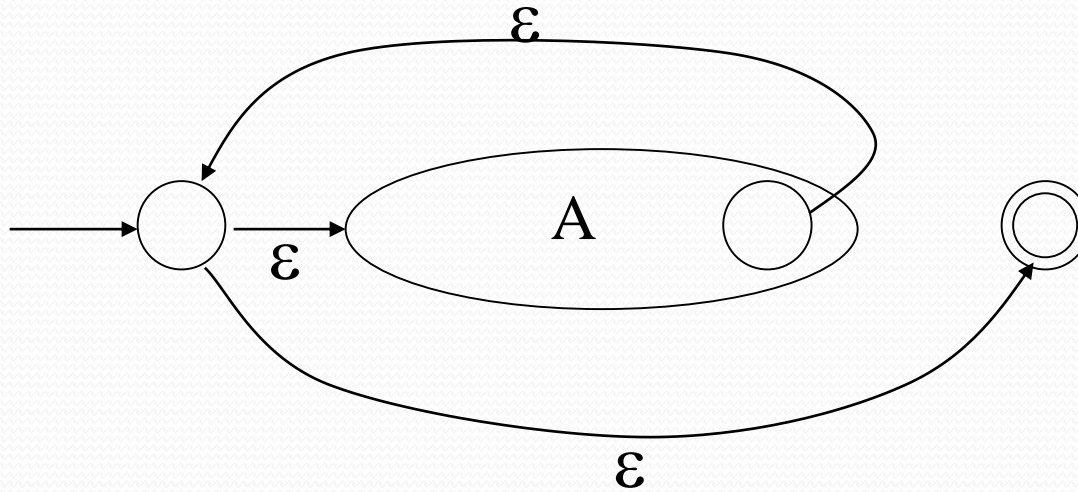


- For $A \mid B$



Regular Expressions to NFA (3)

- For A^*

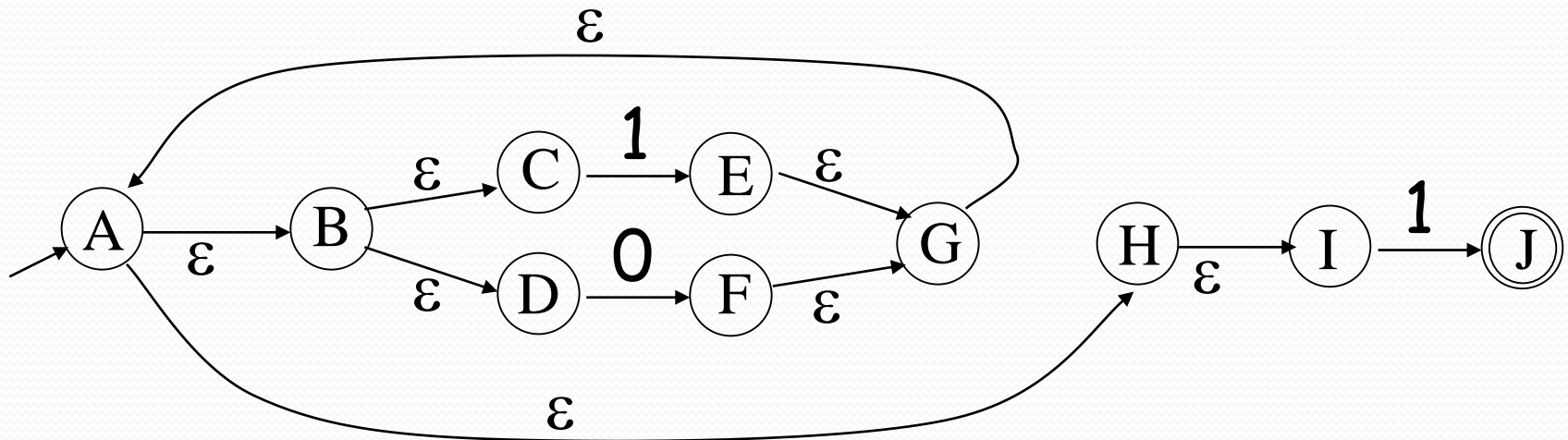


Example of RegExp \rightarrow NFA conversion

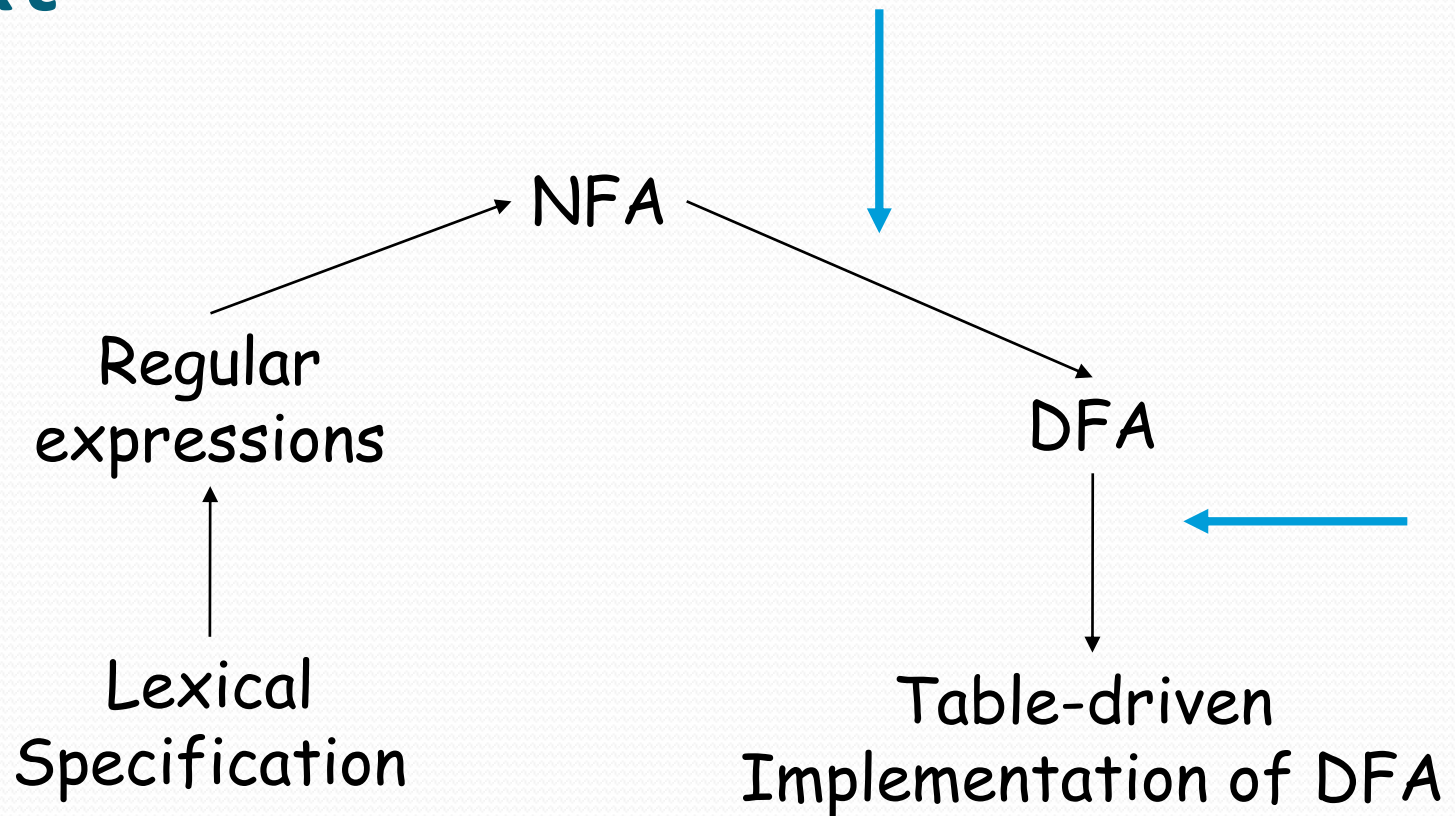
- Consider the regular expression

$$(1 \mid 0)^*1$$

- The NFA is



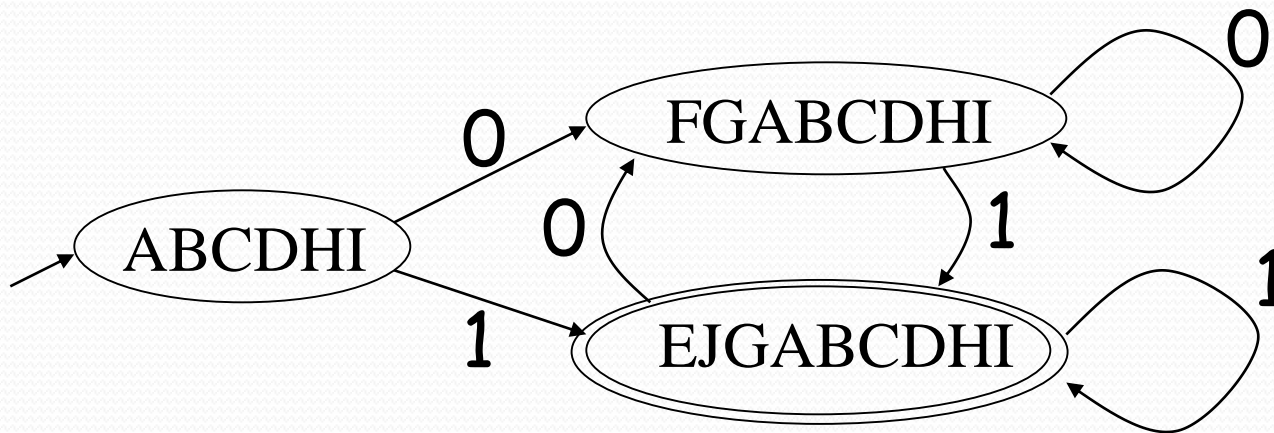
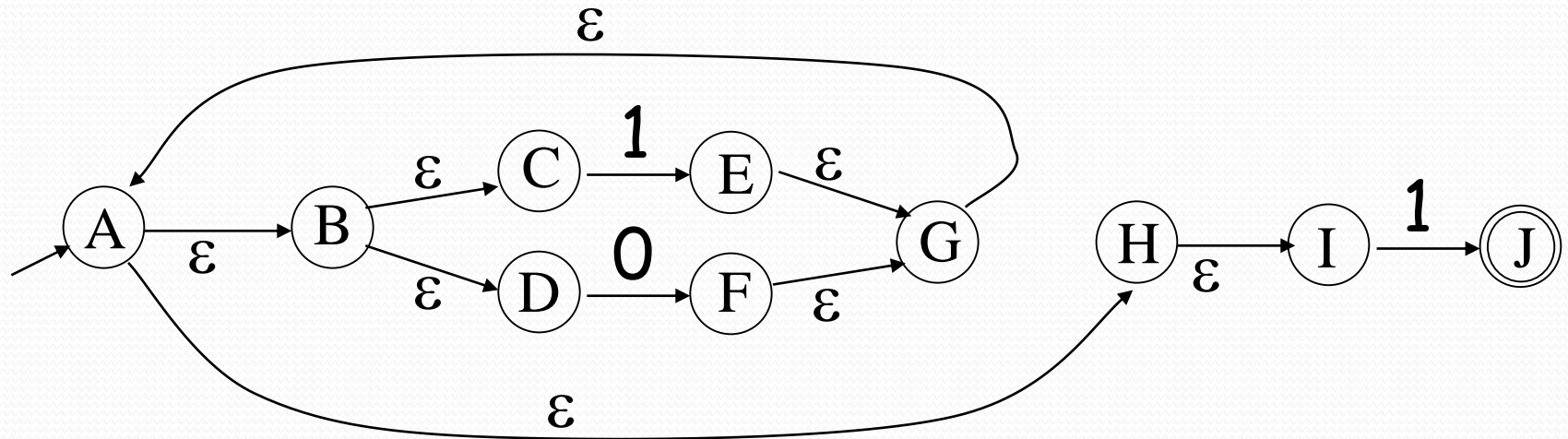
Next



NFA to DFA. The Trick

- Simulate the NFA
- Each state of resulting DFA
 - = a non-empty subset of states of the NFA
- Start state
 - = the set of NFA states reachable through ϵ -moves from NFA start state
- Add a transition $S \xrightarrow{a} S'$ to DFA iff
 - S' is the set of NFA states reachable from the states in S after seeing the input a
 - considering ϵ -moves as well

NFA \rightarrow DFA Example



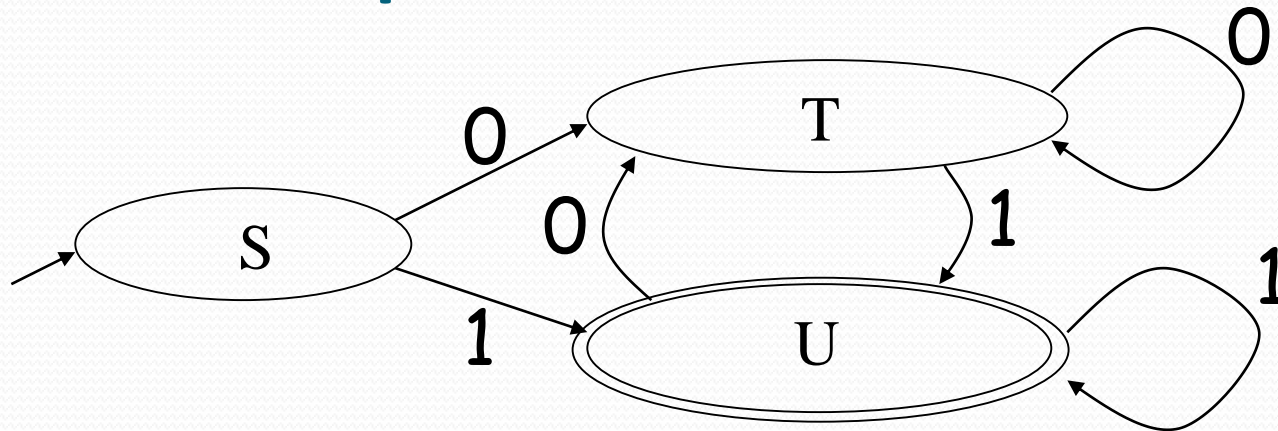
NFA to DFA. Remark

- An NFA may be in many states at any time
- How many different states ?
- If there are N states, the NFA must be in some subset of those N states
- How many non-empty subsets are there?
 - $2^N - 1 =$ finitely many, but exponentially many

Implementation

- A DFA can be implemented by a 2D table T
 - One dimension is “states”
 - Other dimension is “input symbols”
 - For every transition $S_i \xrightarrow{a} S_k$ define $T[i,a] = k$
- DFA “execution”
 - If in state S_i and input a , read $T[i,a] = k$ and skip to state S_k
 - Very efficient

Table Implementation of a DFA



	0	1
S	T	U
T	T	U
U	T	U

Implementation (Cont.)

- NFA \rightarrow DFA conversion is at the heart of tools such as flex or jflex
- But, DFAs can be huge
- In practice, flex-like tools trade off speed for space in the choice of NFA and DFA representations

Readings

- Chapter 3 of the book

Compiler course

Chapter 8
Code Generation

Outline

- Code Generation Issues
- Target language Issues
- Addresses in Target Code
- Basic Blocks and Flow Graphs
- Optimizations of Basic Blocks
- A Simple Code Generator
- Peephole optimization
- Register allocation and assignment
- Instruction selection by tree rewriting

Introduction

- The final phase of a compiler is code generator
- It receives an intermediate representation (IR) with supplementary information in symbol table
- Produces a semantically equivalent target program
- Code generator main tasks:
 - Instruction selection
 - Register allocation and assignment
 - Instruction ordering



Issues in the Design of Code Generator

- The most important criterion is that it produces correct code
- Input to the code generator
 - IR + Symbol table
 - We assume front end produces low-level IR, i.e. values of names in it can be directly manipulated by the machine instructions.
 - Syntactic and semantic errors have been already detected
- The target program
 - Common target architectures are: RISC, CISC and Stack based machines
 - In this chapter we use a very simple RISC-like computer with addition of some CISC-like addressing modes

complexity of mapping

- the level of the IR
- the nature of the instruction-set architecture
- the desired quality of the generated code.

$x=y+z$

```
LD    R0, y
ADD   R0, R0, z
ST    x, R0
```

$a=b+c$

$d=a+e$

```
LD    R0, b
ADD   R0, R0, c
ST    a, R0
LD    R0, a
ADD   R0, R0, e
ST    d, R0
```

Register allocation

- Two subproblems
 - Register allocation: selecting the set of variables that will reside in registers at each point in the program
 - Register assignment: selecting specific register that a variable reside in
- Complications imposed by the hardware architecture
 - Example: register pairs for multiplication and division

t=a+b
t=t*c
T=t/d

L	R1, a
A	R1, b
M	R0, c
D	R0, d
ST	R1, t

t=a+b
t=t+c
T=t/d

L	R0, a
A	R0, b
M	R0, c
SRDA	R0, 32
D	R0, d
ST	R1, t

A simple target machine model

- Load operations: LD r,x and LD r1, r2
- Store operations: ST x,r
- Computation operations: OP dst, src1, src2
- Unconditional jumps: BR L
- Conditional jumps: Bcond r, L like BLTZ r, L

Addressing Modes

- variable name: x
- indexed address: $a(r)$ like $LD R_1, a(R_2)$ means $R_1 = \text{contents}(a + \text{contents}(R_2))$
- integer indexed by a register : like $LD R_1, 100(R_2)$
- Indirect addressing mode: $*r$ and $*100(r)$
- immediate constant addressing mode: like $LD R_1, \#100$

b = a [i]

LD R1, i //R1 = i

MUL R1, R1, 8 //R1 = R1 * 8

LD R2, a(R1)

//R2=contents(a+contents(R1))

ST b, R2 //b = R2

a[j] = c

LD R1, c //R1 = c

LD R2, j // R2 = j

MUL R2, R2, 8 //R2 = R2 * 8

ST a(R2), R1

//contents(a+contents(R2))=R1

$x = *p$

LD R1, p // **R1 = p**

LD R2, o(R1) // **R2 =**
contents(o+contents(R1))

ST x, R2 // **x=R2**

conditional-jump three-address instruction

If $x < y$ goto L

LD R₁, x // R₁ = x

LD R₂, y // R₂ = y

SUB R₁, R₁, R₂ // R₁ = R₁ - R₂

BLTZ R₁, M // i f R₁ < 0 jump t o M

costs associated with the addressing modes

- LD R₀, R₁ cost = 1
- LD R₀, M cost = 2
- LD R₁, *₁₀₀(R₂) cost = 3

Addresses in the Target Code

- A statically determined area Code
- A statically determined data area Static
- A dynamically managed area Heap
- A dynamically managed area Stack

three-address statements for procedure calls and returns

- call callee
- Return
- Halt
- action

Target program for a sample call and return

```
    // code for c
action1
call p
action2
halt

    // code for p
action3
return
```

```
100: ACTION1           // code for c
120: ST 364, #140      // code for action1
132: BR 200           // save return address 140 in location 364
140: ACTION2         // call p
160: HALT             // return to operating system
    ...
    // code for p
200: ACTION3
220: BR *364         // return to address saved in location 364
    ...
    // 300-363 hold activation record for c
300: // return address
304: // local data for c
    ...
    // 364-451 hold activation record for p
364: // return address
368: // local data for p
```

Stack Allocation

```
LD    SP, #stackStart           // initialize the stack
code for the first procedure
HALT                                // terminate execution

ADD   SP, SP, #caller.recordSize // increment stack pointer
ST    *SP, #here + 16           // save return address
BR    callee.codeArea         // return to caller
                                     Branch to called procedure
```

Return to caller

in Callee: BR *0(SP)

in caller: SUB SP, SP, #*caller.recordsize*

Target code for stack allocation

```
                // code for m
action1
call q
action2
halt

                // code for p
action3
return

                // code for q
action4
call p
action5
call q
action6
call q
return
```

```
100: LD SP, #600           // code for m
108: ACTION1             // initialize the stack
128: ADD SP, SP, #msize   // code for action1
136: ST *SP, #152        // call sequence begins
144: BR 300              // push return address
152: SUB SP, SP, #msize  // call q
160: ACTION12           // restore SP
180: HALT
...
200: ACTION3             // code for p
220: BR *0(SP)          // return
...
300: ACTION4             // code for q
320: ADD SP, SP, #qsize // contains a conditional jump to 456
328: ST *SP, #344       // push return address
336: BR 200              // call p
344: SUB SP, SP, #qsize
352: ACTION5
372: ADD SP, SP, #qsize
380: BR *SP, #396       // push return address
388: BR 300              // call q
396: SUB SP, SP, #qsize
404: ACTION6
424: ADD SP, SP, #qsize
432: ST *SP, #440       // push return address
440: BR 300              // call q
448: SUB SP, SP, #qsize
456: BR *0(SP)          // return
...
600:                      // stack starts here
```

Basic blocks and flow graphs

- Partition the intermediate code into basic blocks
 - The flow of control can only enter the basic block through the first instruction in the block. That is, there are no jumps into the middle of the block.
 - Control will leave the block without halting or branching, except possibly at the last instruction in the block.
- The basic blocks become the nodes of a flow graph

rules for finding leaders

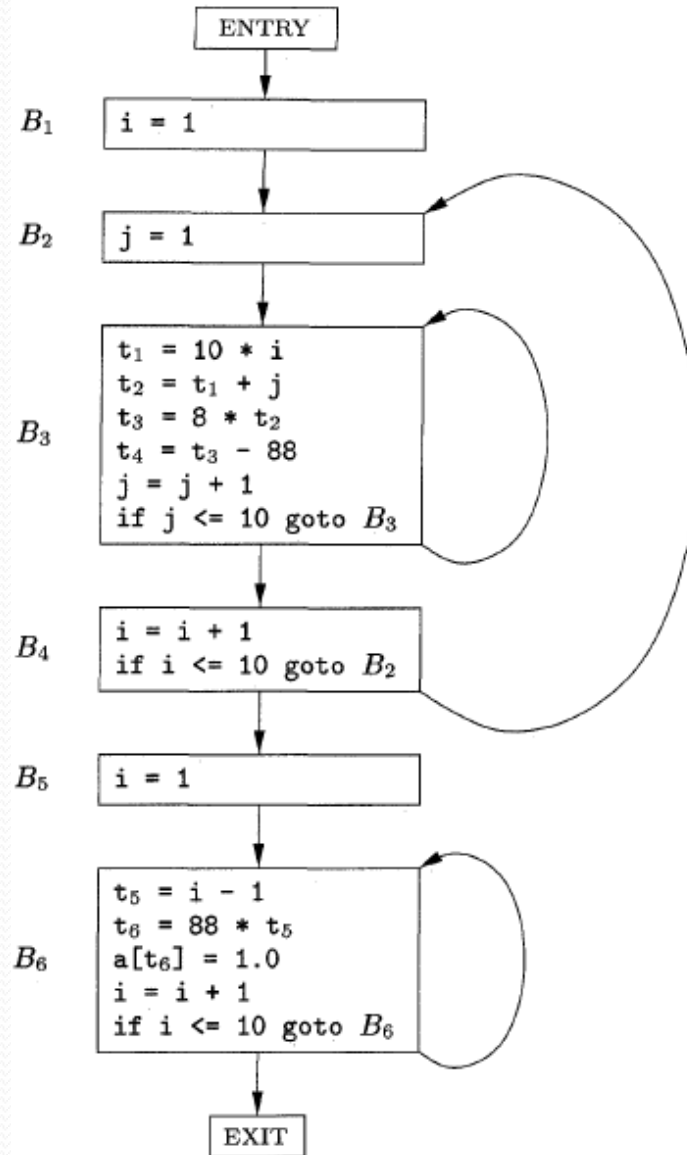
- The first three-address instruction in the intermediate code is a leader.
- Any instruction that is the target of a conditional or unconditional jump is a leader.
- Any instruction that immediately follows a conditional or unconditional jump is a leader.

Intermediate code to set a 10*10 matrix to an identity matrix

```
for i from 1 to 10 do
    for j from 1 to 10 do
        a[i, j] = 0.0;
for i from 1 to 10 do
    a[i, i] = 1.0;
```

```
1)  i = 1
2)  j = 1
3)  t1 = 10 * i
4)  t2 = t1 + j
5)  t3 = 8 * t2
6)  t4 = t3 - 88
7)  a[t4] = 0.0
8)  j = j + 1
9)  if j <= 10 goto (3)
10) i = i + 1
11) if i <= 10 goto (2)
12) i = 1
13) t5 = i - 1
14) t6 = 88 * t5
15) a[t6] = 1.0
16) i = i + 1
17) if i <= 10 goto (13)
```

Flow graph based on Basic Blocks



liveness and next-use information

- We wish to determine for each three address statement $x=y+z$ what the next uses of x , y and z are.
- Algorithm:
 - Attach to statement i the information currently found in the symbol table regarding the next use and liveness of x , y , and z .
 - In the symbol table, set x to "not live" and "no next use."
 - In the symbol table, set y and z to "live" and the next uses of y and z to i .

DAG representation of basic blocks

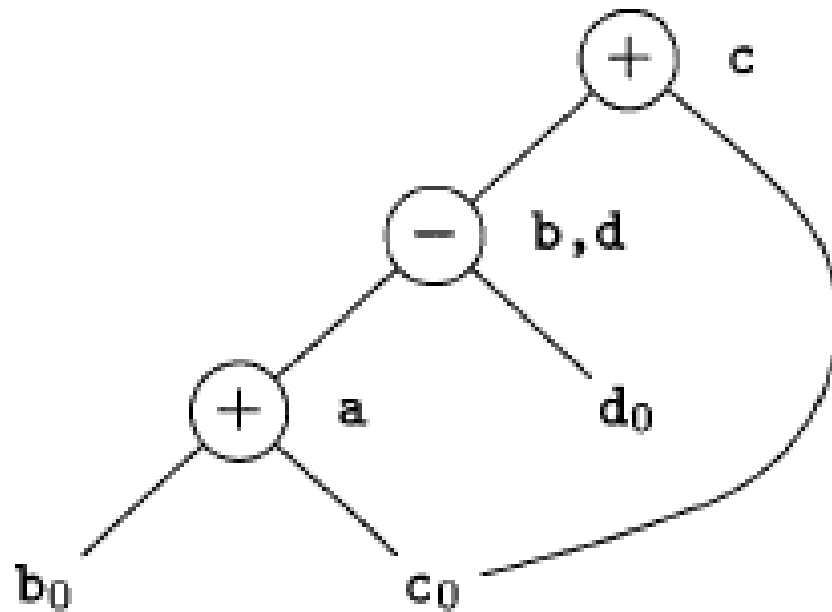
- There is a node in the DAG for each of the initial values of the variables appearing in the basic block.
- There is a node N associated with each statement s within the block. The children of N are those nodes corresponding to statements that are the last definitions, prior to s , of the operands used by s .
- Node N is labeled by the operator applied at s , and also attached to N is the list of variables for which it is the last definition within the block.
- Certain nodes are designated *output nodes*. These are the nodes whose variables are *live on exit* from the block.

Code improving transformations

- We can eliminate *local common subexpressions*, that is, instructions that compute a value that has already been computed.
- We can eliminate *dead code*, that is, instructions that compute a value that is never used.
- We can reorder statements that do not depend on one another; such reordering may reduce the time a temporary value needs to be preserved in a register.
- We can apply algebraic laws to reorder operands of three-address instructions, and sometimes thereby simplify the computation.

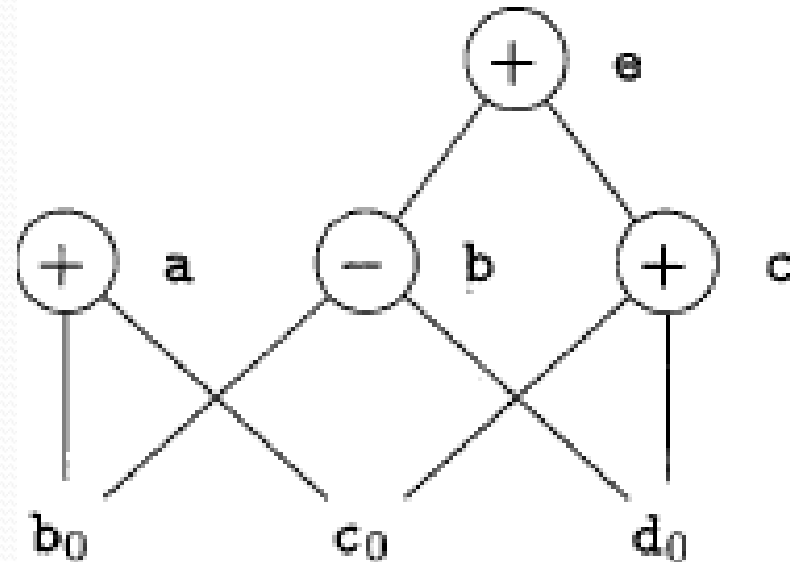
DAG for basic block

```
a = b + c
b = a - d
c = b + c
d = a - d
```



DAG for basic block

```
a = b + c;  
b = b - d  
c = c + d  
e = b + c
```

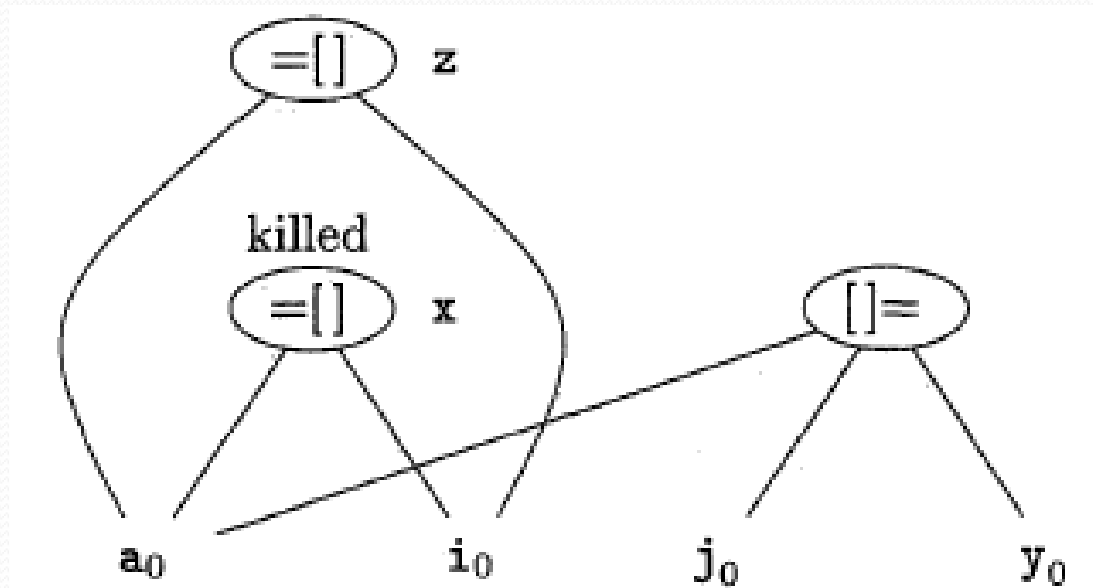


array accesses in a DAG

- An assignment from an array, like $x = a[i]$, is represented by creating a node with operator $=[]$ and two children representing the initial value of the array, ao in this case, and the index i . Variable x becomes a label of this new node.
- An assignment to an array, like $a[j] = y$, is represented by a new node with operator $[]=$ and three children representing ao , j and y . There is no variable labeling this node. What is different is that the creation of this node *kills* all currently constructed nodes whose value depends on ao . A node that has been killed cannot receive any more labels; that is, it cannot become a common subexpression.

DAG for a sequence of array assignments

```
x = a[i]  
a[j] = y  
z = a[i]
```



Rules for reconstructing the basic block from a DAG

- The order of instructions must respect the order of nodes in the DAG. That is, we cannot compute a node's value until we have computed a value for each of its children.
- Assignments to an array must follow all previous assignments to, or evaluations from, the same array, according to the order of these instructions in the original basic block.
- Evaluations of array elements must follow any previous (according to the original block) assignments to the same array. The only permutation allowed is that two evaluations from the same array may be done in either order, as long as neither crosses over an assignment to that array.
- Any use of a variable must follow all previous (according to the original block) procedure calls or indirect assignments through a pointer.
- Any procedure call or indirect assignment through a pointer must follow all previous (according to the original block) evaluations of any variable.

principal uses of registers

- In most machine architectures, some or all of the operands of an operation must be in registers in order to perform the operation.
- Registers make good temporaries - places to hold the result of a subexpression while a larger expression is being evaluated, or more generally, a place to hold a variable that is used only within a single basic block.
- Registers are often used to help with run-time storage management, for example, to manage the run-time stack, including the maintenance of stack pointers and possibly the top elements of the stack itself.

Descriptors for data structure

- For each available register, a register descriptor keeps track of the variable names whose current value is in that register. Since we shall use only those registers that are available for local use within a basic block, we assume that initially, all register descriptors are empty. As the code generation progresses, each register will hold the value of zero or more names.
- For each program variable, an address descriptor keeps track of the location or locations where the current value of that variable can be found. The location might be a register, a memory address, a stack location, or some set of more than one of these. The information can be stored in the symbol-table entry for that variable name.

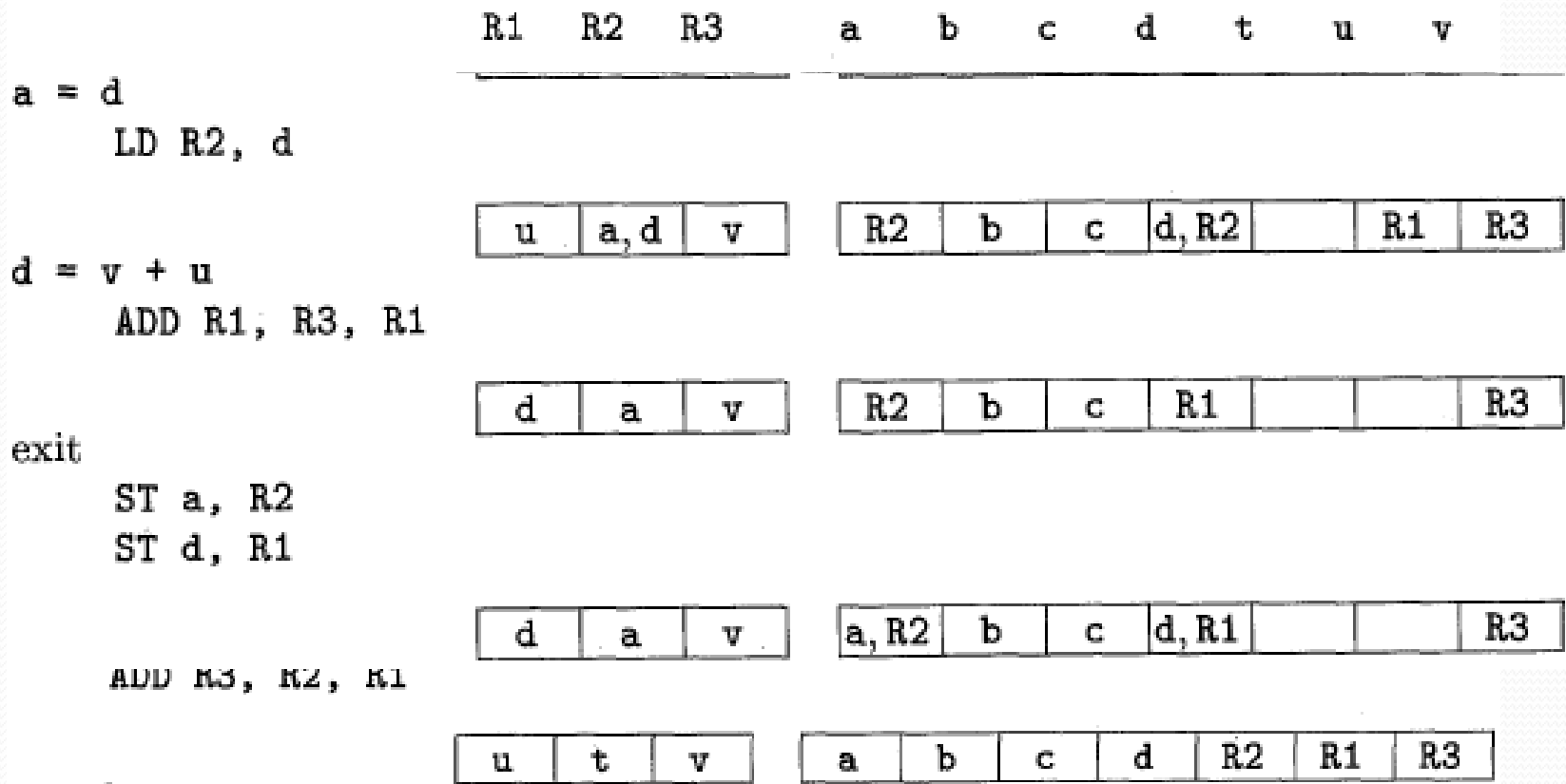
Machine Instructions for Operations

- Use `getReg(x = y + z)` to select registers for x , y , and z . Call these R_x , R_y and R_z .
- If y is not in R_y (according to the register descriptor for R_y), then issue an instruction `LD R_y , y'` , where y' is one of the memory locations for y (according to the address descriptor for y).
- Similarly, if z is not in R_z , issue and instruction `LD R_z , z'` , where z' is a location for x .
- Issue the instruction `ADD R_x , R_y , R_z` .

Rules for updating the register and address descriptors

- For the instruction LD R, x
 - Change the register descriptor for register R so it holds only x.
 - Change the address descriptor for x by adding register R as an additional location.
- For the instruction ST x, R, change the address descriptor for x to include its own memory location.
- For an operation such as ADD R_x, R_y, R_z implementing a three-address instruction $x = y + z$
 - Change the register descriptor for R_x so that it holds only x.
 - Change the address descriptor for x so that its only location is R_x . Note that the memory location for x is *not* now in the address descriptor for x.
 - Remove R_x from the address descriptor of any variable other than x.
- When we process a copy statement $x = y$, after generating the load for y into register R_y , if needed, and after managing descriptors as for all load statements (per rule I):
 - Add x to the register descriptor for R_y .
 - Change the address descriptor for x so that its only location is R_y .

Instructions generated and the changes in the register and address descriptors



Rules for picking register R_y for y

- If y is currently in a register, pick a register already containing y as R_y . Do not issue a machine instruction to load this register, as none is needed.
- If y is not in a register, but there is a register that is currently empty, pick one such register as R_y .
- The difficult case occurs when y is not in a register, and there is no register that is currently empty. We need to pick one of the allowable registers anyway, and we need to make it safe to reuse.

Possibilities for value of R

- If the address descriptor for v says that v is somewhere besides R , then we are OK.
- If v is x , the value being computed by instruction I , and x is not also one of the other operands of instruction I (z in this example), then we are OK. The reason is that in this case, we know this value of x is never again going to be used, so we are free to ignore it.
- Otherwise, if v is not used later (that is, after the instruction I , there are no further uses of v , and if v is live on exit from the block, then v is recomputed within the block), then we are OK.
- If we are not OK by one of the first two cases, then we need to generate the store instruction $ST\ v, R$ to place a copy of v in its own memory location. This operation is called a spill.

Selection of the register Rx

1. Since a new value of x is being computed, a register that holds only x is always an acceptable choice for R_x .
2. If y is not used after instruction I , and R_y holds only y after being loaded, R_y can also be used as R_x . A similar option holds regarding z and R_x .

Possibilities for value of R

- If the address descriptor for v says that v is somewhere besides R , then we are OK.
- If v is x , the value being computed by instruction I , and x is not also one of the other operands of instruction I (z in this example), then we are OK. The reason is that in this case, we know this value of x is never again going to be used, so we are free to ignore it.
- Otherwise, if v is not used later (that is, after the instruction I , there are no further uses of v , and if v is live on exit from the block, then v is recomputed within the block), then we are OK.
- If we are not OK by one of the first two cases, then we need to generate the store instruction $ST\ v, R$ to place a copy of v in its own memory location. This operation is called a spill.

Characteristic of peephole optimizations

- Redundant-instruction elimination
- Flow-of-control optimizations
- Algebraic simplifications
- Use of machine idioms

Redundant-instruction elimination

- LD a, R0
ST R0, a
- if debug == 1 goto L1
goto L2
L1 : print debugging information
L2:

Flow-of-control optimizations

goto L1

...

L1: goto L2

if a<b goto L1

...

L1: goto L2

Can be replaced by:

goto L2

...

L1: goto L2

Can be replaced by:

if a<b goto L2

...

L1: goto L2

Algebraic simplifications

- $X = X + 0$

- $X = X * 1$

Register Allocation and Assignment

- Global Register Allocation
- Usage Counts
- Register Assignment for Outer Loops
- Register Allocation by Graph Coloring

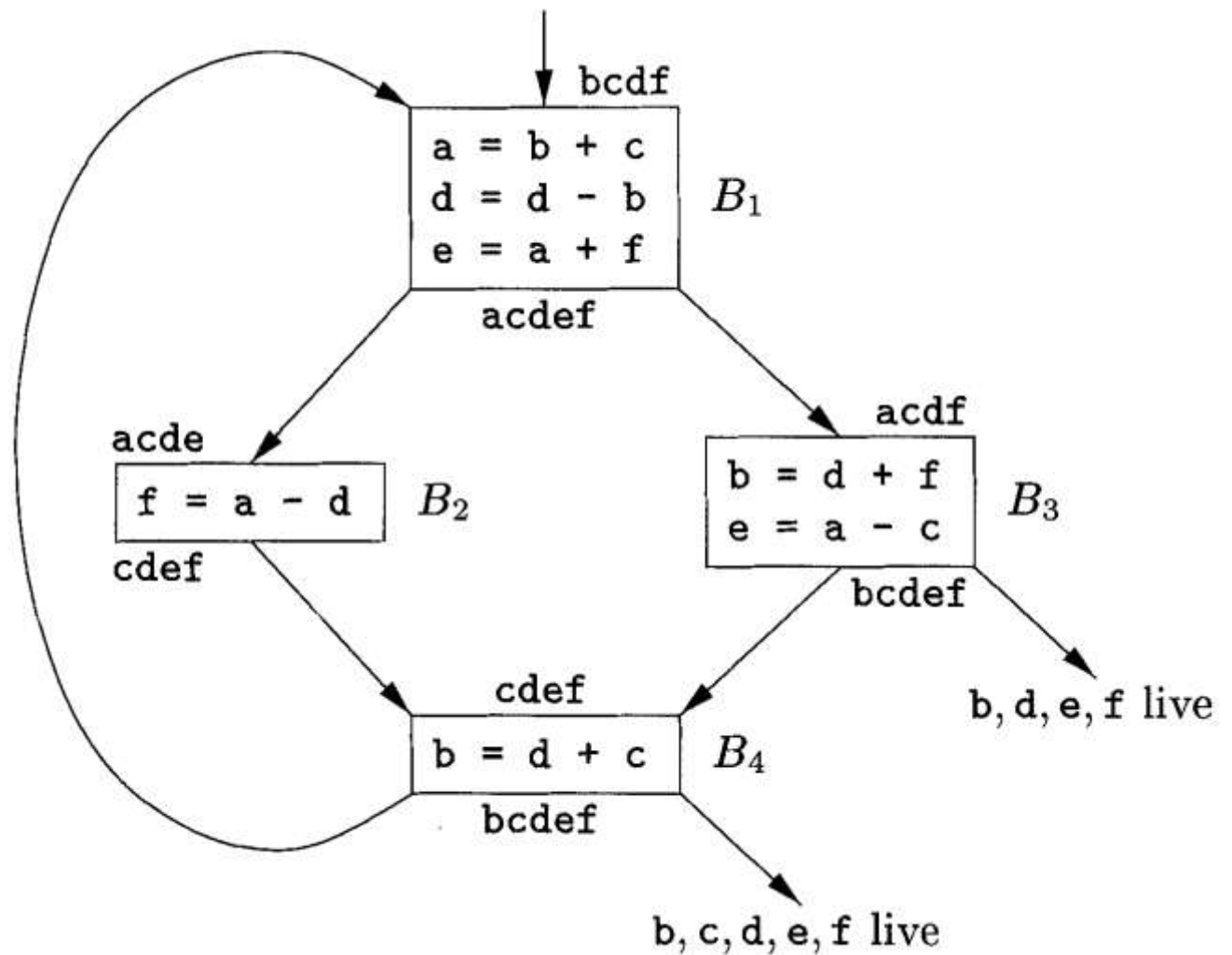
Global register allocation

- Previously explained algorithm does local (block based) register allocation
- This resulted that all live variables be stored at the end of block
- To save some of these stores and their corresponding loads, we might arrange to assign registers to frequently used variables and keep these registers consistent across block boundaries (globally)
- Some options are:
 - Keep values of variables used in loops inside registers
 - Use graph coloring approach for more globally allocation

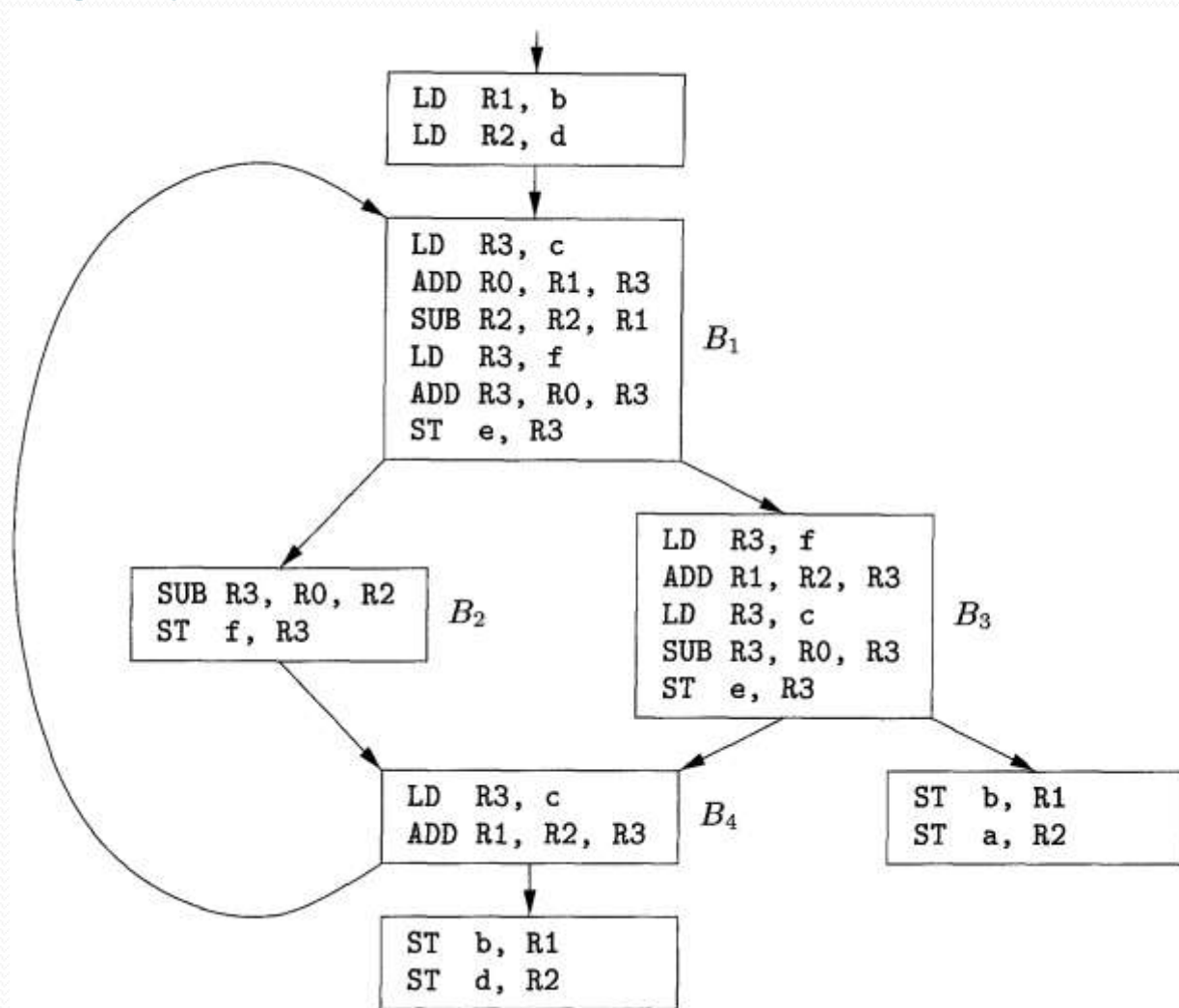
Usage counts

- For the loops we can approximate the saving by register allocation as:
 - Sum over all blocks (B) in a loop (L)
 - For each uses of x before any definition in the block we add one unit of saving
 - If x is live on exit from B and is assigned a value in B, then we ass 2 units of saving

Flow graph of an inner loop



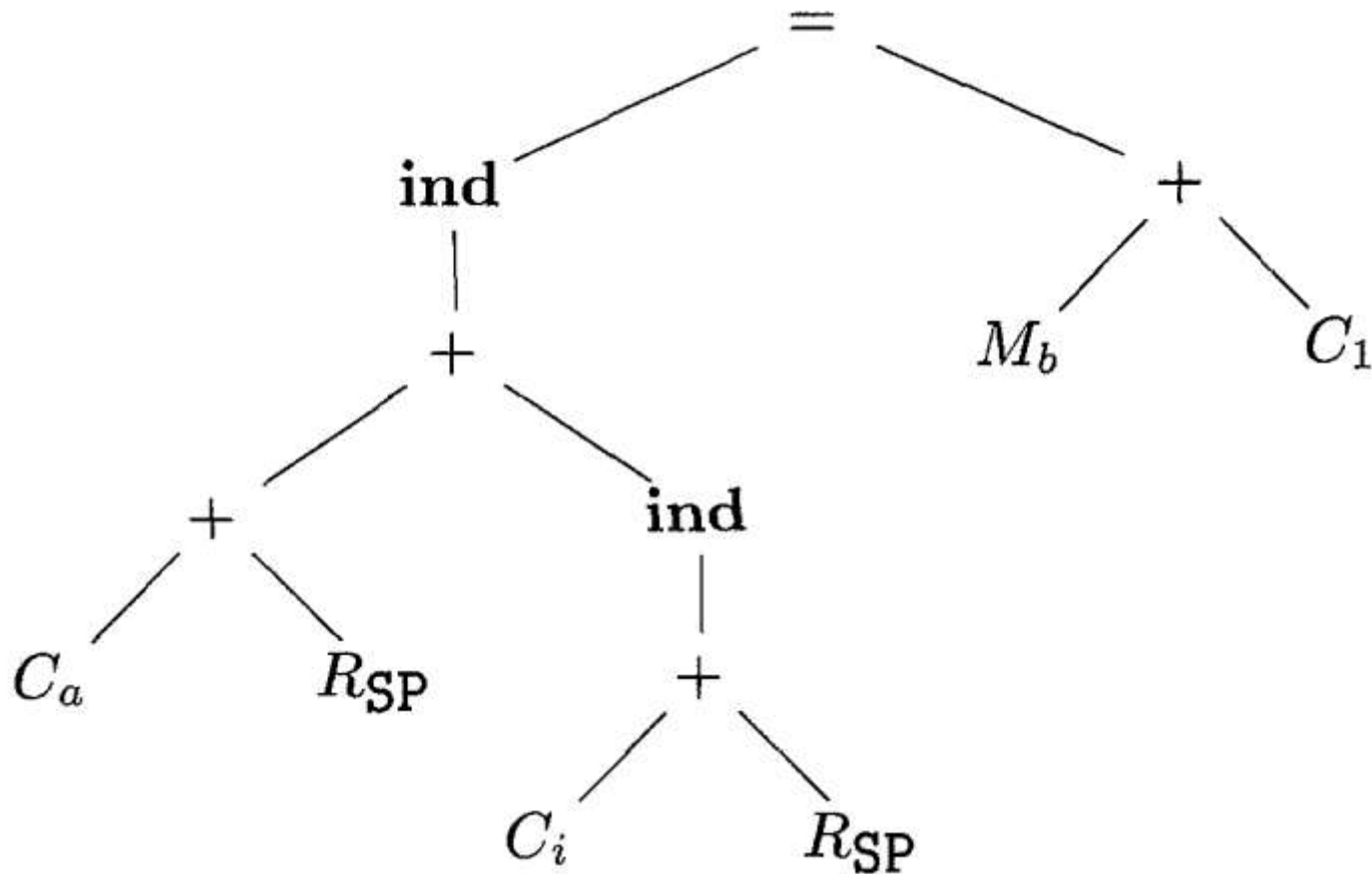
Code sequence using global register assignment



Register allocation by Graph coloring

- Two passes are used
 - Target-machine instructions are selected as though there are an infinite number of symbolic registers
 - Assign physical registers to symbolic ones
 - Create a register-interference graph
 - Nodes are symbolic registers and edges connects two nodes if one is live at a point where the other is defined.
 - For example in the previous example an edge connects a and d in the graph
 - Use a graph coloring algorithm to assign registers.

Intermediate-code tree for $a[i]=b+1$



Tree-rewriting rules

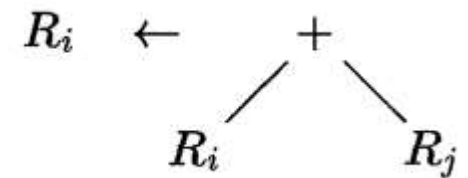
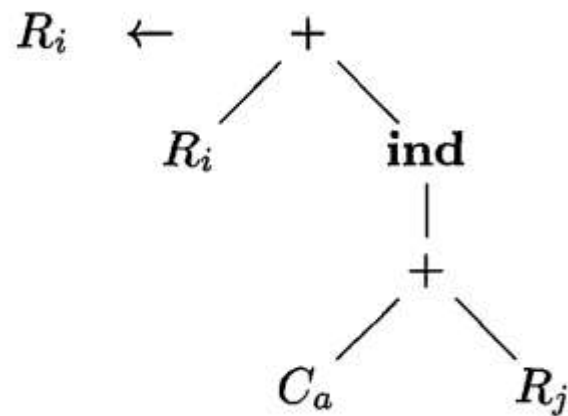
1)	$R_i \leftarrow C_a$	{ LD $R_i, \#a$ }
2)	$R_i \leftarrow M_x$	{ LD R_i, x }
3)	$M \leftarrow \begin{array}{c} = \\ / \quad \backslash \\ M_x \quad R_i \end{array}$	{ ST x, R_i }
4)	$M \leftarrow \begin{array}{c} = \\ / \quad \backslash \\ \mathbf{ind} \quad R_j \\ \\ R_i \end{array}$	{ ST $\#R_i, R_j$ }
5)	$R_i \leftarrow \begin{array}{c} \mathbf{ind} \\ \\ + \\ / \quad \backslash \\ C_a \quad R_j \end{array}$	{ LD $R_i, a(R_j)$ }
6)	$R_i \leftarrow \begin{array}{c} + \\ / \quad \backslash \\ R_i \quad \mathbf{ind} \\ \quad \quad \\ \quad \quad + \\ \quad \quad / \quad \backslash \\ \quad \quad C_a \quad R_j \end{array}$	{ ADD $R_i, R_i, a(R_j)$ }
7)	$R_i \leftarrow \begin{array}{c} + \\ / \quad \backslash \\ R_i \quad R_j \end{array}$	{ ADD R_i, R_i, R_j }
8)	$R_i \leftarrow \begin{array}{c} + \\ / \quad \backslash \\ R_i \quad C_1 \end{array}$	{ INC R_i }

Syntax-directed translation scheme

- | | | |
|-----|---|--------------------------------------|
| 1) | $R_i \rightarrow \mathbf{c}_a$ | $\{ \text{LD } R_i, \#a \}$ |
| 2) | $R_i \rightarrow M_x$ | $\{ \text{LD } R_i, x \}$ |
| 3) | $M \rightarrow = M_x R_i$ | $\{ \text{ST } x, R_i \}$ |
| 4) | $M \rightarrow = \mathbf{ind} R_i R_j$ | $\{ \text{ST } *R_i, R_j \}$ |
| 5) | $R_i \rightarrow \mathbf{ind} + \mathbf{c}_a R_j$ | $\{ \text{LD } R_i, a(R_j) \}$ |
| 6) | $R_i \rightarrow + R_i \mathbf{ind} + \mathbf{c}_a R_j$ | $\{ \text{ADD } R_i, R_i, a(R_j) \}$ |
| 7) | $R_i \rightarrow + R_i R_j$ | $\{ \text{ADD } R_i, R_i, R_j \}$ |
| 8) | $R_i \rightarrow + R_i \mathbf{c}_1$ | $\{ \text{INC } R_i \}$ |
| 9) | $R \rightarrow \mathbf{sp}$ | |
| 10) | $M \rightarrow \mathbf{m}$ | |

An instruction set for tree matching

$R_i \leftarrow C_a$



Ershov Numbers

- Label any leaf 1.
- The label of an interior node with one child is the label of its child.
- The label of an interior node with two children is
 - The larger of the labels of its children, if those labels are different.
 - One plus the label of its children if the labels are the same.

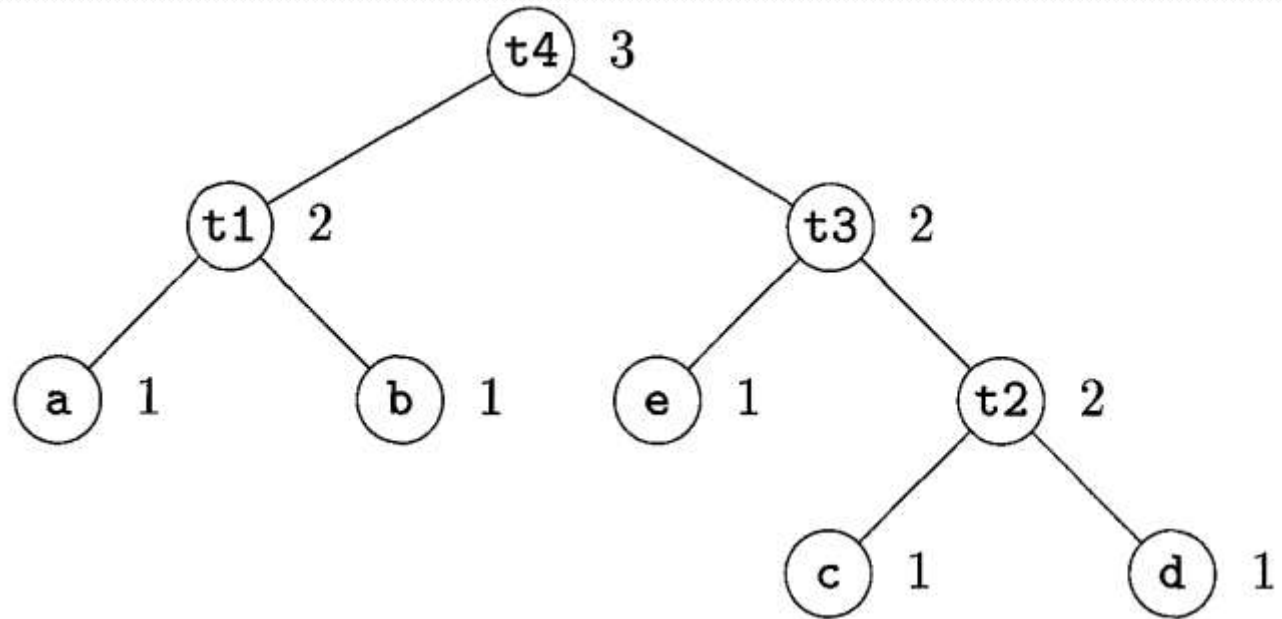
A tree labeled with Ershov numbers

$$t_1 = a - b$$

$$t_2 = c + d$$

$$t_3 = e * t_2$$

$$t_4 = t_1 + t_3$$



Generating code from a labeled expression tree

- To generate machine code for an interior node with label k and two children with equal labels (which must be $k - 1$) do the following:
 - Recursively generate code for the right child, using base $b+1$. The result of the right child appears in register R_{b+k} .
 - Recursively generate code for the left child, using base b ; the result appears in R_{b+k-1} .
 - Generate the instruction $OP R_{b+k}, R_{b+k-1}, R_{b+k}$, where OP is the appropriate operation for the interior node in question.
- Suppose we have an interior node with label k and children with unequal labels. Then one of the children, which we'll call the "big" child, has label k , and the other child, the "little" child, has some label $m < k$. Do the following to generate code for this interior node, using base b :
 - Recursively generate code for the big child, using base b ; the result appears in register R_{b+k-1} .
 - Recursively generate code for the small child, using base b ; the result appears in register R_{b+m-1} . Note that since $m < k$, neither R_{b+k-1} nor any higher-numbered register is used.
 - Generate the instruction $OP R_{b+k-1}, R_{b+m-1}, R_{b+k-1}$ or the instruction $OP R_{b+k-1}, R_{b+m-1}, R_{b+m-1}$, depending on whether the big child is the right or left child, respectively.
- For a leaf representing operand x , if the base is b generate the instruction $LD R_b, x$.

Optimal three-register code

```
LD R3, d
LD R2, c
ADD R3, R2, R3
LD R2, e
MUL R3, R2, R3
LD R2, b
LD R1, a
SUB R2, R1, R2
ADD R3, R2, R3
```


Evaluating Expressions with an Insufficient Supply of Registers

- Node N has at least one child with label r or greater. Pick the larger child (or either if their labels are the same) to be the "big" child and let the other child be the "little" child.
- Recursively generate code for the big child, using base $b = 1$. The result of this evaluation will appear in register R_r .
- Generate the machine instruction $ST\ t_k, R_r$, where t_k is a temporary variable used for temporary results used to help evaluate nodes with label k .
- Generate code for the little child as follows. If the little child has label r or greater, pick base $b=1$. If the label of the little child is $j < r$, then pick $b=r-j$. Then recursively apply this algorithm to the little child; the result appears in R_r .
- Generate the instruction $LD\ R_{r-1}, t_k$.
- If the big child is the right child of N , then generate the instruction $OP\ R_r, R_r, R_{r-1}$. If the big child is the left child, generate $OP\ R_r, R_{r-1}, R_r$.

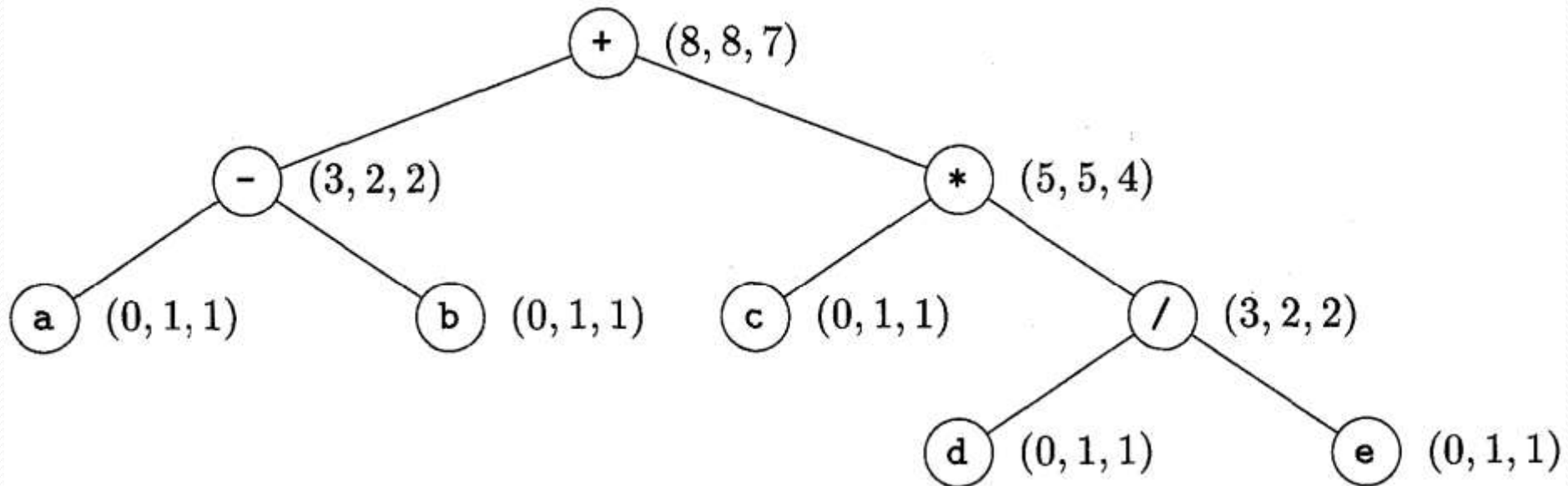
Optimal three-register code using only two registers

```
LD  R2, d
LD  R1, c
ADD R2, R1, R2
LD  R1, e
MUL R2, R1, R2
ST  t3, R2
LD  R2, b
LD  R1, a
SUB R2, R1, R2
LD  R1, t3
ADD R2, R2, R1
```

Dynamic Programming Algorithm

- Compute bottom-up for each node n of the expression tree T an array C of costs, in which the i th component $C[i]$ is the optimal cost of computing the subtree S rooted at n into a register, assuming i registers are available for the computation, for
- Traverse T , using the cost vectors to determine which subtrees of T must be computed into memory.
- Traverse each tree using the cost vectors and associated instructions to generate the final target code. The code for the subtrees computed into memory locations is generated first.

Syntax tree for $(a-b)+c*(d/e)$ with cost vector at each node



minimum cost of evaluating the root with two registers available

- Compute the left subtree with two registers available into register R_0 , compute the right subtree with one register available into register R_1 , and use the instruction $\text{ADD } R_0, R_0, R_1$ to compute the root. This sequence has cost $2+5+1=8$.
- Compute the right subtree with two registers available into R_1 , compute the left subtree with one register available into R_0 , and use the instruction $\text{ADD } R_0, R_0, R_1$. This sequence has cost $4+2+1=7$.
- Compute the right subtree into memory location M , compute the left subtree with two registers available into register R_0 , and use the instruction $\text{ADD } R_0, R_0, M$. This sequence has cost $5+2+1=8$.

Compiler course

Chapter 5

Syntax Directed Translation

Outline

- Syntax Directed Definitions
- Evaluation Orders of SDD's
- Applications of Syntax Directed Translation
- Syntax Directed Translation Schemes

Introduction

- We can associate information with a language construct by attaching attributes to the grammar symbols.
- A syntax directed definition specifies the values of attributes by associating semantic rules with the grammar productions.

Production	Semantic Rule
$E \rightarrow E1 + T$	$E.code = E1.code T.code '+'$

- We may alternatively insert the semantic actions inside the grammar
$$E \rightarrow E1 + T \{ \text{print '+'} \}$$

Syntax Directed Definitions

- A SDD is a context free grammar with attributes and rules
- Attributes are associated with grammar symbols and rules with productions
- Attributes may be of many kinds: numbers, types, table references, strings, etc.
- Synthesized attributes
 - A synthesized attribute at node N is defined only in terms of attribute values of children of N and at N it
- Inherited attributes
 - An inherited attribute at node N is defined only in terms of attribute values at N 's parent, N itself and N 's siblings

Example of S-attributed SDD

Production

- 1) $L \rightarrow E n$
- 2) $E \rightarrow E1 + T$
- 3) $E \rightarrow T$
- 4) $T \rightarrow T1 * F$
- 5) $T \rightarrow F$
- 6) $F \rightarrow (E)$
- 7) $F \rightarrow \text{digit}$

Semantic Rules

- $L.val = E.val$
- $E.val = E1.val + T.val$
- $E.val = T.val$
- $T.val = T1.val * F.val$
- $T.val = F.val$
- $F.val = E.val$
- $F.val = \text{digit.lexval}$

Example of mixed attributes

Production

1) $T \rightarrow FT'$

2) $T' \rightarrow *FT'_1$

3) $T' \rightarrow \epsilon$

1) $F \rightarrow \text{digit}$

Semantic Rules

$T'.inh = F.val$

$T.val = T'.syn$

$T'_1.inh = T'.inh * F.val$

$T'.syn = T'_1.syn$

$T'.syn = T'.inh$

$F.val = F.val = \text{digit.lexval}$

Evaluation orders for SDD's

- A dependency graph is used to determine the order of computation of attributes
- Dependency graph
 - For each parse tree node, the parse tree has a node for each attribute associated with that node
 - If a semantic rule defines the value of synthesized attribute A.b in terms of the value of X.c then the dependency graph has an edge from X.c to A.b
 - If a semantic rule defines the value of inherited attribute B.c in terms of the value of X.a then the dependency graph has an edge from X.c to B.c
- Example!

Ordering the evaluation of attributes

- If dependency graph has an edge from M to N then M must be evaluated before the attribute of N
- Thus the only allowable orders of evaluation are those sequence of nodes N_1, N_2, \dots, N_k such that if there is an edge from N_i to N_j then $i < j$
- Such an ordering is called a topological sort of a graph
- Example!

S-Attributed definitions

- An SDD is S-attributed if every attribute is synthesized
- We can have a post-order traversal of parse-tree to evaluate attributes in S-attributed definitions

```
postorder(N) {  
    for (each child C of N, from the left) postorder(C);  
    evaluate the attributes associated with node N;  
}
```

- S-Attributed definitions can be implemented during bottom-up parsing without the need to explicitly create parse trees

L-Attributed definitions

- A SDD is L-Attributed if the edges in dependency graph goes from Left to Right but not from Right to Left.
- More precisely, each attribute must be either
 - Synthesized
 - Inherited, but if there us a production $A \rightarrow X_1 X_2 \dots X_n$ and there is an inherited attribute X_i computed by a rule associated with this production, then the rule may only use:
 - Inherited attributes associated with the head A
 - Either inherited or synthesized attributes associated with the occurrences of symbols X_1, X_2, \dots, X_{i-1} located to the left of X_i
 - Inherited or synthesized attributes associated with this occurrence of X_i itself, but in such a way that there is no cycle in the graph

Application of Syntax Directed Translation

- Type checking and intermediate code generation (chapter 6)
- Construction of syntax trees
 - Leaf nodes: $\text{Leaf}(\text{op}, \text{val})$
 - Interior node: $\text{Node}(\text{op}, c_1, c_2, \dots, c_k)$
- Example:

Production

- 1) $E \rightarrow E_1 + T$
- 2) $E \rightarrow E_1 - T$
- 3) $E \rightarrow T$
- 4) $T \rightarrow (E)$
- 5) $T \rightarrow \text{id}$
- 6) $T \rightarrow \text{num}$

Semantic Rules

- $E.\text{node} = \text{new node}(\text{'+'}, E_1.\text{node}, T.\text{node})$
- $E.\text{node} = \text{new node}(\text{'-'}, E_1.\text{node}, T.\text{node})$
- $E.\text{node} = T.\text{node}$
- $T.\text{node} = E.\text{node}$
- $T.\text{node} = \text{new Leaf}(\text{id}, \text{id.entry})$
- $T.\text{node} = \text{new Leaf}(\text{num}, \text{num.val})$

Syntax tree for L-attributed definition

Production

1) $E \rightarrow TE'$

2) $E' \rightarrow + TE1'$

3) $E' \rightarrow -TE1'$

4) $E' \rightarrow \epsilon$

5) $T \rightarrow (E)$

6) $T \rightarrow id$

7) $T \rightarrow num$

Semantic Rules

$E.node = E'.syn$ +

$E'.inh = T.node$

$E1'.inh = \text{new node}('+', E'.inh, T.node)$

$E'.syn = E1'.syn$

$E1'.inh = \text{new node}('-', E'.inh, T.node)$

$E'.syn = E1'.syn$

$E'.syn = E'.inh$

$T.node = E.node$

$T.node = \text{new Leaf}(id, id.entry)$

$T.node = \text{new Leaf}(num, num.val)$

Syntax directed translation schemes

- An SDT is a Context Free grammar with program fragments embedded within production bodies
- Those program fragments are called semantic actions
- They can appear at any position within production body
- Any SDT can be implemented by first building a parse tree and then performing the actions in a left-to-right depth first order
- Typically SDT's are implemented during parsing without building a parse tree

Postfix translation schemes

- Simplest SDDs are those that we can parse the grammar bottom-up and the SDD is s-attributed
- For such cases we can construct SDT where each action is placed at the end of the production and is executed along with the reduction of the body to the head of that production
- SDT's with all actions at the right ends of the production bodies are called postfix SDT's

Example of postfix SDT

- 1) $L \rightarrow E n$ $\{\text{print}(E.\text{val});\}$
- 2) $E \rightarrow E1 + T$ $\{E.\text{val}=E1.\text{val}+T.\text{val};\}$
- 3) $E \rightarrow T$ $\{E.\text{val} = T.\text{val};\}$
- 4) $T \rightarrow T1 * F$ $\{T.\text{val}=T1.\text{val}*F.\text{val};\}$
- 5) $T \rightarrow F$ $\{T.\text{val}=F.\text{val};\}$
- 6) $F \rightarrow (E)$ $\{F.\text{val}=E.\text{val};\}$
- 7) $F \rightarrow \text{digit}$ $\{F.\text{val}=\text{digit}.\text{lexval};\}$

Parse-Stack implementation of postfix SDT's

- In a shift-reduce parser we can easily implement semantic action using the parser stack
- For each nonterminal (or state) on the stack we can associate a record holding its attributes
- Then in a reduction step we can execute the semantic action at the end of a production to evaluate the attribute(s) of the non-terminal at the leftside of the production
- And put the value on the stack in replace of the rightside of production

Example

L \rightarrow E n { print(stack[top-1].val);
 top=top-1;}

E \rightarrow E1 + T { stack[top-2].val=stack[top-2].val+stack.val;
 top=top-2;}

E \rightarrow T

T \rightarrow T1 * F { stack[top-2].val=stack[top-2].val+stack.val;
 top=top-2;}

T \rightarrow F

F \rightarrow (E) { stack[top-2].val=stack[top-1].val
 top=top-2;}

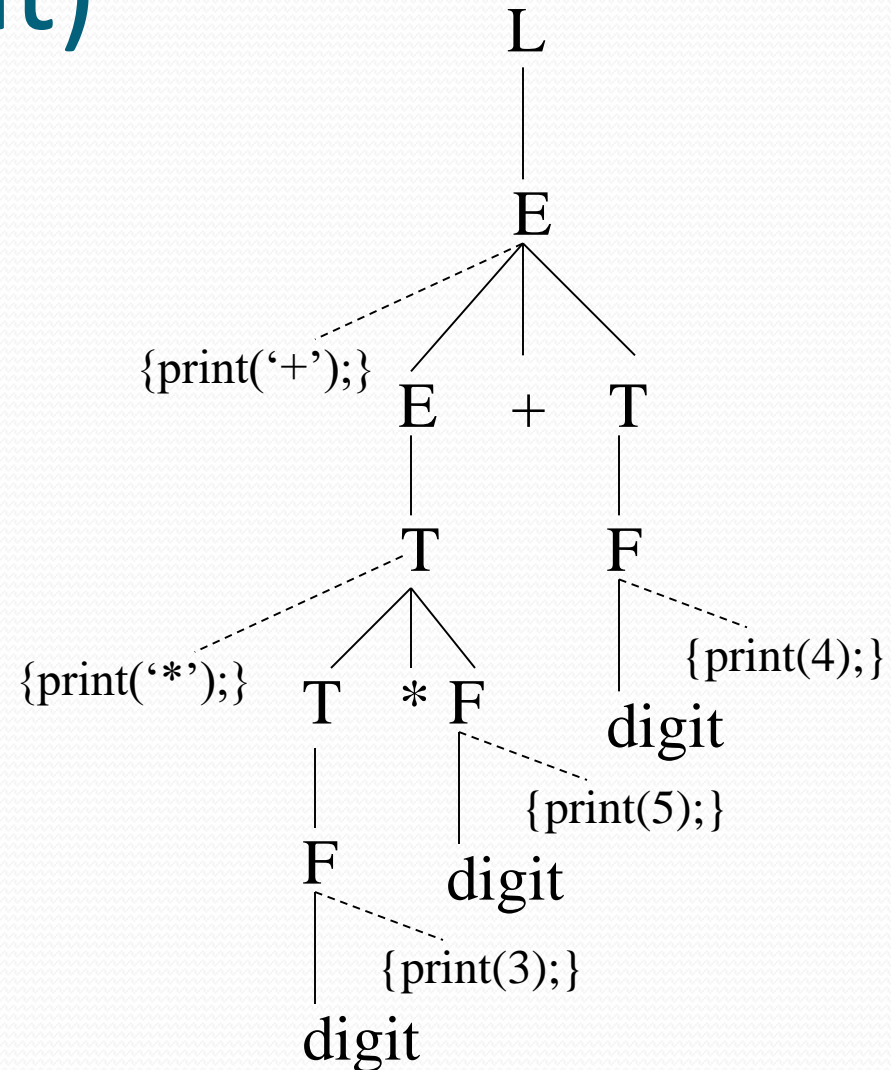
F \rightarrow digit

SDT's with actions inside productions

- For a production $B \rightarrow X \{a\} Y$
 - If the parse is bottom-up then we perform action “a” as soon as this occurrence of X appears on the top of the parser stack
 - If the parser is top down we perform “a” just before we expand Y
 - Sometimes we cant do things as easily as explained above
 - One example is when we are parsing this SDT with a bottom-
- 1) $L \rightarrow E n$
 - 2) $E \rightarrow \{\text{print('+');}\} E1 + T$
 - 3) $E \rightarrow T$
 - 4) $T \rightarrow \{\text{print('*') ;}\} T1 * F$
 - 5) $T \rightarrow F$
 - 6) $F \rightarrow (E)$
 - 7) $F \rightarrow \text{digit } \{\text{print(digit.lexval);}\}$

SDT's with actions inside productions (cont)

- Any SDT can be implemented as follows
 1. Ignore the actions and produce a parse tree
 2. Examine each interior node N and add actions as new children at the correct position
 3. Perform a postorder traversal and execute actions when their nodes are visited



SDT's for L-Attributed definitions

- We can convert an L-attributed SDD into an SDT using following two rules:
 - Embed the action that computes the inherited attributes for a nonterminal A immediately before that occurrence of A . if several inherited attributes of A are dependent on one another in an acyclic fashion, order them so that those needed first are computed first
 - Place the action of a synthesized attribute for the head of a production at the end of the body of the production

Example

```
S -> while (C) S1      L1=new();
                       L2=new();
                       S1.next=L1;
                       C.false=S.next;
                       C.true=L2;
                       S.code=label||L1||C.code||label||L2||S1.code
```

```
S -> while ( {L1=new();L2=new();C.false=S.next;C.true=L2;}
C) {S1.next=L1;}
S1 {S.code=label||L1||C.code||label||L2||S1.code;}
```

Readings

- Chapter 5 of the book

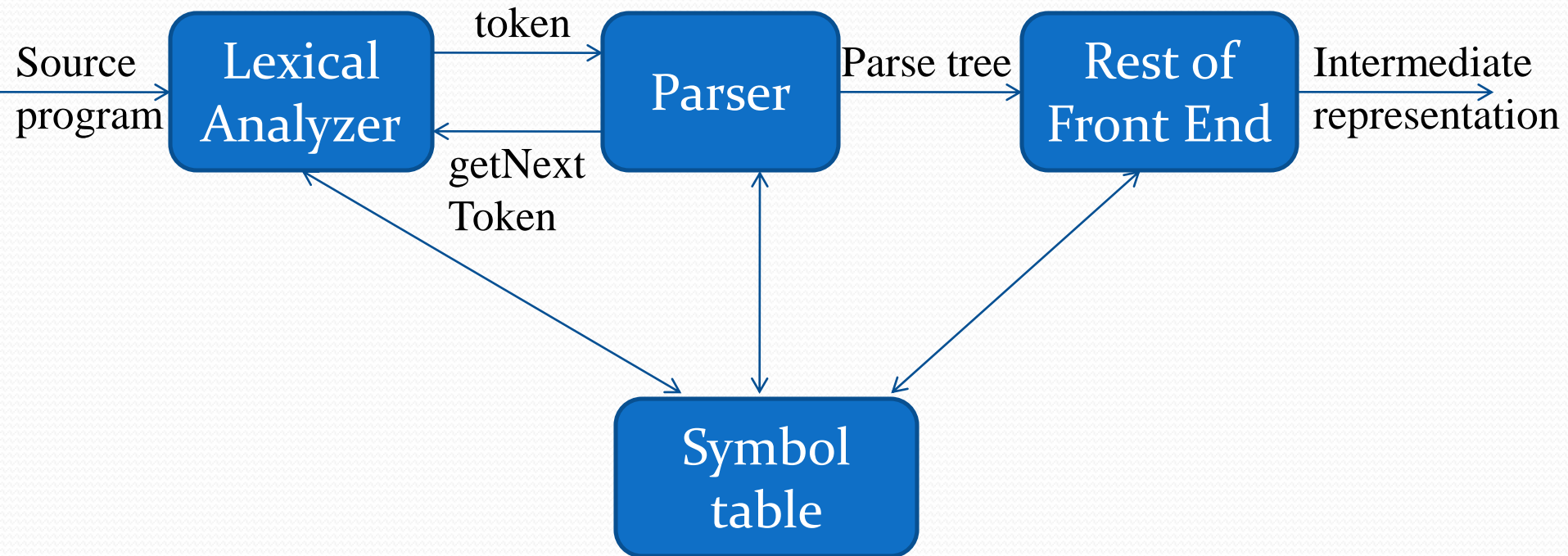
Compiler course

Chapter 4
Syntax Analysis

Outline

- Role of parser
- Context free grammars
- Top down parsing
- Bottom up parsing
- Parser generators

The role of parser



Uses of grammars

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

Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
 - Lexical errors
- Error handler goals
 - Report the presence of errors clearly and accurately
 - Recover from each error quickly enough to detect subsequent errors
 - Add minimal overhead to the processing of correct programs

Error-recover strategies

- Panic mode recovery
 - Discard input symbol one at a time until one of designated set of synchronization tokens is found
- Phrase level recovery
 - Replacing a prefix of remaining input by some string that allows the parser to continue
- Error productions
 - Augment the grammar with productions that generate the erroneous constructs
- Global correction
 - Choosing minimal sequence of changes to obtain a globally least-cost correction

Context free grammars

- Terminals
- Nonterminals
- Start symbol
- productions

expression \rightarrow expression + term

expression \rightarrow expression – term

expression \rightarrow term

term \rightarrow term * factor

term \rightarrow term / factor

term \rightarrow factor

factor \rightarrow (expression)

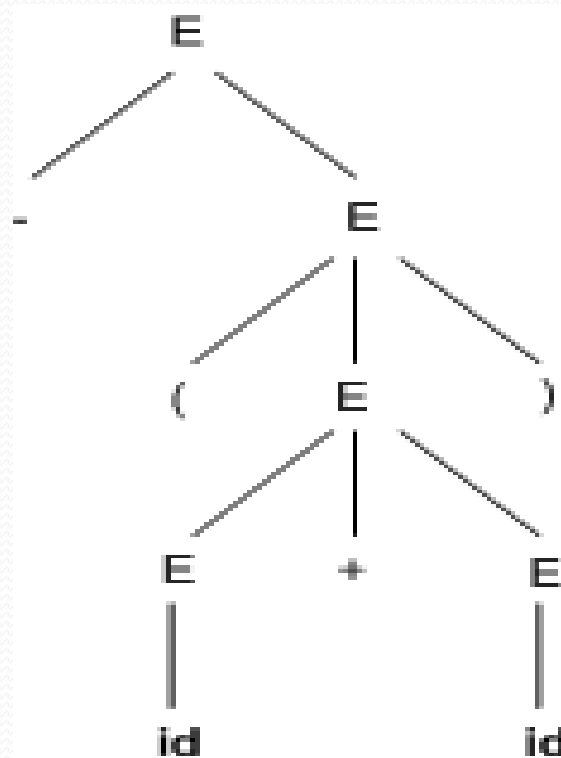
factor \rightarrow **id**

Derivations

- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
 - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \mathbf{id}$
 - Derivations for $\mathbf{-(id+id)}$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow \mathbf{-(id+E)} \Rightarrow \mathbf{-(id+id)}$

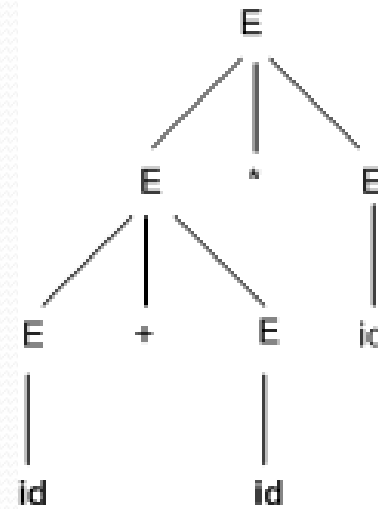
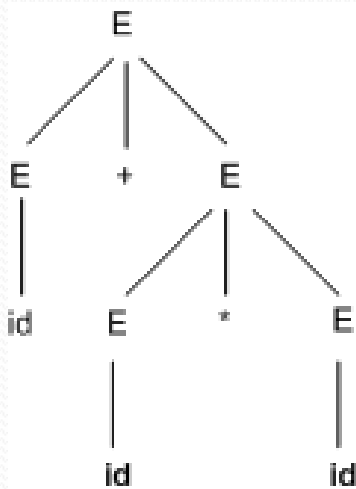
Parse trees

- $-(\text{id}+\text{id})$
- $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\text{id}+E) \Rightarrow -(\text{id}+\text{id})$



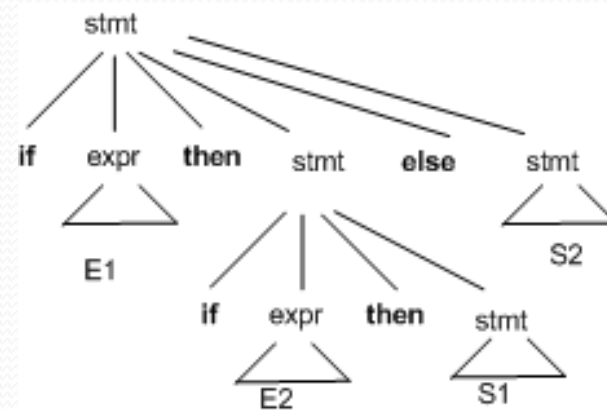
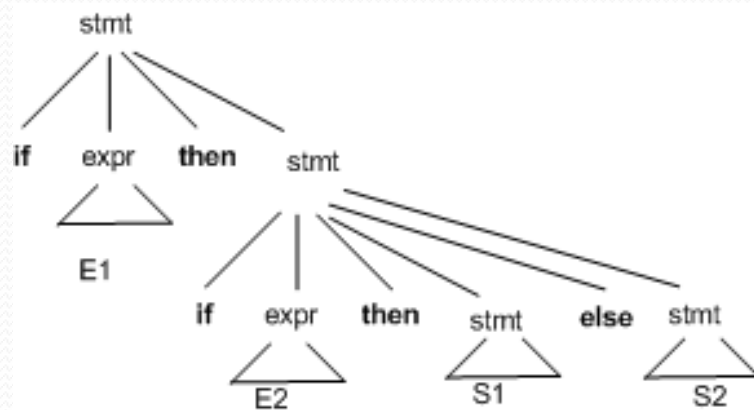
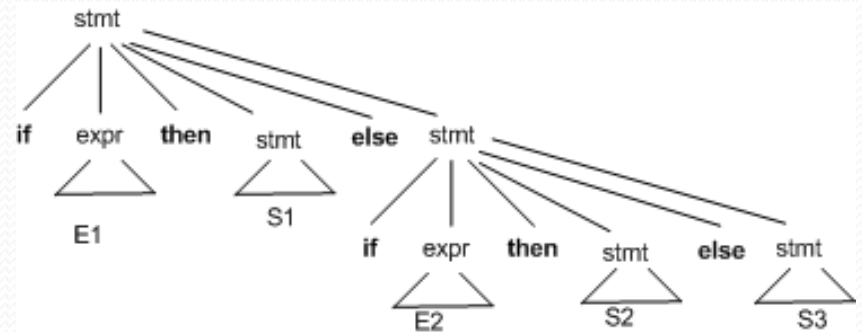
Ambiguity

- For some strings there exist more than one parse tree
- Or more than one leftmost derivation
- Or more than one rightmost derivation
- Example: $\text{id}+\text{id}*\text{id}$



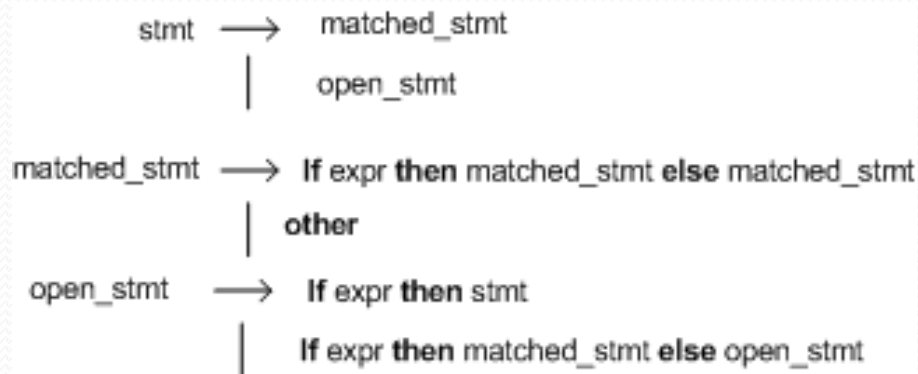
Elimination of ambiguity

stmt \rightarrow **If** expr **then** stmt
| **If** expr **then** stmt **else** stmt
| **other**



Elimination of ambiguity (cont.)

- Idea:
 - A statement appearing between a **then** and an **else** must be matched



Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \xRightarrow{+} A \alpha$
- Top down parsing methods cant handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:
 - $A \rightarrow A \alpha \mid \beta$
 - We may replace it with
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' \mid \epsilon$

Left recursion elimination (cont.)

- There are cases like following
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \varepsilon$
- Left recursion elimination algorithm:
 - Arrange the nonterminals in some order A_1, A_2, \dots, A_n .
 - For (each i from 1 to n) {
 - For (each j from 1 to $i-1$) {
 - Replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j productions
 - }
 - Eliminate left recursion among the A_i -productions
 - }

Left factoring

- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.
- Consider following grammar:
 - Stmt \rightarrow **if** expr **then** stmt **else** stmt
 - | **if** expr **then** stmt
- On seeing input **if** it is not clear for the parser which production to use
- We can easily perform left factoring:
 - If we have $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$ then we replace it with
 - $A \rightarrow \alpha A'$
 - $A' \rightarrow \beta_1 \mid \beta_2$

Left factoring (cont.)

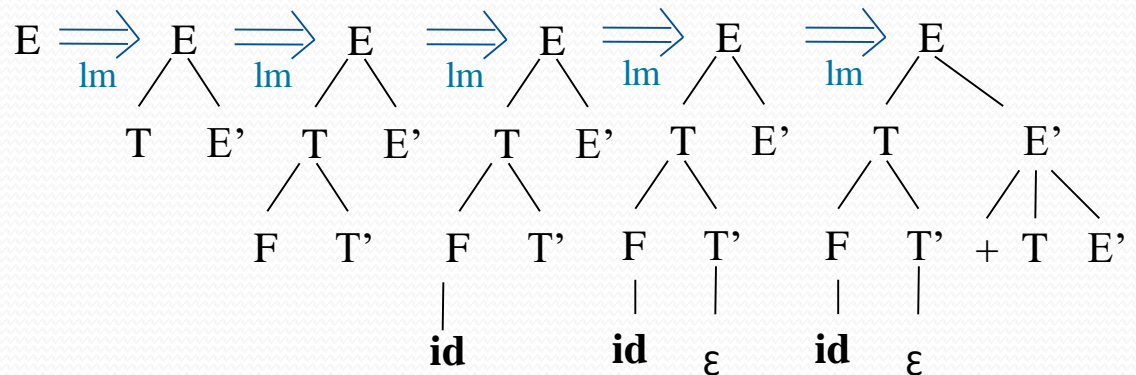
- Algorithm
 - For each non-terminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$, then replace all of A -productions $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$ by
 - $A \rightarrow \alpha A' \mid \gamma$
 - $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
- Example:
 - $S \rightarrow I E t S \mid i E t S e S \mid a$
 - $E \rightarrow b$

Top Down Parsing

Introduction

- A Top-down parser tries to create a parse tree from the root towards the leafs scanning input from left to right
- It can be also viewed as finding a leftmost derivation for an input string
- Example: $id+id*id$

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



Recursive descent parsing

- Consists of a set of procedures, one for each nonterminal
- Execution begins with the procedure for start symbol
- A typical procedure for a non-terminal

```
void A() {  
    choose an A-production, A->X1X2..Xk  
    for (i=1 to k) {  
        if (Xi is a nonterminal  
            call procedure Xi();  
        else if (Xi equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```

Recursive descent parsing (cont)

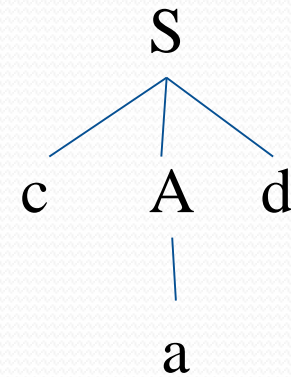
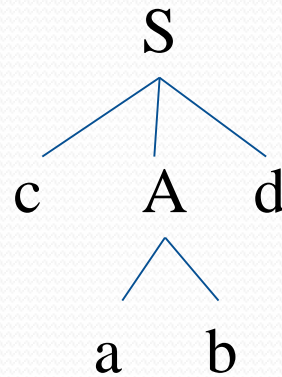
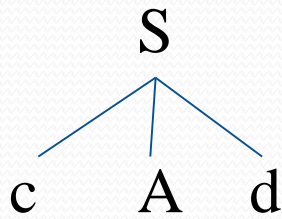
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it cant choose an A -production easily.
- So we need to try all alternatives
- If one failed the input pointer needs to be reset and another alternative should be tried
- Recursive descent parsers cant be used for left-recursive grammars

Example

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

Input: cad



First and Follow

- $\text{First}()$ is set of terminals that begins strings derived from
- If $\alpha \xRightarrow{*} \epsilon$ then ϵ is also in $\text{First}(\epsilon)$
- In predictive parsing when we have $A \rightarrow \alpha \mid \beta$, if $\text{First}(\alpha)$ and $\text{First}(\beta)$ are disjoint sets then we can select appropriate A-production by looking at the next input
- $\text{Follow}(A)$, for any nonterminal A, is set of terminals a that can appear immediately after A in some sentential form
 - If we have $S \xRightarrow{*} \alpha A a \beta$ for some α and β then a is in $\text{Follow}(A)$
- If A can be the rightmost symbol in some sentential form, then $\$$ is in $\text{Follow}(A)$

Computing First

- To compute $\text{First}(X)$ for all grammar symbols X , apply following rules until no more terminals or ϵ can be added to any First set:
 1. If X is a terminal then $\text{First}(X) = \{X\}$.
 2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$, then place a in $\text{First}(X)$ if for some i a is in $\text{First}(Y_i)$ and ϵ is in all of $\text{First}(Y_1), \dots, \text{First}(Y_{i-1})$ that is $Y_1 \dots Y_{i-1} \xRightarrow{*} \epsilon$. if ϵ is in $\text{First}(Y_j)$ for $j=1, \dots, k$ then add ϵ to $\text{First}(X)$.
 3. If $X \rightarrow \epsilon$ is a production then add ϵ to $\text{First}(X)$
- Example!

Computing follow

- To compute $\text{First}(A)$ for all nonterminals A , apply following rules until nothing can be added to any follow set:
 1. Place $\$$ in $\text{Follow}(S)$ where S is the start symbol
 2. If there is a production $A \rightarrow \alpha B \beta$ then everything in $\text{First}(\beta)$ except ϵ is in $\text{Follow}(B)$.
 3. If there is a production $A \rightarrow B$ or a production $A \rightarrow \alpha B \beta$ where $\text{First}(\beta)$ contains ϵ , then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$
- Example!

LL(1) Grammars

- Predictive parsers are those recursive descent parsers needing no backtracking
- Grammars for which we can create predictive parsers are called LL(1)
 - The first L means scanning input from left to right
 - The second L means leftmost derivation
 - And 1 stands for using one input symbol for lookahead
- A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha \mid \beta$ are two distinct productions of G , the following conditions hold:
 - For no terminal a do α and β both derive strings beginning with a
 - At most one of α or β can derive empty string
 - If $\alpha \Rightarrow \varepsilon$ then β does not derive any string beginning with a terminal in $\text{Follow}(A)$.

Construction of predictive parsing table

- For each production $A \rightarrow \alpha$ in grammar do the following:
 1. For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow$ in $M[A,a]$
 2. If ϵ is in $\text{First}(\alpha)$, then for each terminal b in $\text{Follow}(A)$ add $A \rightarrow \epsilon$ to $M[A,b]$. If ϵ is in $\text{First}(\alpha)$ and $\$$ is in $\text{Follow}(A)$, add $A \rightarrow \epsilon$ to $M[A,\$]$ as well
- If after performing the above, there is no production in $M[A,a]$ then set $M[A,a]$ to error

Example

$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow (E) \mid \mathbf{id}$

	First	Follow
F	{(,id}	{+, *,), \$}
T	{(,id}	{+,), \$}
E	{(,id}	{), \$}
E'	{+, ϵ }	{), \$}
T'	{*, ϵ }	{+,), \$}

Non - terminal	Input Symbol					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \mathbf{id}$			$F \rightarrow (E)$		

Another example

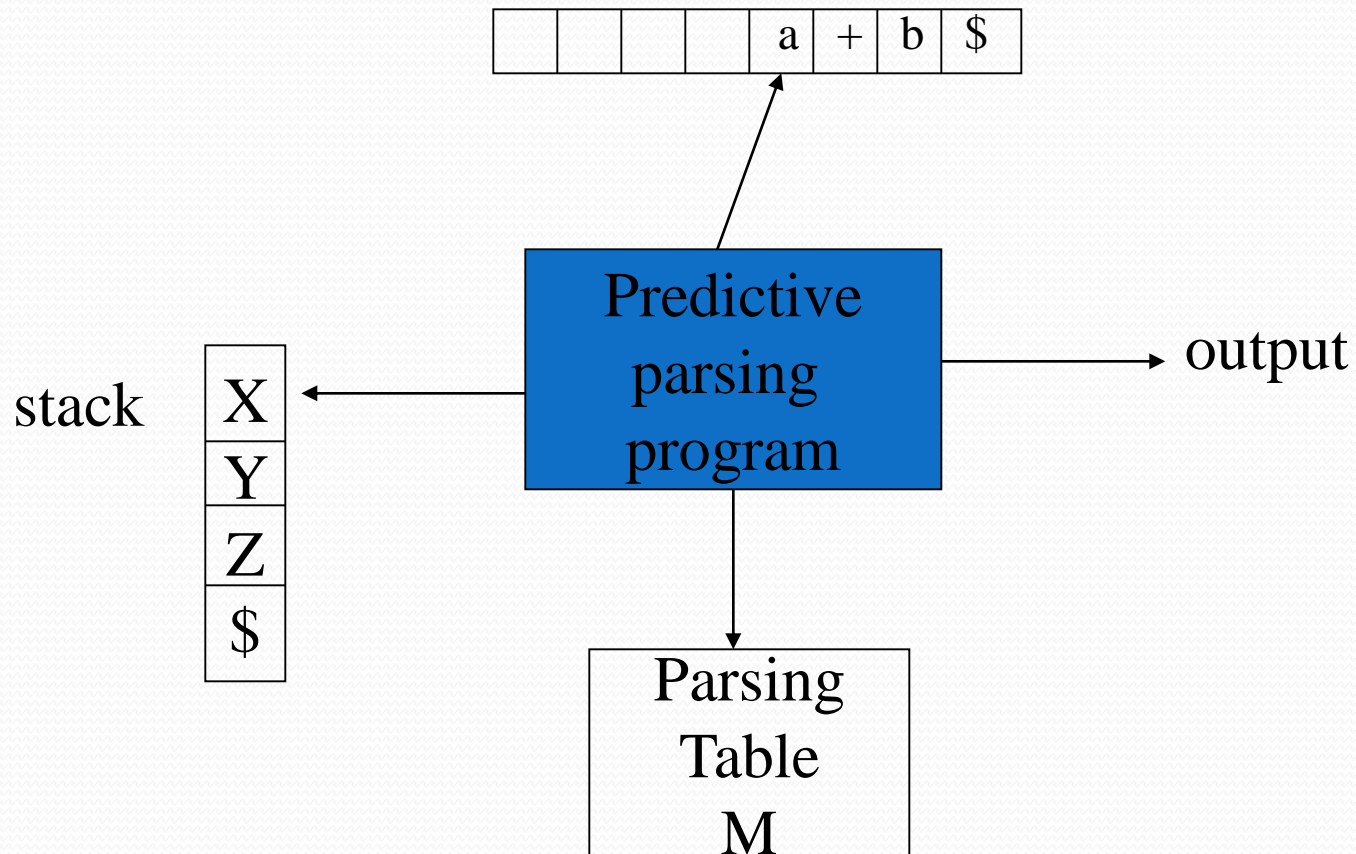
$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \epsilon$

$E \rightarrow b$

Non - terminal	Input Symbol					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Non-recursive predicting parsing



Predictive parsing algorithm

Set ip point to the first symbol of w ;

Set X to the top stack symbol;

While $(X \neq \$)$ { /* stack is not empty */

 if (X is a) pop the stack and advance ip;

 else if (X is a terminal) error();

 else if ($M[X,a]$ is an error entry) error();

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

 output the production $X \rightarrow Y_1 Y_2 \dots Y_k$;

 pop the stack;

 push Y_k, \dots, Y_2, Y_1 on to the stack with Y_1 on top;

 }

 set X to the top stack symbol;

}

Example

- $\text{id+id*id\$}$

Matched	Stack	Input	Action
	E\$	$\text{id+id*id\$}$	

Error recovery in predictive parsing

- Panic mode
 - Place all symbols in $\text{Follow}(A)$ into synchronization set for nonterminal A : skip tokens until an element of $\text{Follow}(A)$ is seen and pop A from stack.
 - Add to the synchronization set of lower level construct the symbols that begin higher level constructs
 - Add symbols in $\text{First}(A)$ to the synchronization set of nonterminal A
 - If a nonterminal can generate the empty string then the production deriving can be used as a default
 - If a terminal on top of the stack cannot be matched, pop the terminal, issue a message saying that the terminal was inserted

Example

Non-terminal	Input Symbol					
	id	+	*	()	\$
E	E → TE'			E → TE'	synch	synch
E'		E' → +TE'			E' → ε	E' → ε
T	T → FT'	synch		T → FT'	synch	synch
T'		T' → ε	T' → *FT'		T' → ε	T' → ε
F	F → id	synch	synch	F → (E)	synch	synch

Stack	Input	Action
E\$)id*+id\$	Error, Skip)
E\$	id*+id\$	id is in First(E)
TE'\$	id*+id\$	
FT'E'\$	id*+id\$	
idT'E'\$	id*+id\$	
T'E'\$	*+id\$	
*FT'E'\$	*+id\$	
FT'E'\$	+id\$	Error, M[F,+]=synch
T'E'\$	+id\$	F has been popped



Bottom-up Parsing

Introduction

- Constructs parse tree for an input string beginning at the leaves (the bottom) and working towards the root (the top)
- Example: $\text{id}*\text{id}$

$E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \mathbf{id}$

$\text{id}*\text{id}$

$F * \text{id}$

$T * \text{id}$

$T * F$

F

E

id

F
 id

F id
 id

$T * F$
 F id
 id

F
 $T * F$
 F id
 id

Shift-reduce parser

- The general idea is to shift some symbols of input to the stack until a reduction can be applied
- At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of the production
- The key decisions during bottom-up parsing are about when to reduce and about what production to apply
- A reduction is a reverse of a step in a derivation
- The goal of a bottom-up parser is to construct a derivation in reverse:
 - $E \Rightarrow T \Rightarrow T * F \Rightarrow T * id \Rightarrow F * id \Rightarrow id * id$

Handle pruning

- A Handle is a substring that matches the body of a production and whose reduction represents one step along the reverse of a rightmost derivation

Right sentential form	Handle	Reducing production
id*id	id	F->id
F*id	F	T->F
T*id	id	F->id
T*F	T*F	E->T*F

Shift reduce parsing

- A stack is used to hold grammar symbols
- Handle always appear on top of the stack
- Initial configuration:

Stack	Input
\$	w\$

- Acceptance configuration

Stack	Input
\$S	\$

Shift reduce parsing (cont.)

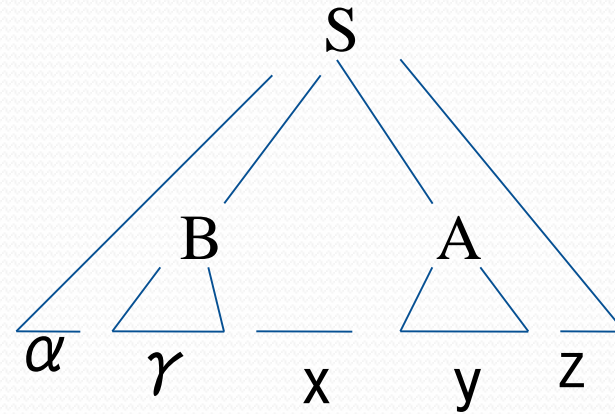
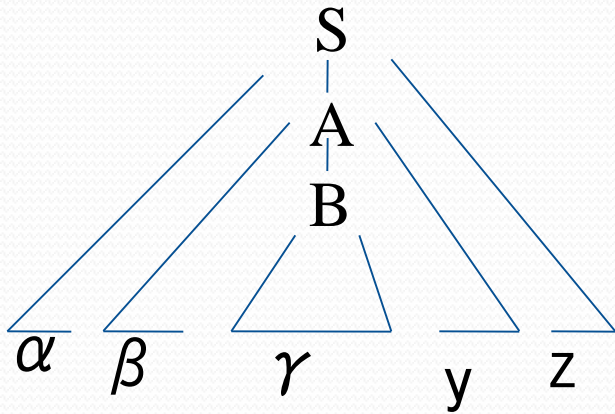
- Basic operations:

- Shift
- Reduce
- Accept
- Error

- Example: $id*id$

Stack	Input	Action
\$	$id*id\$$	shift
$\$id$	$*id\$$	reduce by $F \rightarrow id$
$\$F$	$*id\$$	reduce by $T \rightarrow F$
$\$T$	$*id\$$	shift
$\$T*$	$id\$$	shift
$\$T*id$	$\$$	reduce by $F \rightarrow id$
$\$T*F$	$\$$	reduce by $T \rightarrow T*F$
$\$T$	$\$$	reduce by $E \rightarrow T$
$\$E$	$\$$	accept

Handle will appear on top of the stack



Stack	Input
$\$ \alpha \beta \gamma$	$yz\$$
$\$ \alpha \beta B$	$yz\$$
$\$ \alpha \beta By$	$z\$$

Stack	Input
$\$ \alpha \gamma$	$xyz\$$
$\$ \alpha Bxy$	$z\$$

Conflicts during shift reduce parsing

- Two kind of conflicts
 - Shift/reduce conflict
 - Reduce/reduce conflict
- Example:

```
stmt → If expr then stmt
      | If expr then stmt else stmt
      | other
```

Stack

... if expr then stmt

Input

else ...\$

Reduce/reduce conflict

stmt -> id(parameter_list)

stmt -> expr:=expr

parameter_list->parameter_list, parameter

parameter_list->parameter

parameter->id

expr->id(expr_list)

expr->id

expr_list->expr_list, expr

expr_list->expr

Stack

... id(id

Input

,id) ...\$

LR Parsing

- The most prevalent type of bottom-up parsers
- LR(k), mostly interested on parsers with $k \leq 1$
- Why LR parsers?
 - Table driven
 - Can be constructed to recognize all programming language constructs
 - Most general non-backtracking shift-reduce parsing method
 - Can detect a syntactic error as soon as it is possible to do so
 - Class of grammars for which we can construct LR parsers are superset of those which we can construct LL parsers

States of an LR parser

- States represent set of items
- An LR(o) item of G is a production of G with the dot at some position of the body:
 - For $A \rightarrow XYZ$ we have following items
 - $A \rightarrow .XYZ$
 - $A \rightarrow X.YZ$
 - $A \rightarrow XY.Z$
 - $A \rightarrow XYZ.$
 - In a state having $A \rightarrow .XYZ$ we hope to see a string derivable from XYZ next on the input.
 - What about $A \rightarrow X.YZ$?

Constructing canonical LR(0) item sets

- Augmented grammar:
 - G with addition of a production: $S' \rightarrow S$
- Closure of item sets:
 - If I is a set of items, $\text{closure}(I)$ is a set of items constructed from I by the following rules:
 - Add every item in I to $\text{closure}(I)$
 - If $A \rightarrow \alpha.B\beta$ is in $\text{closure}(I)$ and $B \rightarrow \gamma$ is a production then add the item $B \rightarrow \gamma$ to $\text{closure}(I)$.

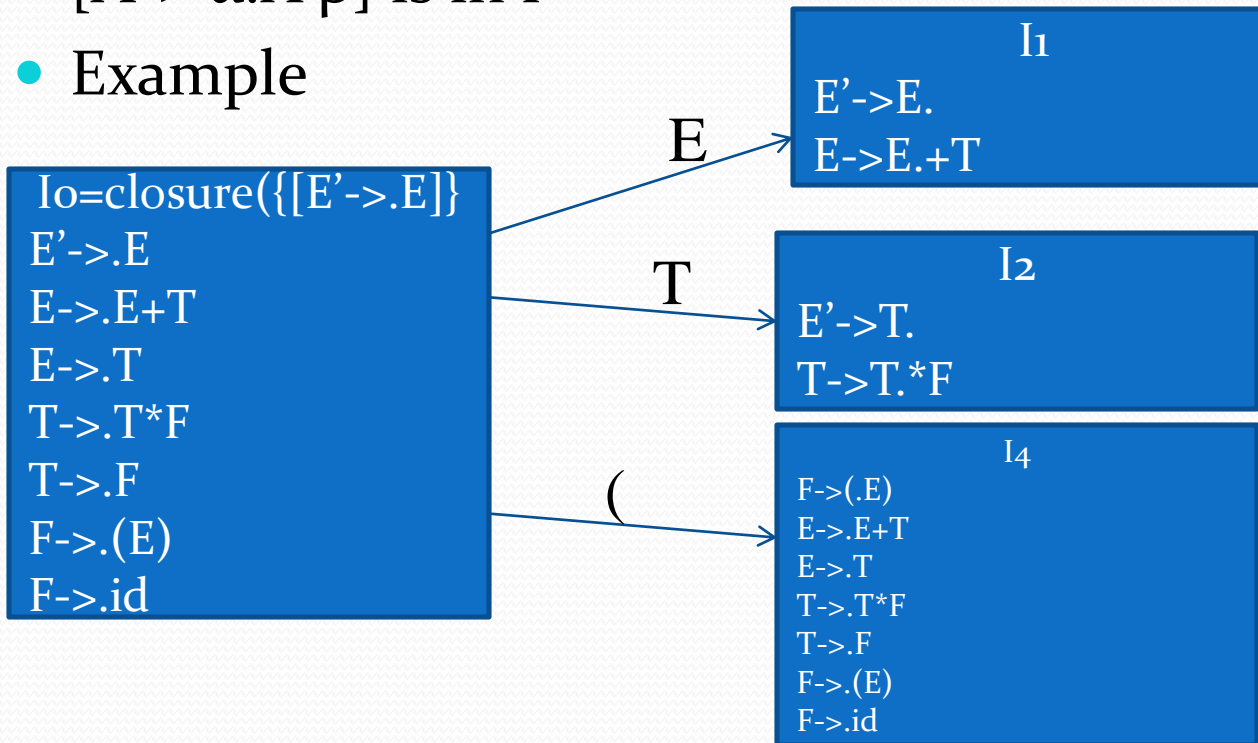
- Example:
 $E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \mathbf{id}$

$I_0 = \text{closure}(\{[E' \rightarrow \cdot E]\})$

$E' \rightarrow \cdot E$
 $E \rightarrow \cdot E + T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot \mathbf{id}$

Constructing canonical LR(0) item sets (cont.)

- Goto (I,X) where I is an item set and X is a grammar symbol is closure of set of all items $[A \rightarrow \alpha X \beta]$ where $[A \rightarrow \alpha.X \beta]$ is in I
- Example



Closure algorithm

SetOfItems CLOSURE(I) {

 J=I;

 repeat

 for (each item $A \rightarrow \alpha.B\beta$ in J)

 for (each production $B \rightarrow \gamma$ of G)

 if ($B \rightarrow \cdot\gamma$ is not in J)

 add $B \rightarrow \cdot\gamma$ to J;

 until no more items are added to J on one round;

 return J;

GOTO algorithm

SetOfItems GOTO(I,X) {

 J=empty;

 if (A- \rightarrow α .X β is in I)

 add CLOSURE(A- \rightarrow α X. β) to J;

 return J;

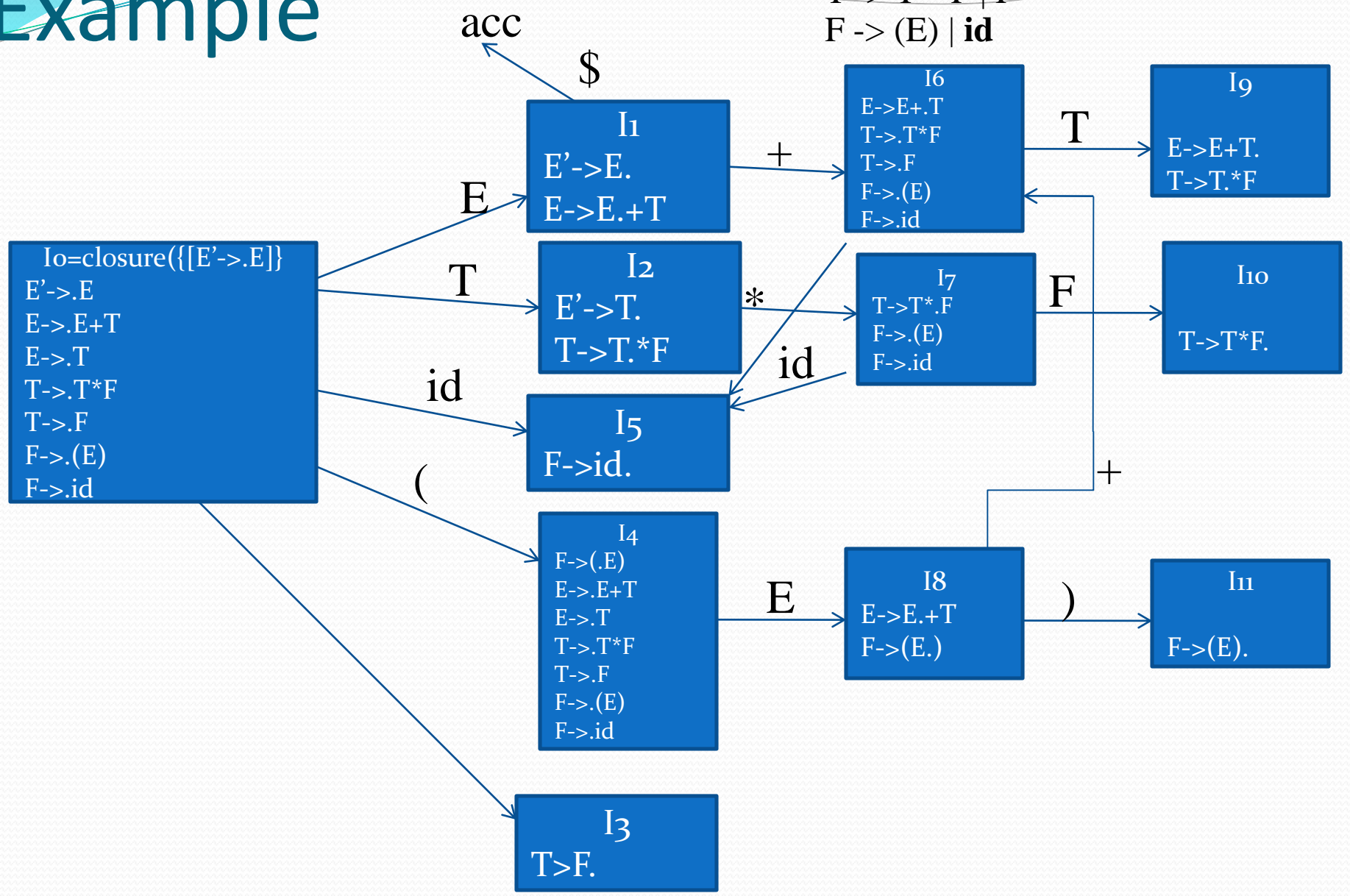
}

Canonical LR(0) items

```
Void items(G') {  
    C = CLOSURE({[S' -> .S]});  
    repeat  
        for (each set of items I in C)  
            for (each grammar symbol X)  
                if (GOTO(I,X) is not empty and not in C)  
                    add GOTO(I,X) to C;  
    until no new set of items are added to C on a round;  
}
```

Example

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

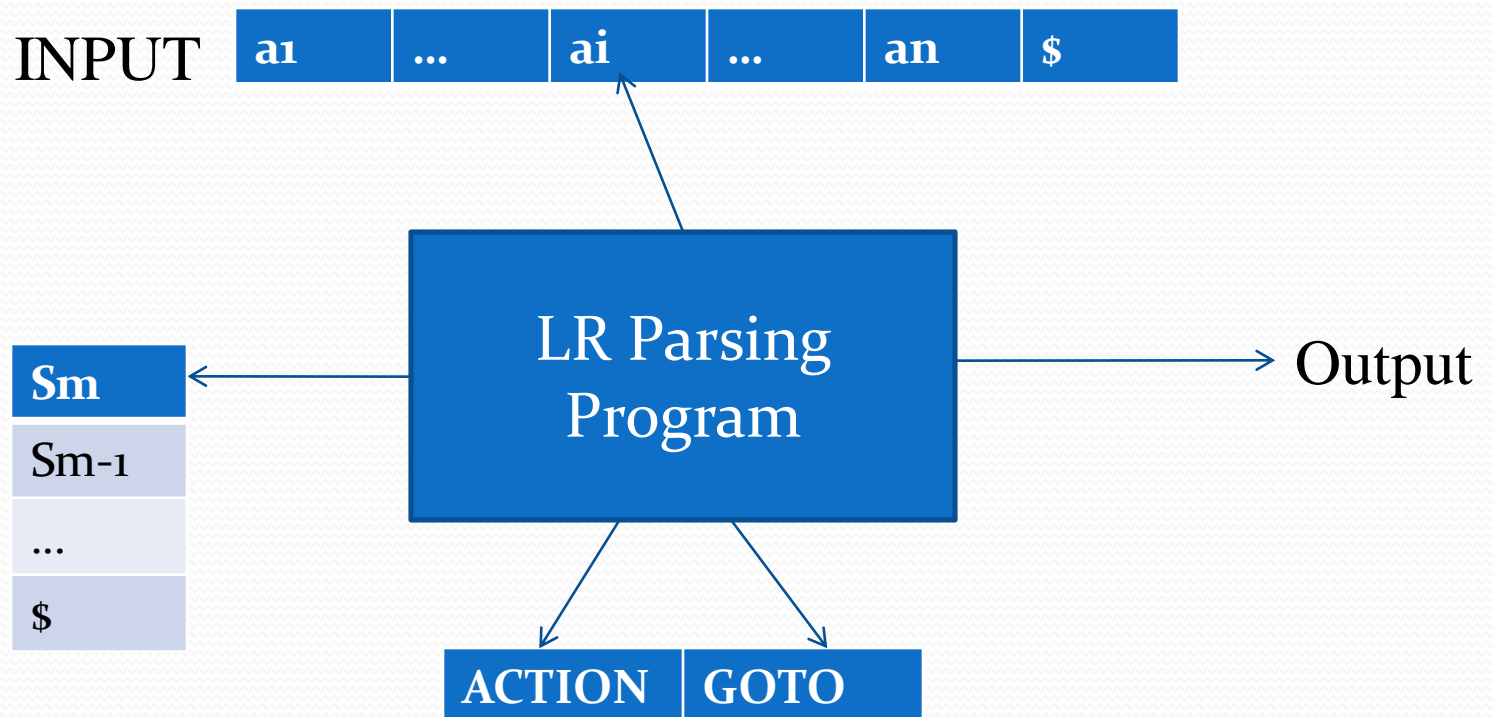


Use of LR(0) automaton

- Example: $id*id$

Line	Stack	Symbols	Input	Action
(1)	0	\$	id*id\$	Shift to 5
(2)	05	\$id	*id\$	Reduce by $F \rightarrow id$
(3)	03	\$F	*id\$	Reduce by $T \rightarrow F$
(4)	02	\$T	*id\$	Shift to 7
(5)	027	\$T*	id\$	Shift to 5
(6)	0275	\$T*id	\$	Reduce by $F \rightarrow id$
(7)	02710	\$T*F	\$	Reduce by $T \rightarrow T*F$
(8)	02	\$T	\$	Reduce by $E \rightarrow T$
(9)	01	\$E	\$	accept

LR-Parsing model



LR parsing algorithm

```
let a be the first symbol of w$;
while(1) { /*repeat forever */
    let s be the state on top of the stack;
    if (ACTION[s,a] = shift t) {
        push t onto the stack;
        let a be the next input symbol;
    } else if (ACTION[s,a] = reduce A-> $\beta$ ) {
        pop  $|\beta|$  symbols of the stack;
        let state t now be on top of the stack;
        push GOTO[t,A] onto the stack;
        output the production A-> $\beta$ ;
    } else if (ACTION[s,a]=accept) break; /* parsing is done */
    else call error-recovery routine;
}
```


Example

- (0) $E' \rightarrow E$
- (1) $E \rightarrow E + T$
- (2) $E \rightarrow T$
- (3) $T \rightarrow T * F$
- (4) $T \rightarrow F$
- (5) $F \rightarrow (E)$
- (6) $F \rightarrow id$

id*id+id?

STATE	ACTON						GOTO		
	id	+	*	()	\$	E	T	F
0	S5			S4			1	2	3
1		S6				Acc			
2		R2	S7		R2	R2			
3		R4	R7		R4	R4			
4	S5			S4			8	2	3
5		R6	R6		R6	R6			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		R1	S7		R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

Line	Stack	Symbols	Input	Action
(1)	o		id*id+id\$	Shift to 5
(2)	o5	id	*id+id\$	Reduce by $F \rightarrow id$
(3)	o3	F	*id+id\$	Reduce by $T \rightarrow F$
(4)	o2	T	*id+id\$	Shift to 7
(5)	o27	T*	id+id\$	Shift to 5
(6)	o275	T*id	+id\$	Reduce by $F \rightarrow id$
(7)	o2710	T*F	+id\$	Reduce by $T \rightarrow T*F$
(8)	o2	T	+id\$	Reduce by $E \rightarrow T$
(9)	o1	E	+id\$	Shift
(10)	o16	E+	id\$	Shift
(11)	o165	E+id	\$	Reduce by $F \rightarrow id$
(12)	o163	E+F	\$	Reduce by $T \rightarrow F$
(13)	o169	E+T'	\$	Reduce by $E \rightarrow E+T$
(14)	o1	E	\$	accept

Constructing SLR parsing table

- Method

- Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of LR(0) items for G'
- State i is constructed from state I_i :
 - If $[A \rightarrow \alpha.a\beta]$ is in I_i and $\text{Goto}(I_i, a) = I_j$, then set $\text{ACTION}[i, a]$ to “shift j ”
 - If $[A \rightarrow \alpha.]$ is in I_i , then set $\text{ACTION}[i, a]$ to “reduce $A \rightarrow \alpha$ ” for all a in $\text{follow}(A)$
 - If $\{S' \rightarrow .S\}$ is in I_i , then set $\text{ACTION}[i, \$]$ to “Accept”
- If any conflicts appears then we say that the grammar is not SLR(1).
- If $\text{GOTO}(I_i, A) = I_j$ then $\text{GOTO}[i, A] = j$
- All entries not defined by above rules are made “error”
- The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S]$

Example grammar which is not SLR(1)

$S \rightarrow L=R \mid R$

$L \rightarrow *R \mid id$

$R \rightarrow L$

I0
 $S' \rightarrow .S$
 $S \rightarrow .L=R$
 $S \rightarrow .R$
 $L \rightarrow .*R \mid$
 $L \rightarrow .id$
 $R \rightarrow .L$

I1
 $S' \rightarrow S.$

I2
 $S \rightarrow L.=R$
 $R \rightarrow L.$

I3
 $S \rightarrow R.$

I4
 $L \rightarrow *.R$
 $R \rightarrow .L$
 $L \rightarrow .*R$
 $L \rightarrow .id$

I5
 $L \rightarrow id.$

I6
 $S \rightarrow L=.R$
 $R \rightarrow .L$
 $L \rightarrow .*R$
 $L \rightarrow .id$

I7
 $L \rightarrow *R.$

I8
 $R \rightarrow L.$

I9
 $S \rightarrow L=R.$

Action

=

Shift 6
 Reduce $R \rightarrow L$

More powerful LR parsers

- Canonical-LR or just LR method
 - Use lookahead symbols for items: LR(1) items
 - Results in a large collection of items
- LALR: lookaheads are introduced in LR(0) items

Canonical LR(1) items

- In LR(1) items each item is in the form: $[A \rightarrow \alpha.\beta, a]$
 - An LR(1) item $[A \rightarrow \alpha.\beta, a]$ is valid for a viable prefix γ if there is a derivation $S \xRightarrow{*} \delta A w \xRightarrow{rm} \delta \alpha \beta w$, where
 - $\Gamma = \delta \alpha$
 - Either a is the first symbol of w , or w is ϵ and a is $\$$
 - Example:
 - $S \rightarrow BB$
 - $B \rightarrow aB \mid b$
- $S \xRightarrow{*} aaBab \xRightarrow{rm} aaaBab$
- Item $[B \rightarrow a.B, a]$ is valid for $\gamma = aaa$ and $w = ab$

Constructing LR(1) sets of items

```
SetOfItems Closure(I) {  
  repeat  
    for (each item  $[A \rightarrow \alpha.B\beta, a]$  in I)  
      for (each production  $B \rightarrow \gamma$  in  $G'$ )  
        for (each terminal  $b$  in  $\text{First}(\beta a)$ )  
          add  $[B \rightarrow \cdot\gamma, b]$  to set I;  
  
  until no more items are added to I;  
  return I;  
}
```

```
SetOfItems Goto(I, X) {  
  initialize J to be the empty set;  
  for (each item  $[A \rightarrow \alpha.X\beta, a]$  in I)  
    add item  $[A \rightarrow \alpha X.\beta, a]$  to set J;  
  return closure(J);  
}
```

```
void items( $G'$ ) {  
  initialize C to  $\text{Closure}(\{[S' \rightarrow \cdot S, \$]\})$ ;  
  repeat  
    for (each set of items I in C)  
      for (each grammar symbol X)  
        if ( $\text{Goto}(I, X)$  is not empty and not in C)  
          add  $\text{Goto}(I, X)$  to C;  
  
  until no new sets of items are added to C;  
}
```

Example

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

Canonical LR(1) parsing table

- Method
 - Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of LR(1) items for G'
 - State i is constructed from state I_i :
 - If $[A \rightarrow \alpha.a\beta, b]$ is in I_i and $\text{Goto}(I_i, a) = I_j$, then set $\text{ACTION}[i, a]$ to “shift j ”
 - If $[A \rightarrow \alpha., a]$ is in I_i , then set $\text{ACTION}[i, a]$ to “reduce $A \rightarrow \alpha$ ”
 - If $\{S' \rightarrow .S, \$\}$ is in I_i , then set $\text{ACTION}[i, \$]$ to “Accept”
 - If any conflicts appears then we say that the grammar is not LR(1).
 - If $\text{GOTO}(I_i, A) = I_j$ then $\text{GOTO}[i, A] = j$
 - All entries not defined by above rules are made “error”
 - The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S, \$]$

Example

$S' \rightarrow S$

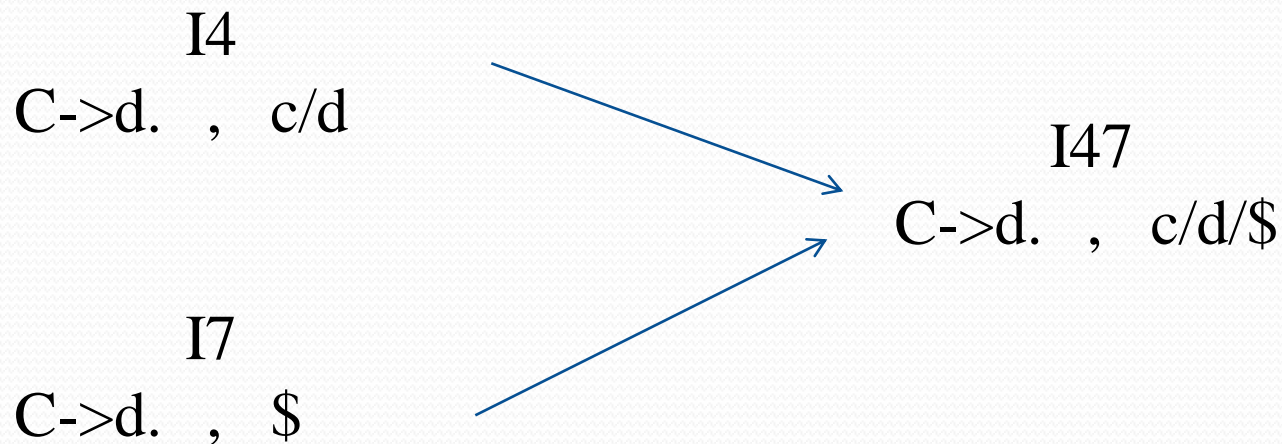
$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

LALR Parsing Table

- For the previous example we had:



- State merges can't produce Shift-Reduce conflicts. Why?
- But it may produce reduce-reduce conflict

Example of RR conflict in state merging

$S' \rightarrow S$

$S \rightarrow aAd \mid bBd \mid aBe \mid bAe$

$A \rightarrow c$

$B \rightarrow c$

An easy but space-consuming LALR table construction

- Method:
 1. Construct $C = \{I_0, I_1, \dots, I_n\}$ the collection of LR(1) items.
 2. For each core among the set of LR(1) items, find all sets having that core, and replace these sets by their union.
 3. Let $C' = \{J_0, J_1, \dots, J_m\}$ be the resulting sets. The parsing actions for state i , is constructed from J_i as before. If there is a conflict grammar is not LALR(1).
 4. If J is the union of one or more sets of LR(1) items, that is $J = I_1 \cup I_2 \dots \cup I_k$ then the cores of $\text{Goto}(I_1, X)$, ..., $\text{Goto}(I_k, X)$ are the same and is a state like K , then we set $\text{Goto}(J, X) = k$.
- This method is not efficient, a more efficient one is discussed in the book

Compaction of LR parsing table

- Many rows of action tables are identical
 - Store those rows separately and have pointers to them from different states
 - Make lists of (terminal-symbol, action) for each state
 - Implement Goto table by having a link list for each nonterminal in the form (current state, next state)

Using ambiguous grammars

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

I0: $E' \rightarrow .E$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

I1: $E' \rightarrow E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I4: $E \rightarrow E + .E$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

$E \rightarrow E. (E)$

$E \rightarrow E. id$

I2: $E \rightarrow (.E)$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

I5: $E \rightarrow E * .E$

$E \rightarrow (.E)$

$E \rightarrow .E + E$

$E \rightarrow .E * E$

$E \rightarrow .(E)$

$E \rightarrow .id$

STATE	ACTON						GO TO
	id	+	*	()	\$	E
0	S ₃			S ₂			1
1		S ₄	S ₅			Acc	
2	S ₃		S ₂				6
3		R ₄	R ₄		R ₄	R ₄	
4	S ₃			S ₂			7
5	S ₃			S ₂			8
6		S ₄	S ₅				
7		R ₁	S ₅		R ₁	R ₁	
8		R ₂	R ₂		R ₂	R ₂	
9		R ₃	R ₃		R ₃	R ₃	

I3: $E \rightarrow .id$

I6: $E \rightarrow (E.)$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I8: $E \rightarrow E * E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I7: $E \rightarrow E + E.$

$E \rightarrow E. + E$

$E \rightarrow E. * E$

I9: $E \rightarrow (E).$

Readings

- Chapter 4 of the book

Compiler course

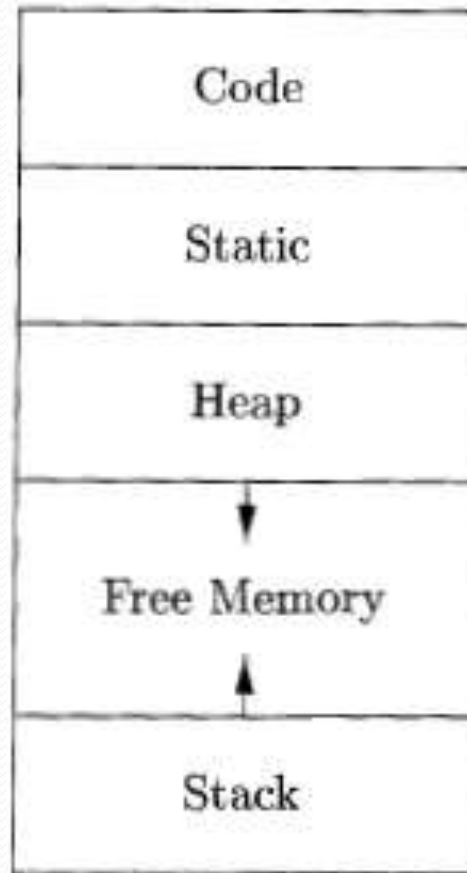
Chapter 7

Run-Time Environments

Outline

- Compiler must do the storage allocation and provide access to variables and data
- Memory management
 - Stack allocation
 - Heap management
 - Garbage collection

Storage Organization



Static vs. Dynamic Allocation

- Static: Compile time, Dynamic: Runtime allocation
- Many compilers use some combination of following
 - Stack storage: for local variables, parameters and so on
 - Heap storage: Data that may outlive the call to the procedure that created it
- Stack allocation is a valid allocation for procedures since procedure calls are nested

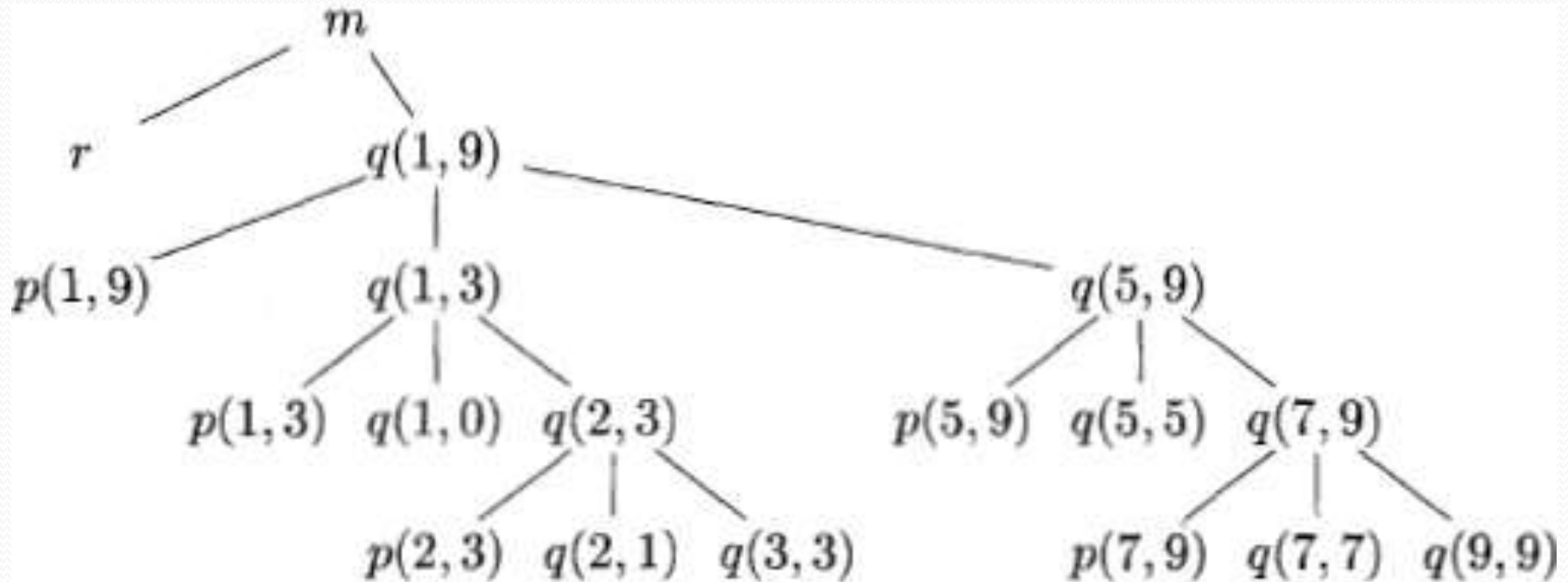
Sketch of a quicksort program

```
int a[11];
void readArray() { /* Reads 9 integers into a[1], ..., a[9]. */
    int i;
    ...
}
int partition(int m, int n) {
    /* Picks a separator value  $v$ , and partitions  $a[m..n]$  so that
        $a[m..p-1]$  are less than  $v$ ,  $a[p] = v$ , and  $a[p+1..n]$  are
       equal to or greater than  $v$ . Returns  $p$ . */
    ...
}
void quicksort(int m, int n) {
    int i;
    if (n > m) {
        i = partition(m, n);
        quicksort(m, i-1);
        quicksort(i+1, n);
    }
}
main() {
    readArray();
    a[0] = -9999;
    a[10] = 9999;
    quicksort(1,9);
}
```

Activation for Quicksort

```
enter main()
  enter readArray()
  leave readArray()
  enter quicksort(1,9)
    enter partition(1,9)
    leave partition(1,9)
    enter quicksort(1,3)
      ...
    leave quicksort(1,3)
    enter quicksort(5,9)
      ...
    leave quicksort(5,9)
  leave quicksort(1,9)
leave main()
```

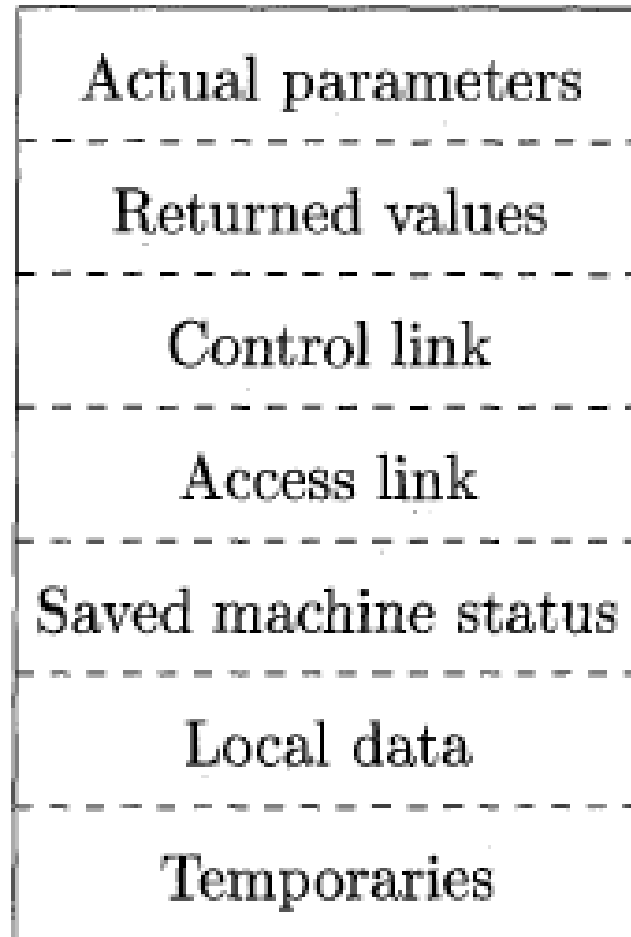
Activation tree representing calls during an execution of quicksort



Activation records

- Procedure calls and returns are usually managed by a run-time stack called the control stack.
- Each live activation has an activation record (sometimes called a frame)
- The root of activation tree is at the bottom of the stack
- The current execution path specifies the content of the stack with the last activation has record in the top of the stack.

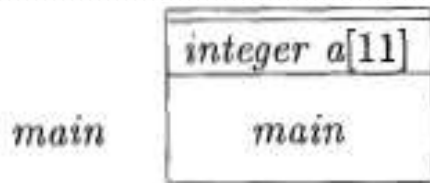
A General Activation Record



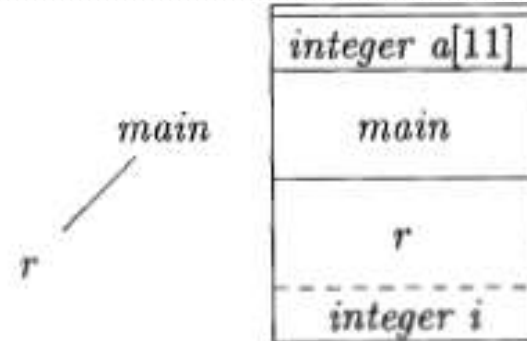
Activation Record

- Temporary values
- Local data
- A saved machine status
- An “access link”
- A control link
- Space for the return value of the called function
- The actual parameters used by the calling procedure

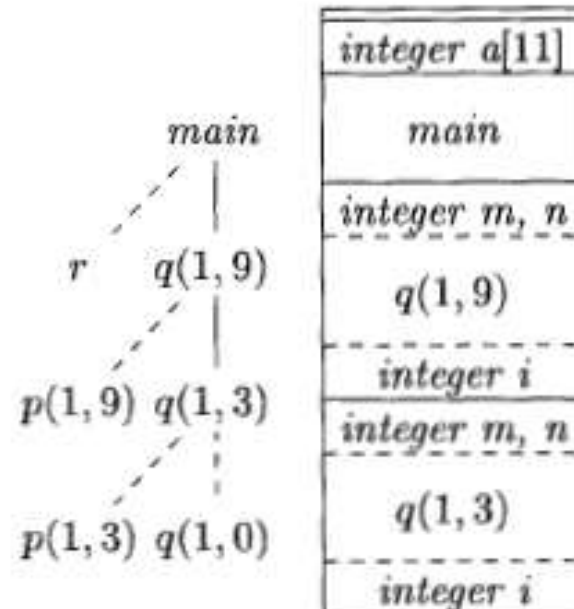
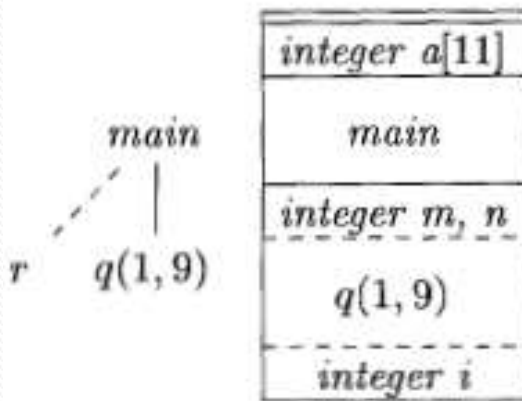
Downward-growing stack of activation records



(a) Frame for *main*



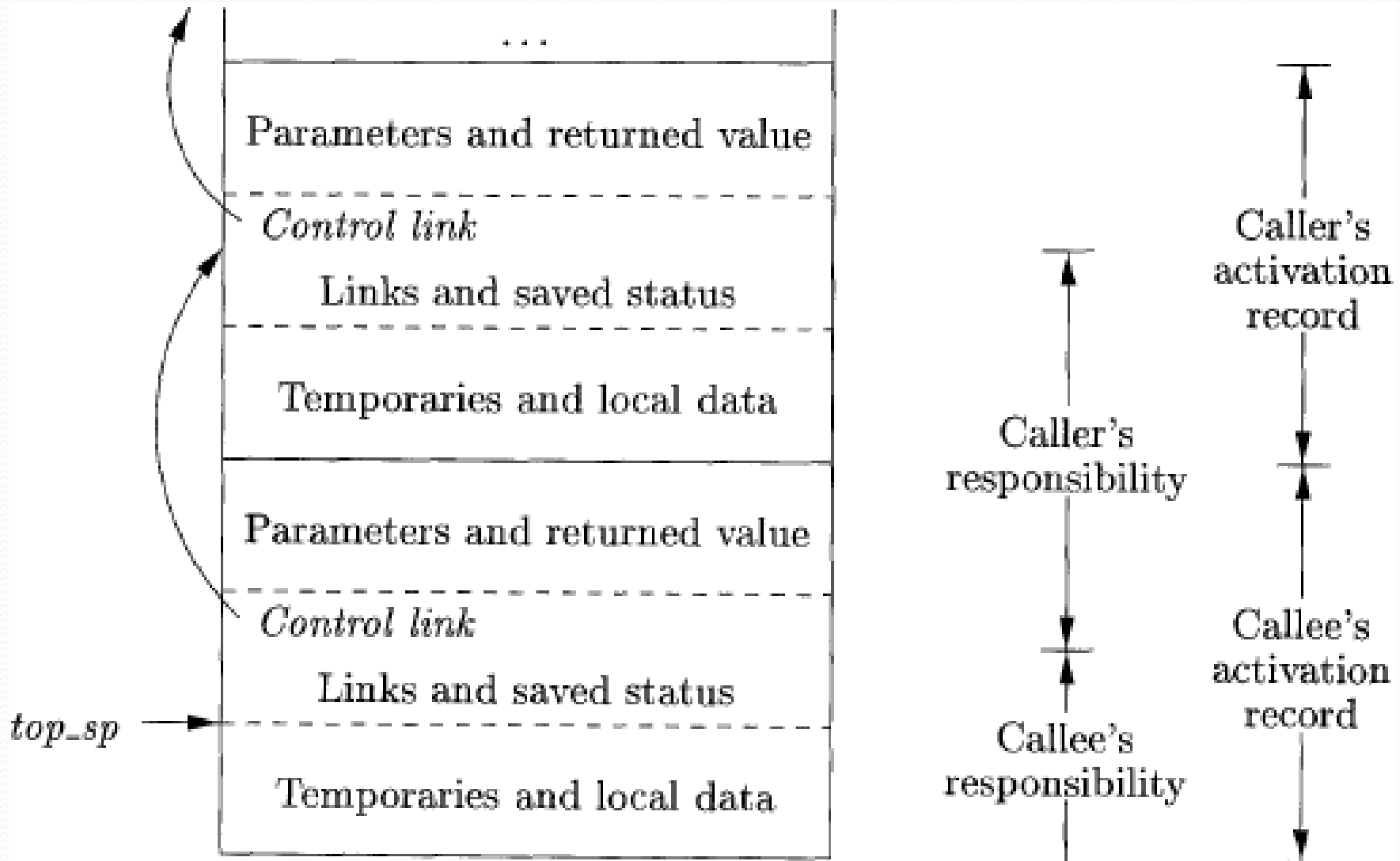
(b) *r* is activated



Designing Calling Sequences

- Values communicated between caller and callee are generally placed at the beginning of callee's activation record
- Fixed-length items: are generally placed at the middle
- Items whose size may not be known early enough: are placed at the end of activation record
- We must locate the top-of-stack pointer judiciously: a common approach is to have it point to the end of fixed length fields.

Division of tasks between caller and callee



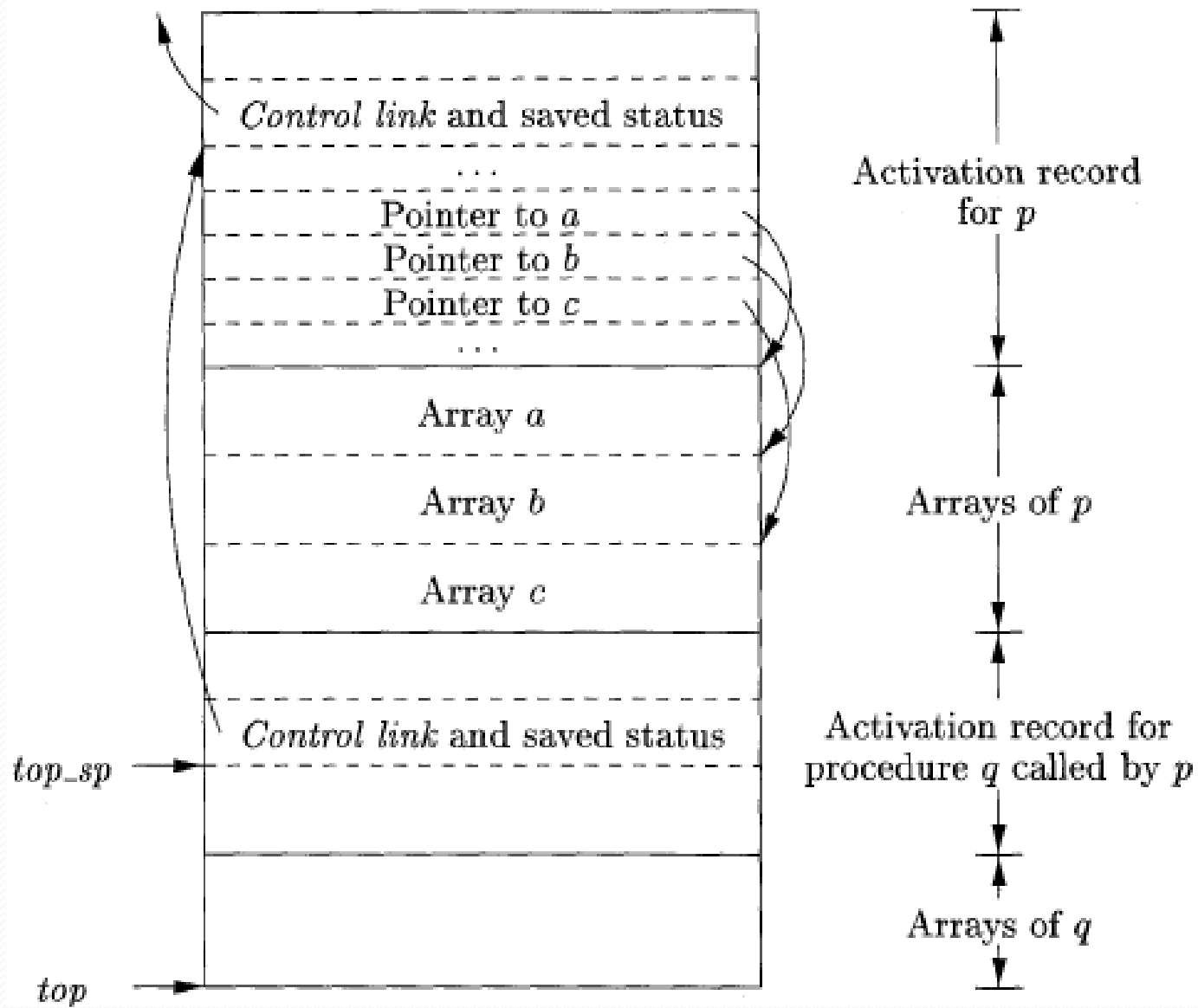
calling sequence

- The caller evaluates the actual parameters
- The caller stores a return address and the old value of *top-sp* into the callee's activation record.
- The callee saves the register values and other status information.
- The callee initializes its local data and begins execution.

corresponding return sequence

- The callee places the return value next to the parameters
- Using information in the machine-status field, the callee restores *top-sp* and other registers, and then branches to the return address that the caller placed in the status field.
- Although *top-sp* has been decremented, the caller knows where the return value is, relative to the current value of *top-sp*; the caller therefore may use that value.

Access to dynamically allocated arrays



ML

- ML is a functional language
- Variables are defined, and have their unchangeable values initialized, by a statement of the form:

```
val (name) = (expression)
```
- Functions are defined using the syntax:

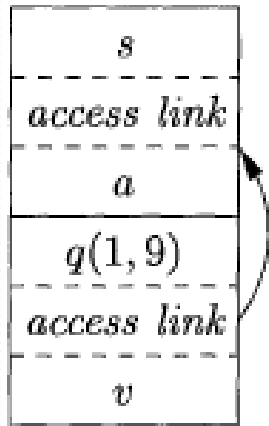
```
fun (name) ( (arguments) ) = (body)
```
- For function bodies we shall use let-statements of the form:

```
let (list of definitions) in (statements) end
```

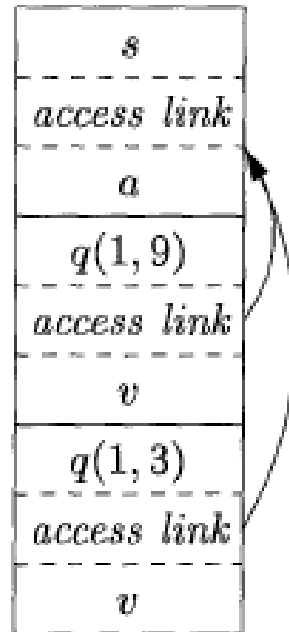

A version of quicksort, in ML style, using nested functions

```
1) fun sort(inputFile, outputFile) =  
    let  
2)     val a = array(11,0);  
3)     fun readArray(inputFile) = ... ;  
4)         ... a ... ;  
5)     fun exchange(i,j) =  
6)         ... a ... ;  
7)     fun quicksort(m,n) =  
        let  
8)         val v = ... ;  
9)         fun partition(y,z) =  
10)            ... a ... v ... exchange ...  
        in  
11)            ... a ... v ... partition ... quicksort  
        end  
    in  
12)        ... a ... readArray ... quicksort ...  
    end;
```

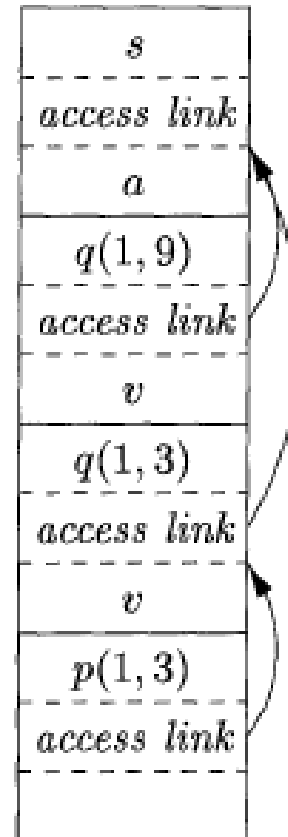
Access links for finding nonlocal data



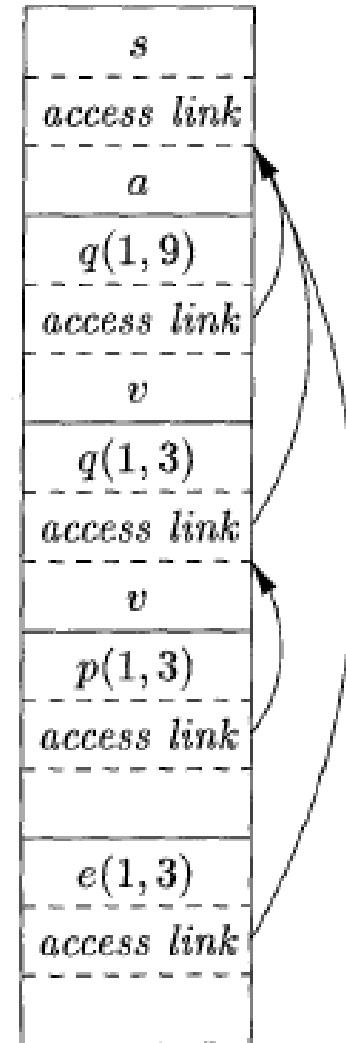
(a)



(b)



(c)

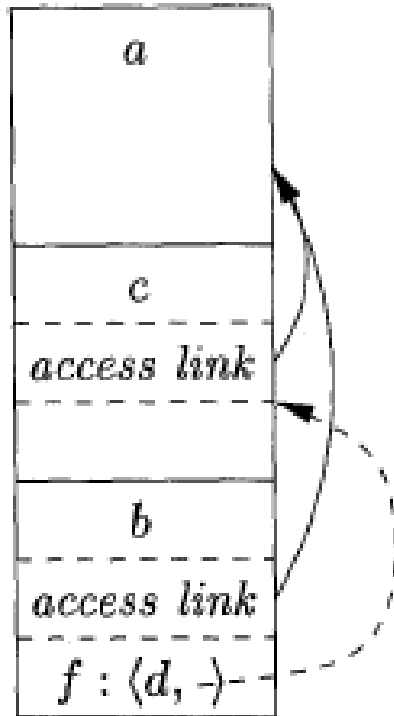


(d)

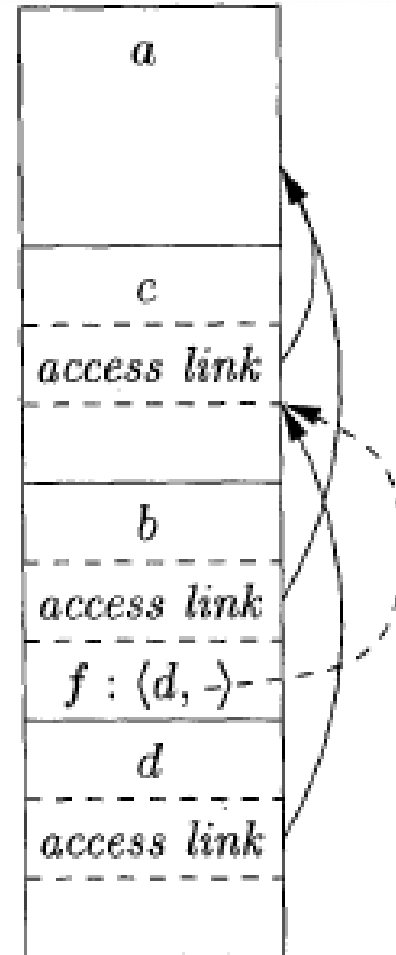
Sketch of ML program that uses function-parameters

```
fun a(x) =  
  let  
    fun b(f) =  
      ... f ... ;  
    fun c(y) =  
      let  
        fun d(z) = ...  
        in  
          ... b(d) ...  
        end  
      in  
        ... c(1) ...  
      end;  
  end;
```

Actual parameters carry their access link with them

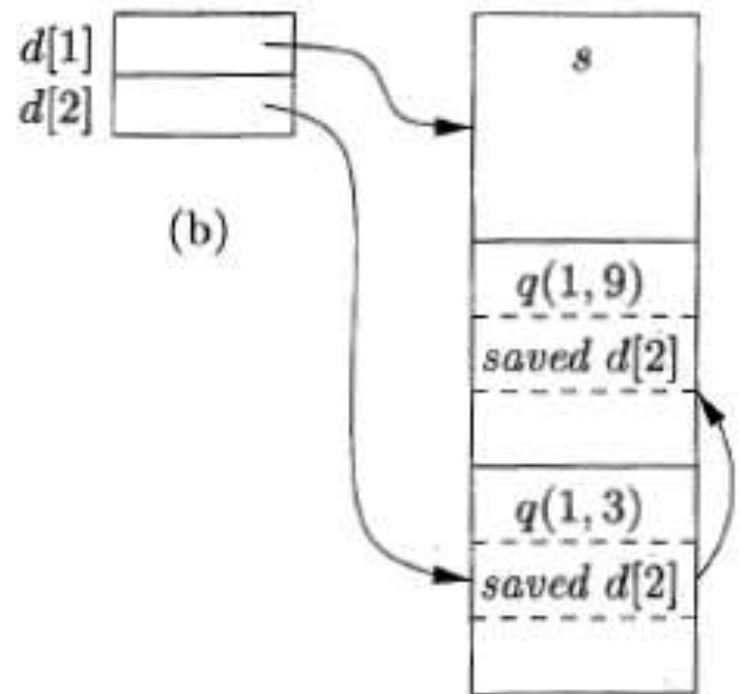
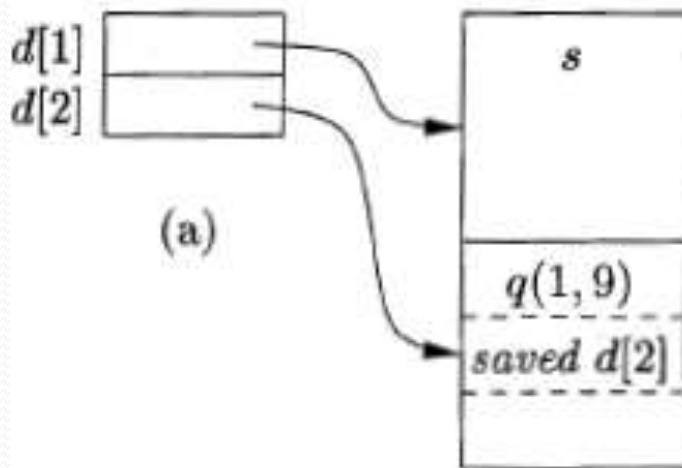


(a)

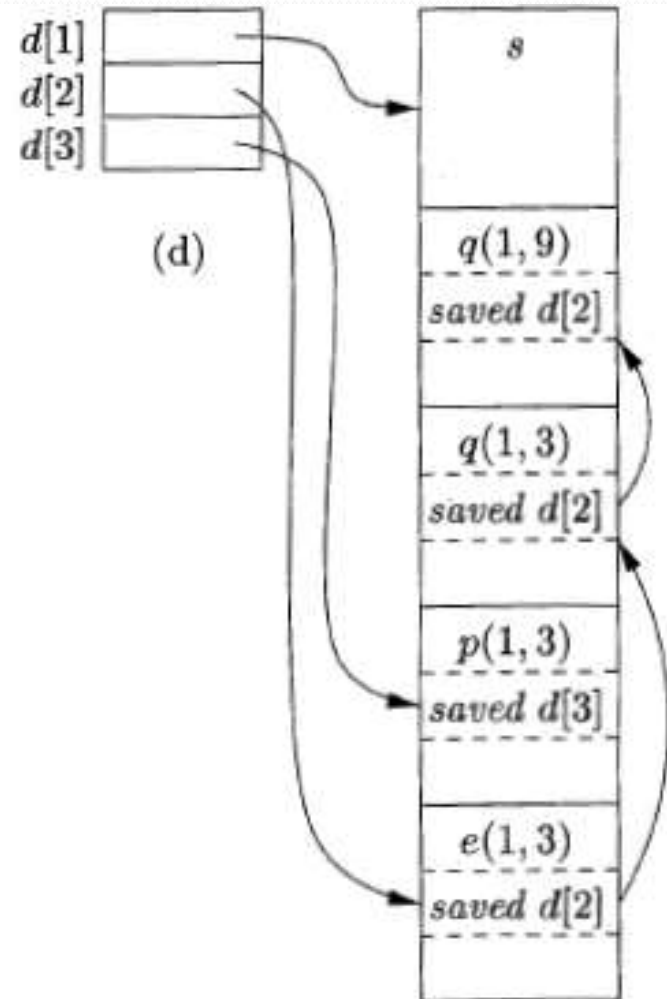
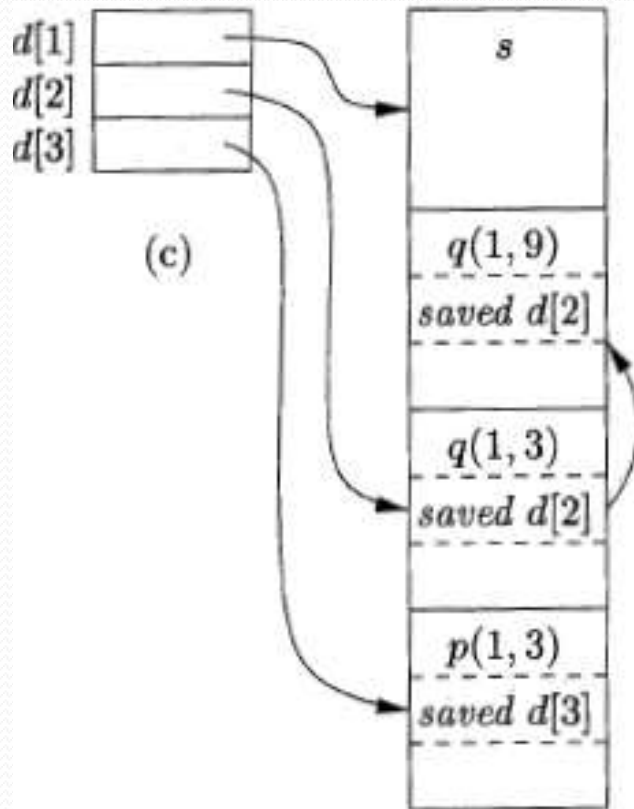


(b)

Maintaining the Display



Maintaining the Display (Cont.)



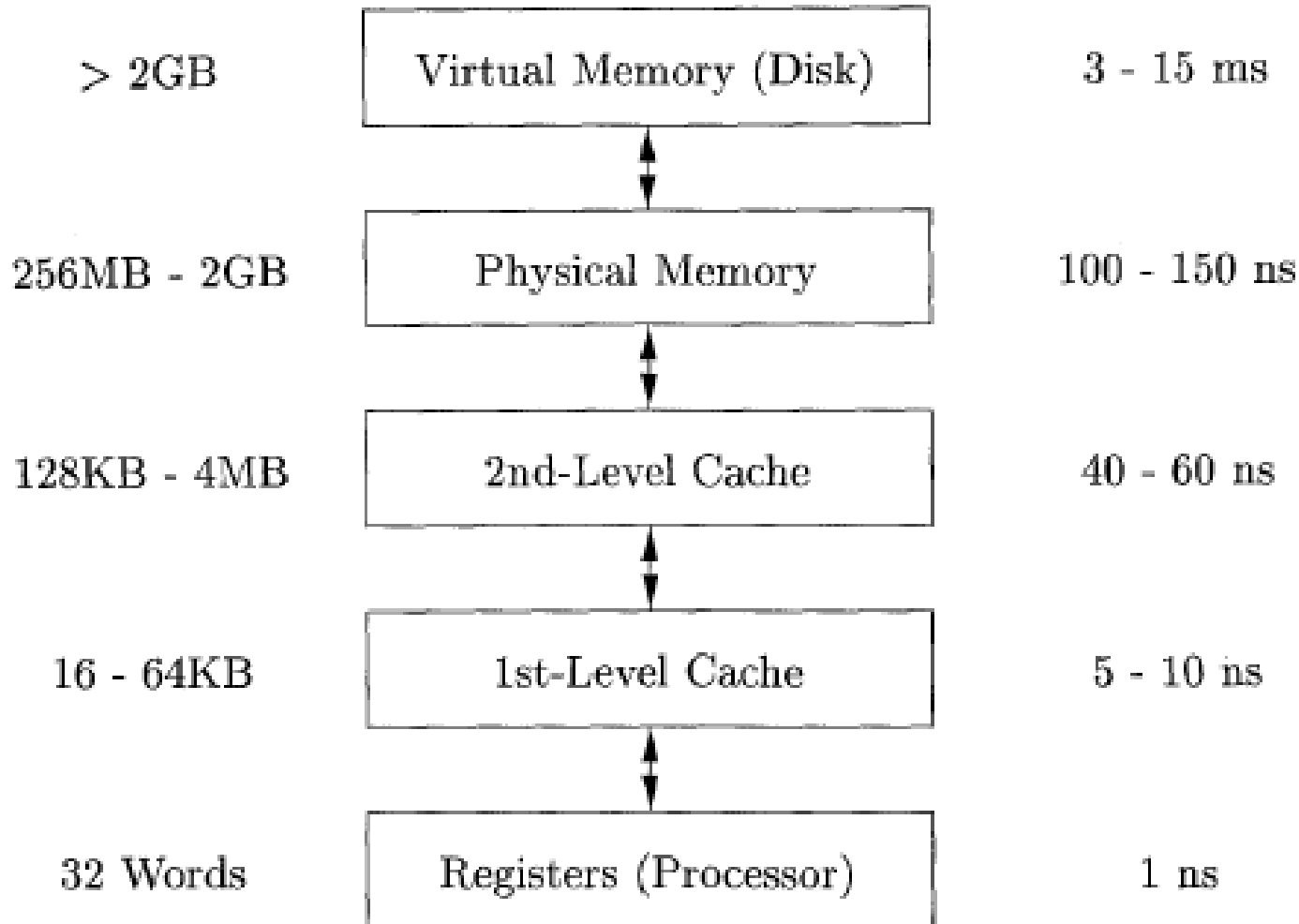
Memory Manager

- Two basic functions:
 - Allocation
 - Deallocation
- Properties of memory managers:
 - Space efficiency
 - Program efficiency
 - Low overhead

Typical Memory Hierarchy Configurations

Typical Sizes

Typical Access Times

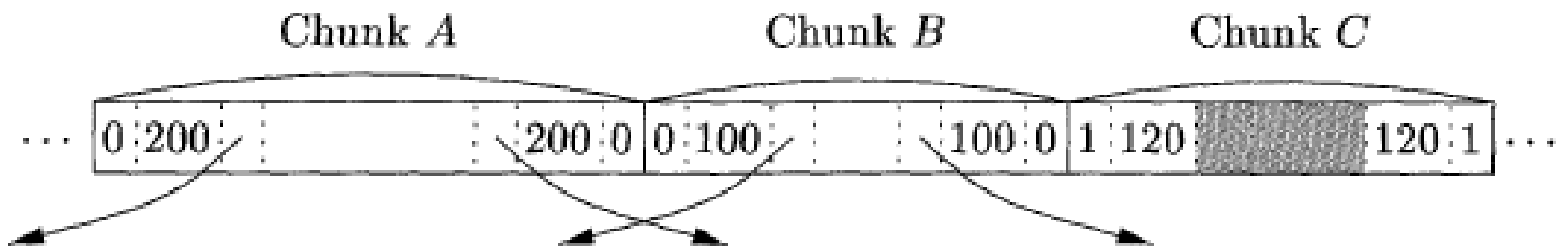


Locality in Programs

The conventional wisdom is that programs spend 90% of their time executing 10% of the code:

- Programs often contain many instructions that are never executed.
- Only a small fraction of the code that could be invoked is actually executed in a typical run of the program.
- The typical program spends most of its time executing innermost loops and tight recursive cycles in a program.

Part of a Heap



Garbage Collection

Reference Counting

Mark-and-Sweep

Short-Pause Methods

The Essence

- Programming is easier if the run-time system “garbage-collects” --- makes space belonging to unusable data available for reuse.
 - Java does it; C does not.
 - But stack allocation in C gets some of the advantage.

Desiderata

1. Speed --- low overhead for garbage collector.
2. Little program interruption.
 - Many collectors shut down the program to hunt for garbage.
3. *Locality* --- data that is used together is placed together on pages, cache-lines.

The Model --- (1)

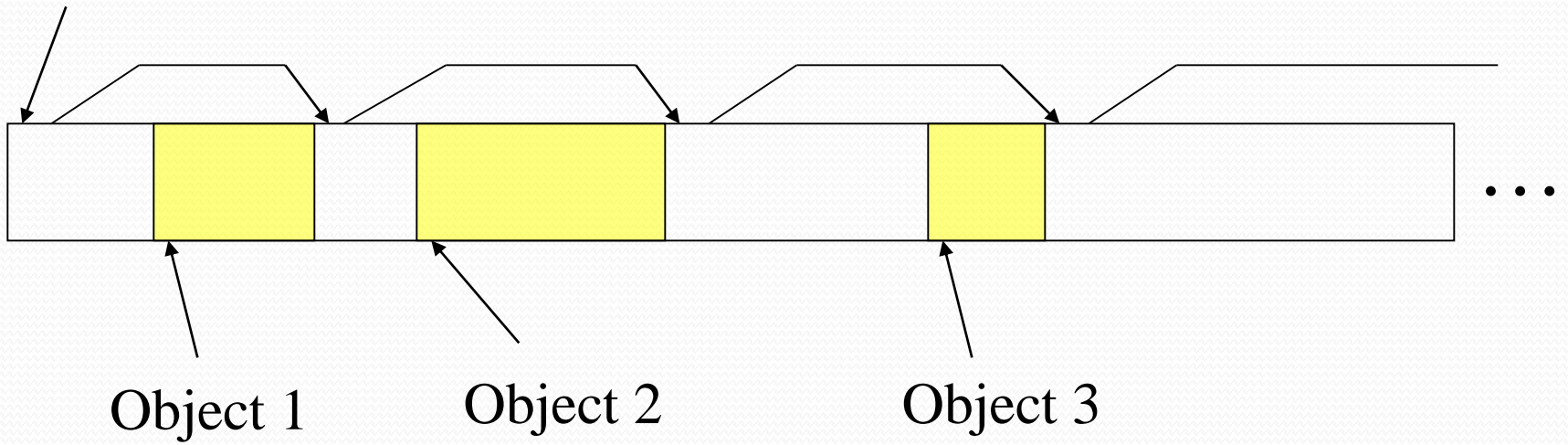
- There is a *root set* of data that is a-priori reachable.
 - *Example*: In Java, root set = static class variables plus variables on run-time stack.
- *Reachable data* : root set plus anything referenced by something reachable.

The Model --- (2)

- Things requiring space are “objects.”
- Available space is in a *heap* --- large area managed by the run-time system.
 - *Allocator* finds space for new objects.
 - Space for an object is a *chunk*.
 - *Garbage collector* finds unusable objects, returns their space to the heap, and maybe moves objects around in the heap.

A Heap

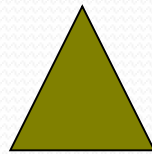
Free List



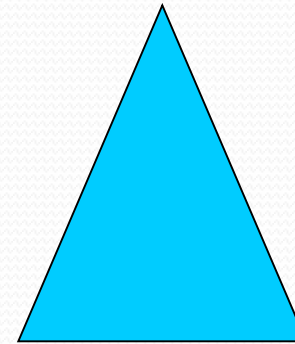
Taxonomy

Garbage Collectors

Reference-
Counters



Trace-
Based



Reference Counting

- The simplest (but imperfect) method is to give each object a *reference count* = number of references to this object.
 - OK if objects have no internal references.
- Initially, object has one reference.
- If reference count becomes 0, object is garbage and its space becomes available.

Examples

```
Integer i = new Integer(10);
```

- Integer object is created with RC = 1.

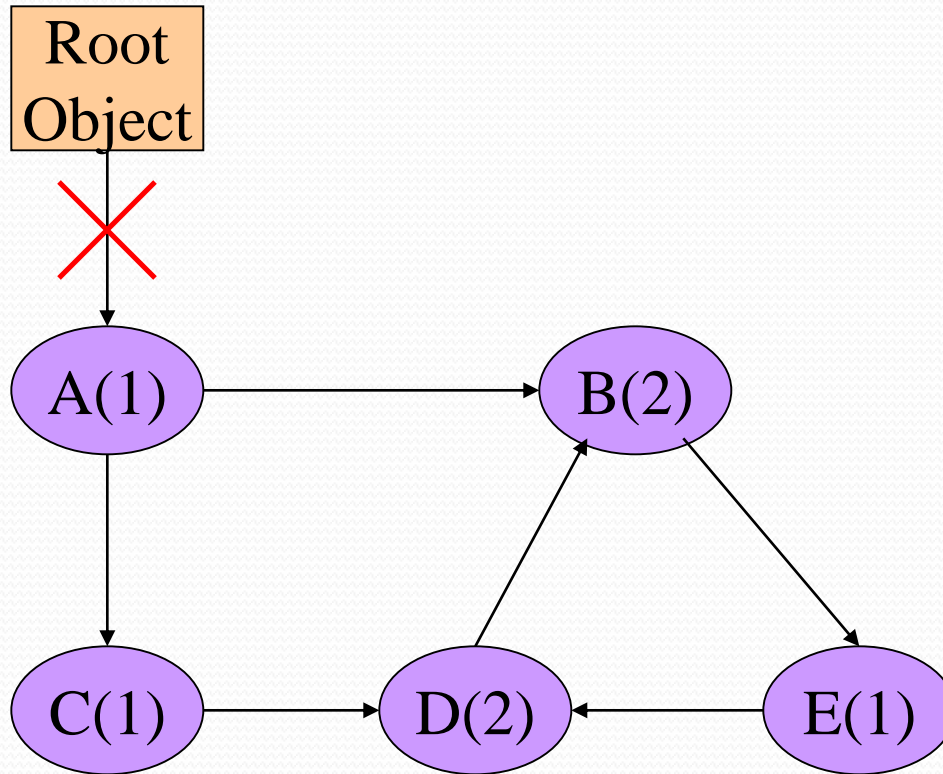
```
j = k; (j, k are Integer references.)
```

- Object referenced by j has RC--.
- Object referenced by k has RC++.

Transitive Effects

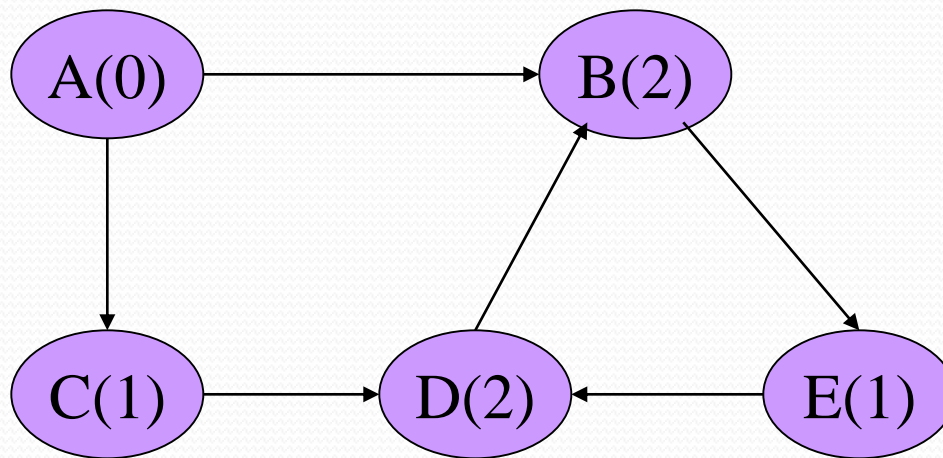
- If an object reaches $RC=0$ and is collected, the references within that object disappear.
- Follow these references and decrement RC in the objects reached.
- That may result in more objects with $RC=0$, leading to recursive collection.

Example: Reference Counting



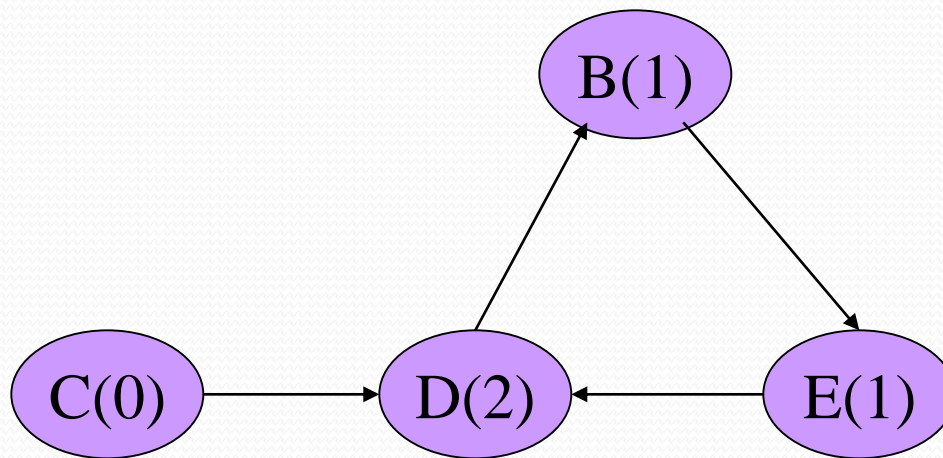
Example: Reference Counting

Root
Object



Example: Reference Counting

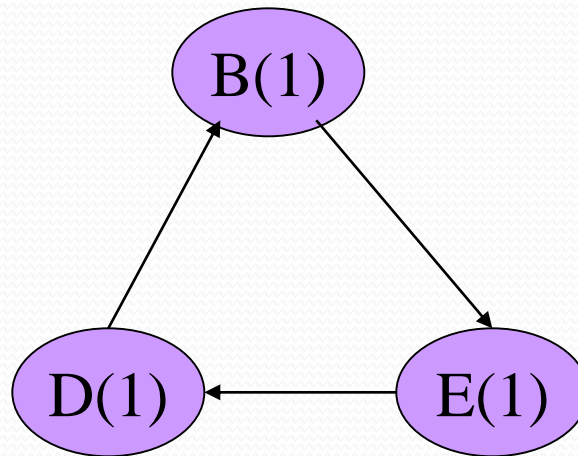
Root
Object



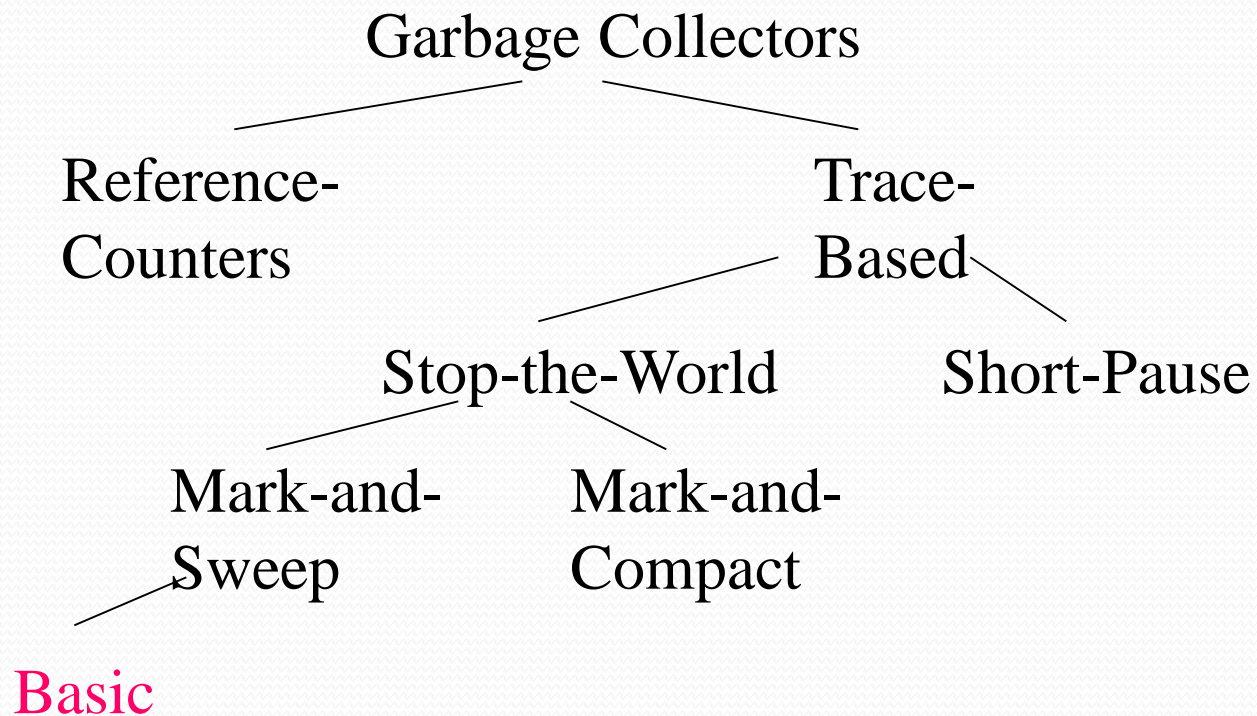
Example: Reference Counting

Root
Object

B, D, and E are garbage, but their reference counts are all > 0 . They never get collected.



Taxonomy



Four States of Memory Chunks

1. *Free* = not holding an object; available for allocation.
2. *Unreached* = Holds an object, but has not yet been reached from the root set.
3. *Unscanned* = Reached from the root set, but its references not yet followed.
4. *Scanned* = Reached and references followed.

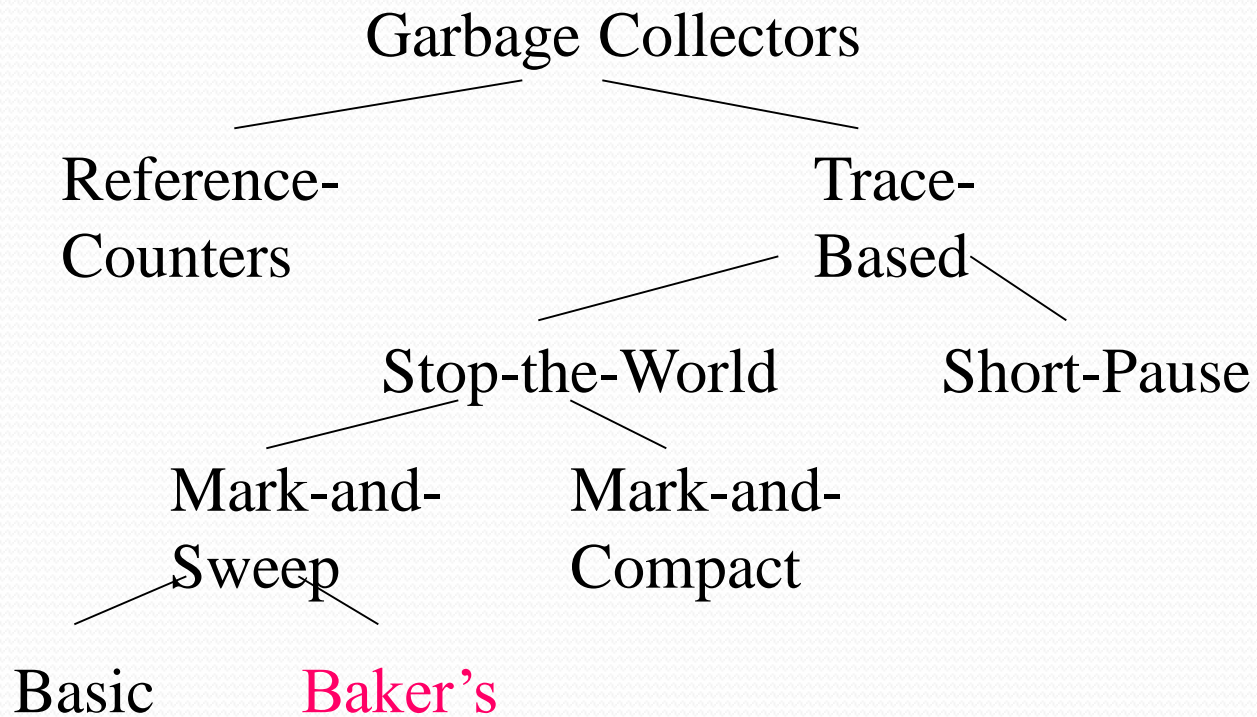
Marking

1. Assume all objects in **Unreached** state.
2. Start with the root set. Put them in state **Unscanned**.
3. **while** **Unscanned** objects remain **do**
 examine one of these objects;
 make its state be **Scanned**;
 add all referenced objects to **Unscanned**
 if they have not been there;
end;

Sweeping

- Place all objects still in the **Unreached** state into the **Free** state.
- Place all objects in **Scanned** state into the **Unreached** state.
 - To prepare for the next mark-and-sweep.

Taxonomy



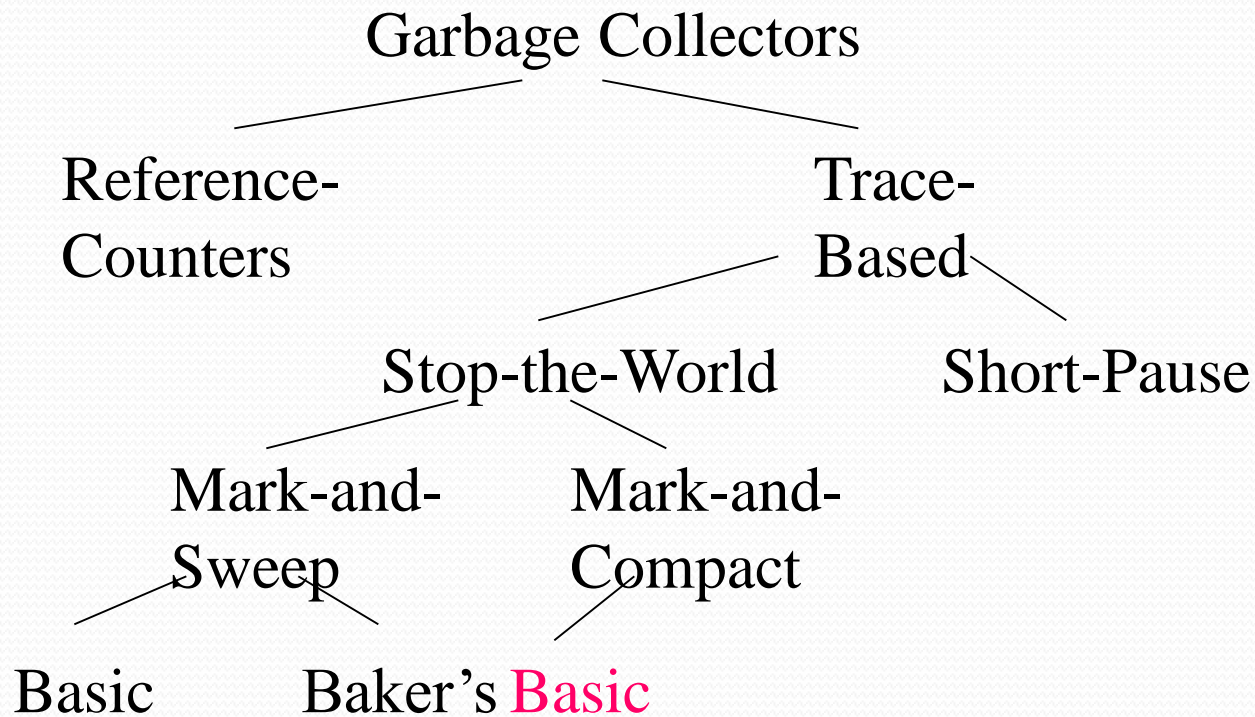
Baker's Algorithm --- (1)

- **Problem:** The basic algorithm takes time proportional to the heap size.
 - Because you must visit all objects to see if they are **Unreached**.
- Baker's algorithm keeps a list of all allocated chunks of memory, as well as the **Free** list.

Baker's Algorithm --- (2)

- **Key change:** In the sweep, look only at the list of allocated chunks.
- Those that are not marked as **Scanned** are garbage and are moved to the **Free** list.
- Those in the **Scanned** state are put in the **Unreached** state.
 - For the next collection.

Taxonomy



Issue: Why Compact?

- *Compact* = move reachable objects to contiguous memory.
- *Locality* --- fewer pages or cache-lines needed to hold the active data.
- *Fragmentation* --- available space must be managed so there is space to store large objects.

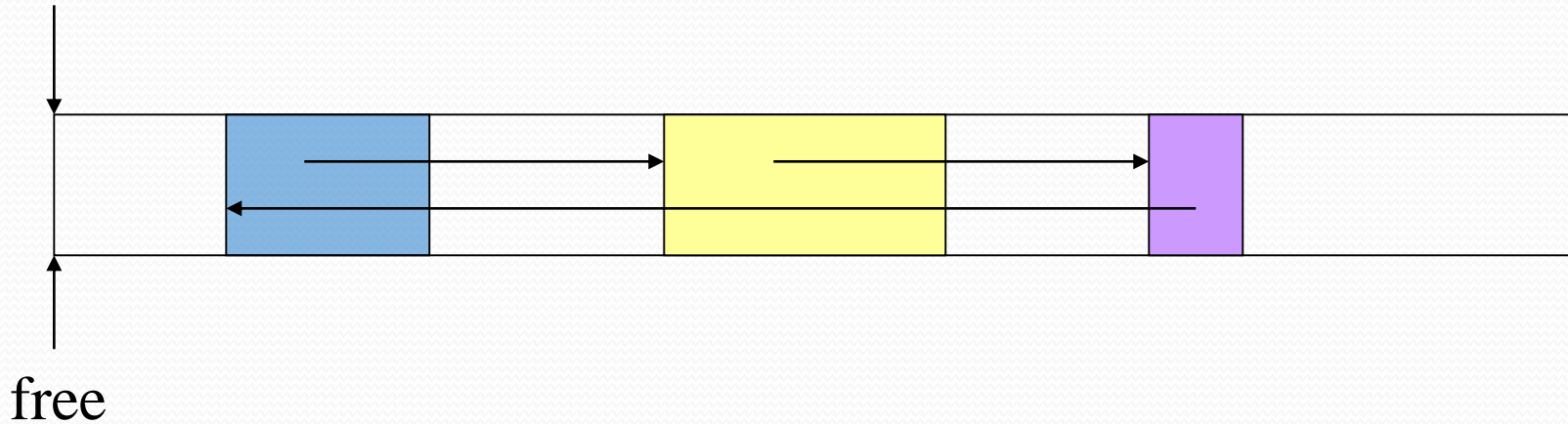
Mark-and-Compact

1. Mark reachable objects as before.
2. Maintain a table (hash?) from reached chunks to new locations for the objects in those chunks.
 - Scan chunks from low end of heap.
 - Maintain pointer *free* that counts how much space is used by reached objects so far.

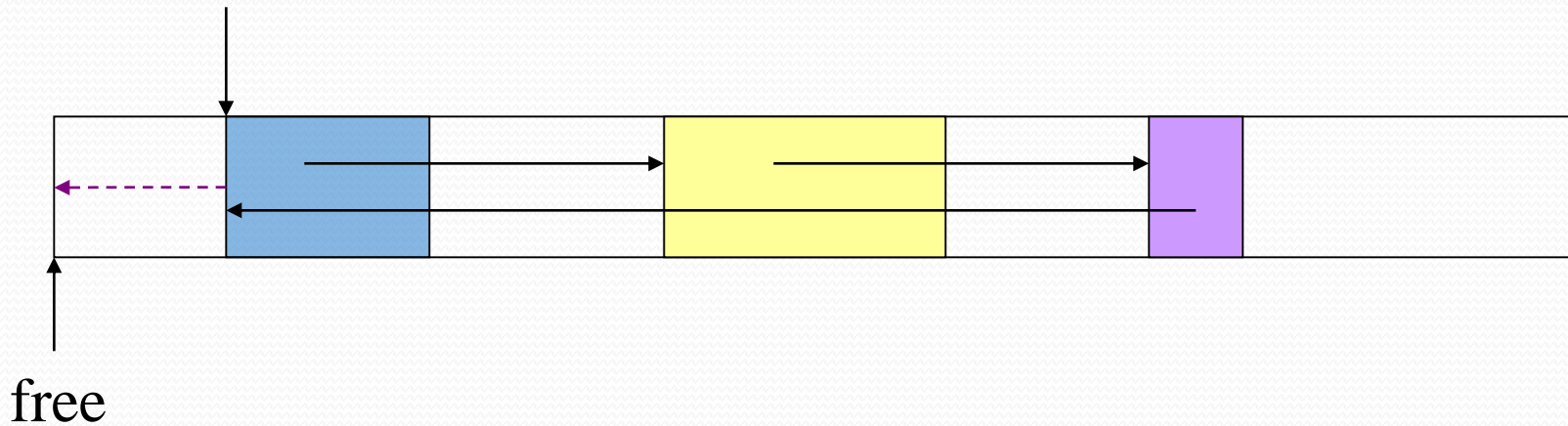
Mark-and-Compact --- (2)

3. Move all reached objects to their new locations, and also retarget all references in those objects to the new locations.
 - Use the table of new locations.
4. Retarget root references.

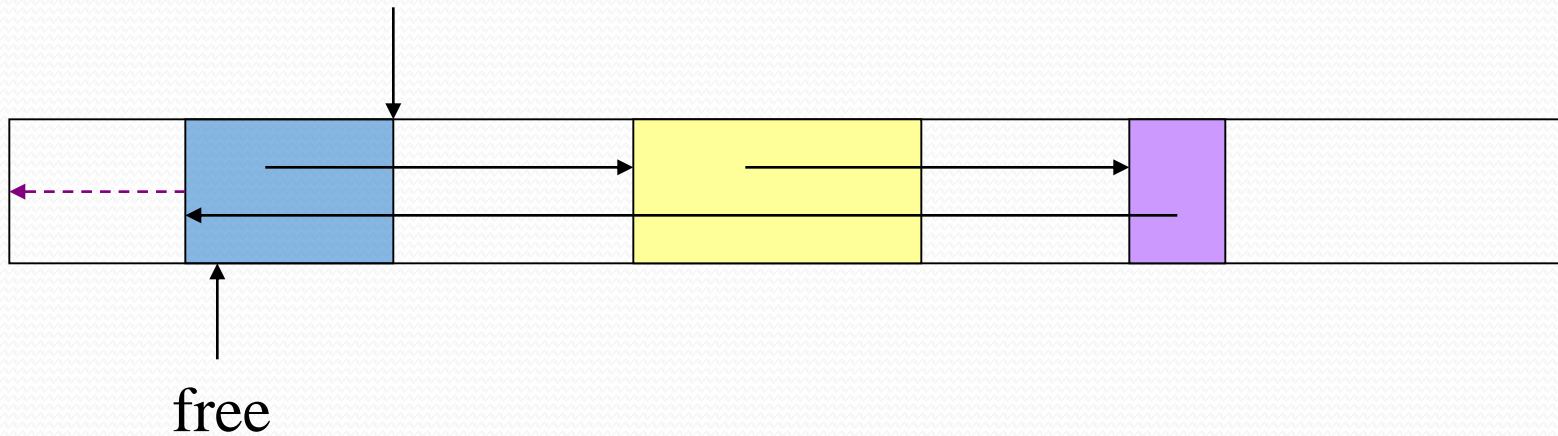
Example: Mark-and-Compact



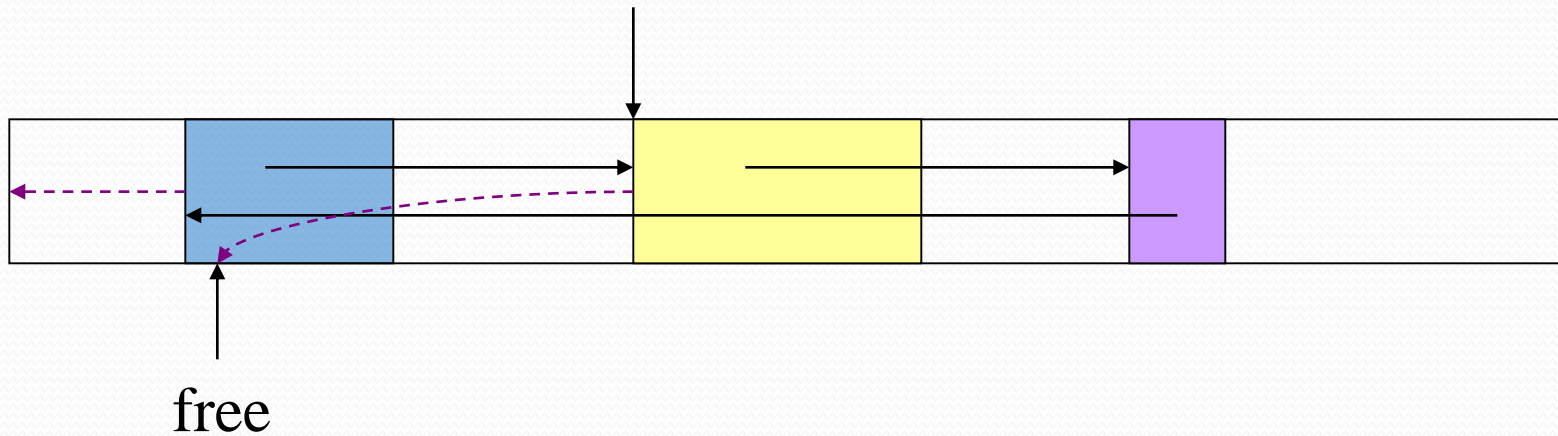
Example: Mark-and-Compact



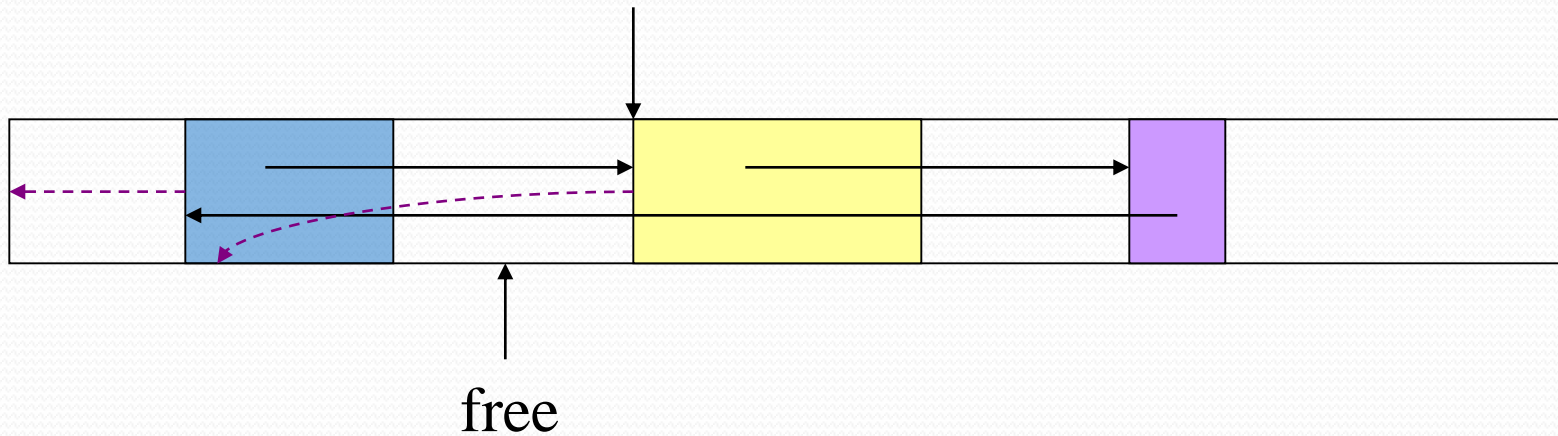
Example: Mark-and-Compact



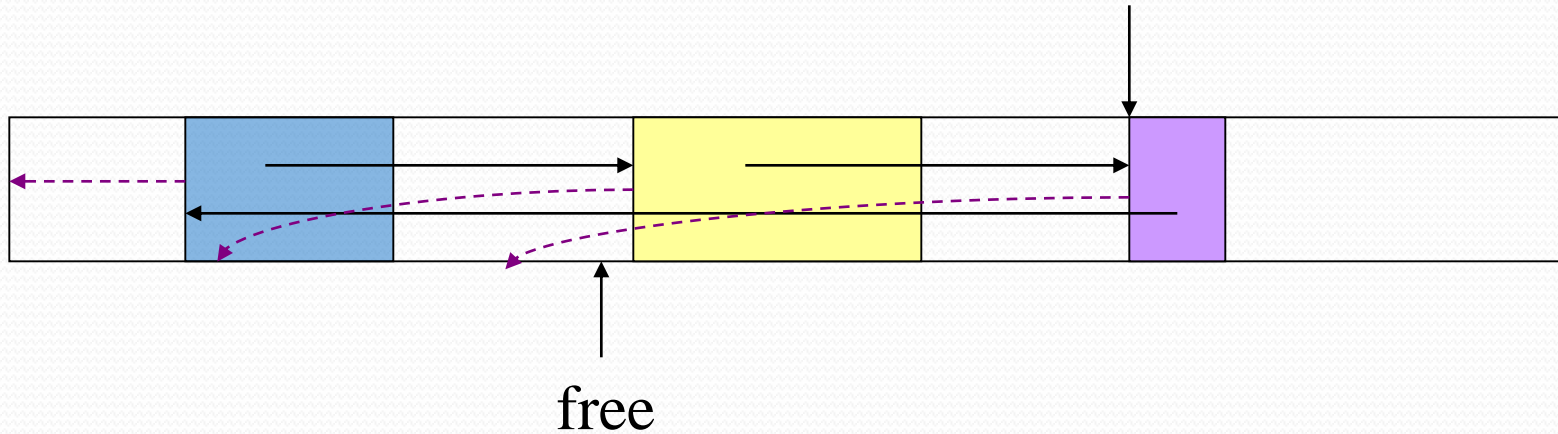
Example: Mark-and-Compact



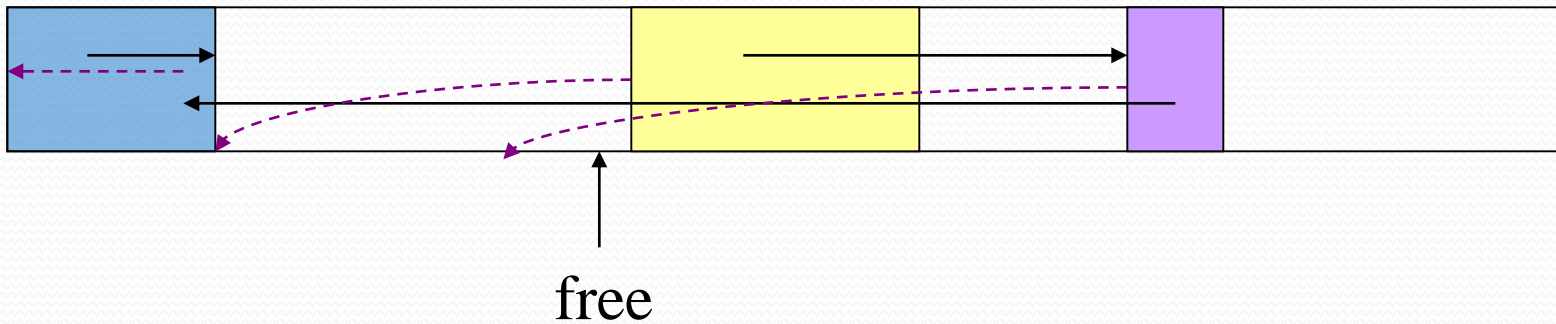
Example: Mark-and-Compact



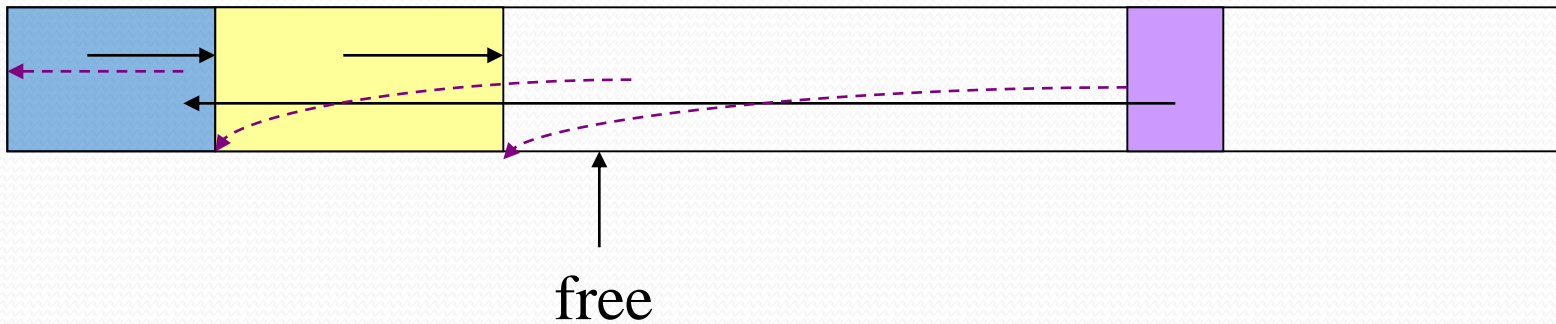
Example: Mark-and-Compact



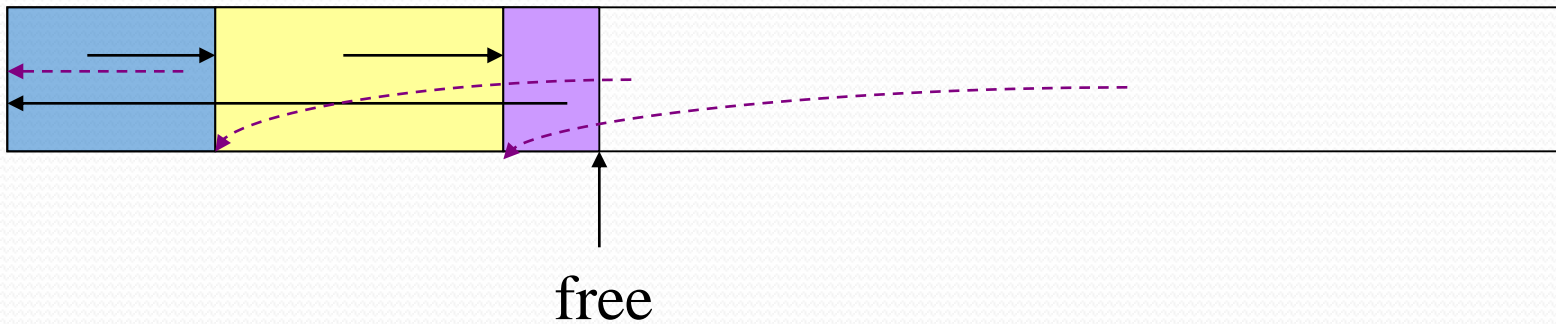
Example: Mark-and-Compact



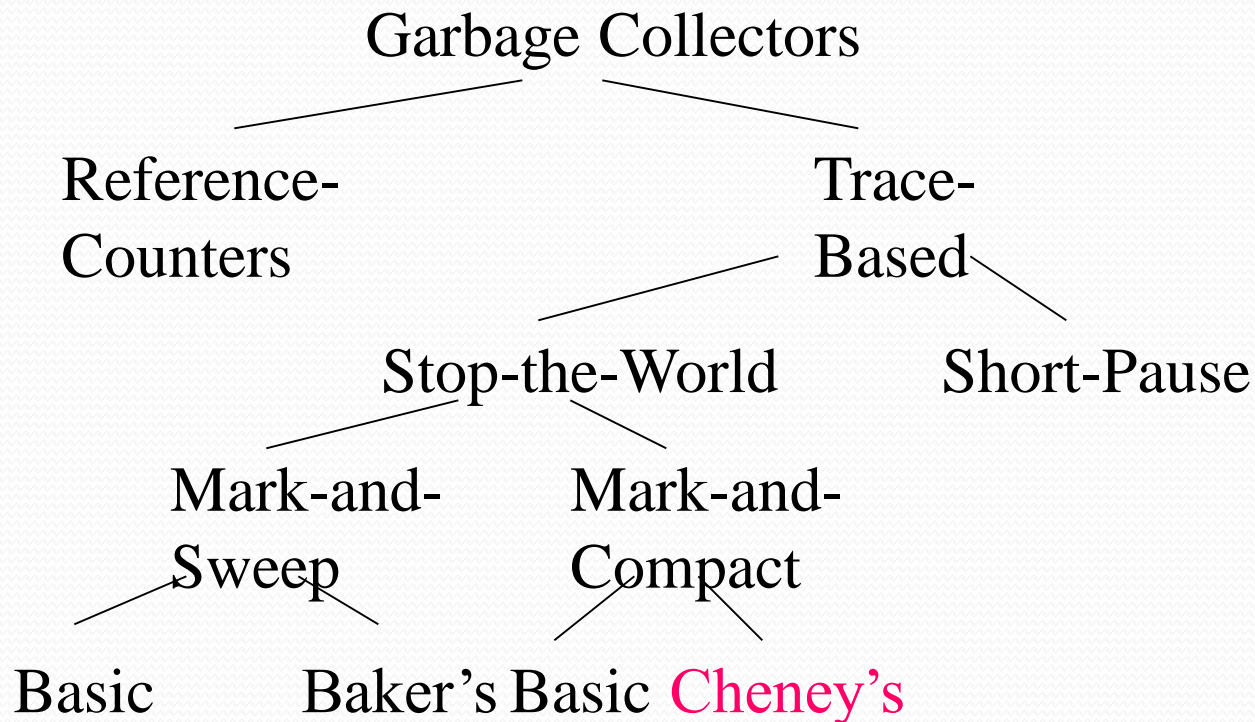
Example: Mark-and-Compact



Example: Mark-and-Compact



Taxonomy

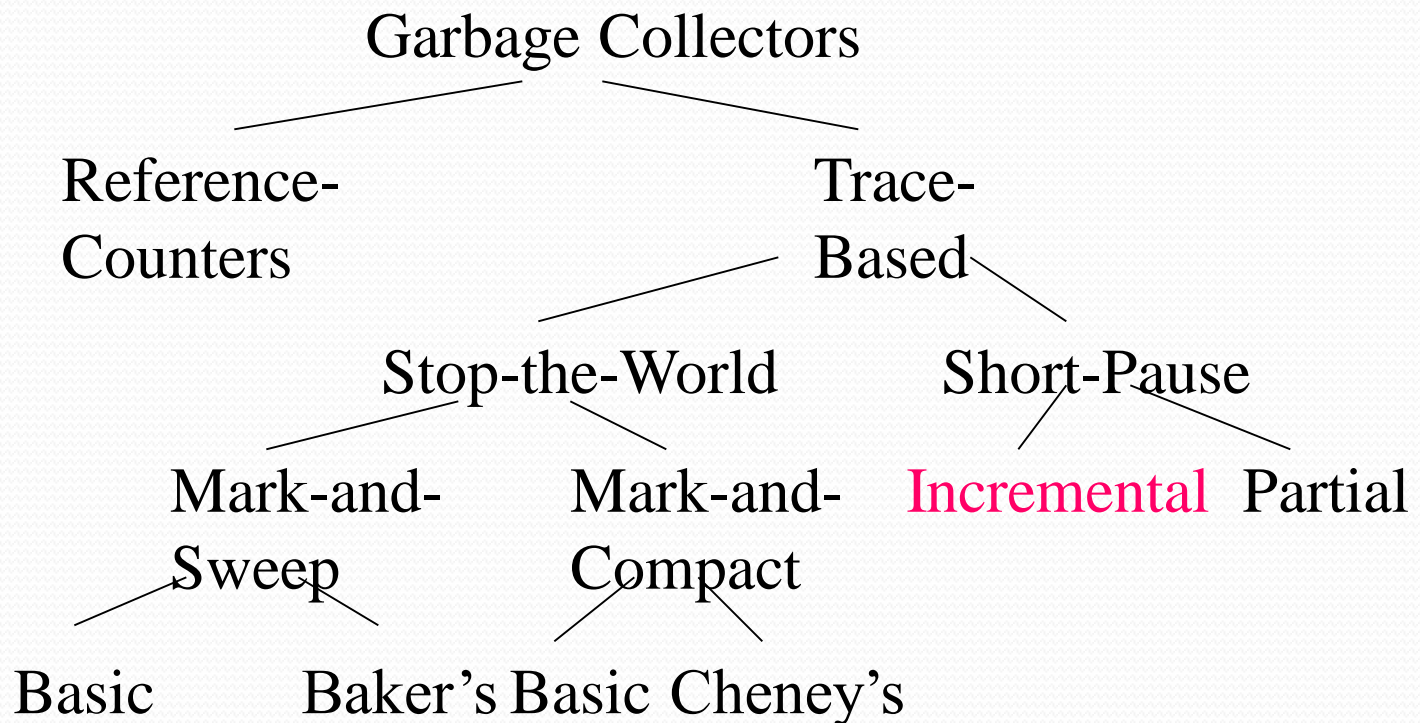


A different Cheney, BTW, so no jokes, please.

Cheney's Copying Collector

- A shotgun approach to GC.
- 2 heaps: Allocate space in one, copy to second when first is full, then swap roles.
- Maintain table of new locations.
- As soon as an object is reached, give it the next free chunk in the second heap.
- As you scan objects, adjust their references to point to second heap.

Taxonomy



Short-Pause Garbage-Collection

1. *Incremental* --- run garbage collection in parallel with *mutation* (operation of the program).
2. *Partial* --- stop the mutation, but only briefly, to garbage collect a *part* of the heap.

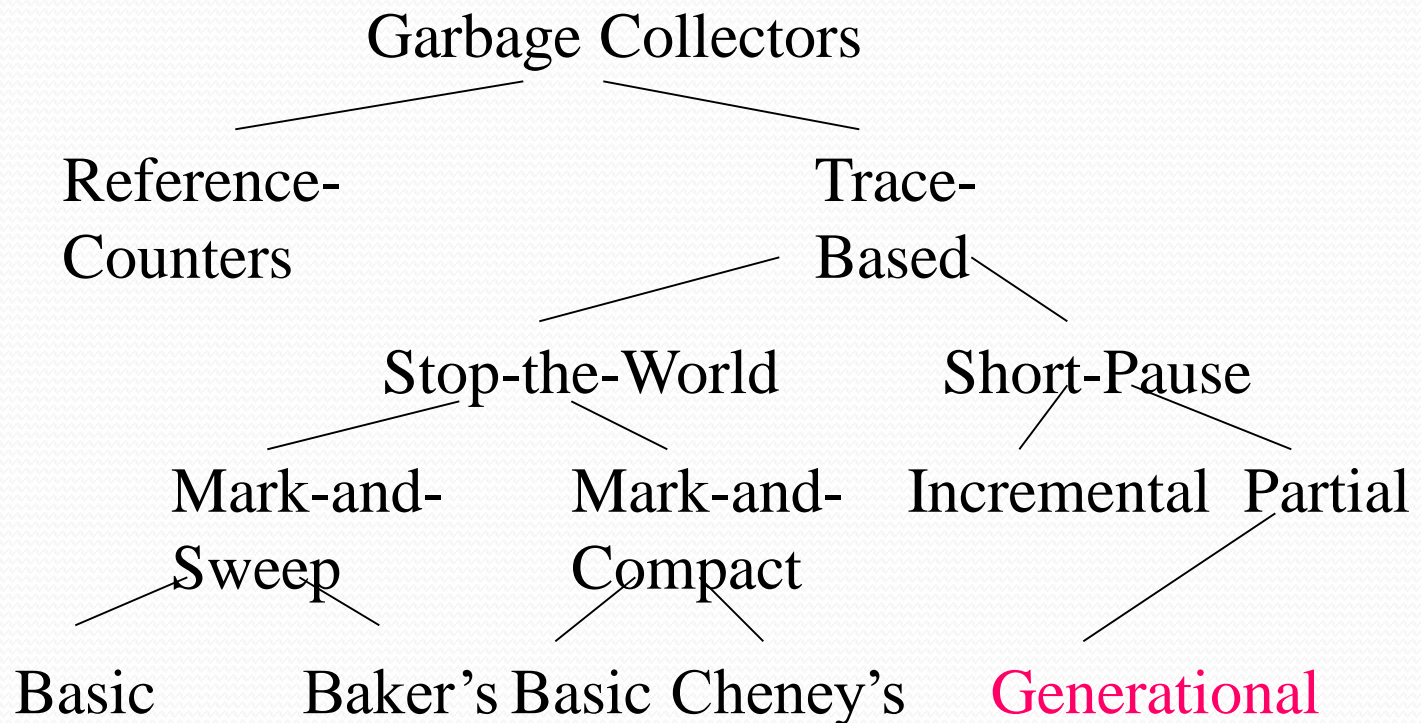
Problem With Incremental GC

- OK to mark garbage as reachable.
- Not OK to GC a reachable object.
- If a reference **r** within a **Scanned** object is mutated to point to an **Unreached** object, the latter may be garbage-collected anyway.
 - **Subtle point**: How do you point to an **Unreached** object?

One Solution: *Write Barriers*

- Intercept every write of a reference in a scanned object.
- Place the new object referred to on the **Unscanned** list.
- **A trick**: protect all pages containing **Scanned** objects.
 - A hardware interrupt will invoke the fixup.

Taxonomy



The Object Life-Cycle

- “Most objects die young.”
 - But those that survive one GC are likely to survive many.
- Tailor GC to spend more time on regions of the heap where objects have just been created.
 - Gives a better ratio of reclaimed space per unit time.

Partial Garbage Collection

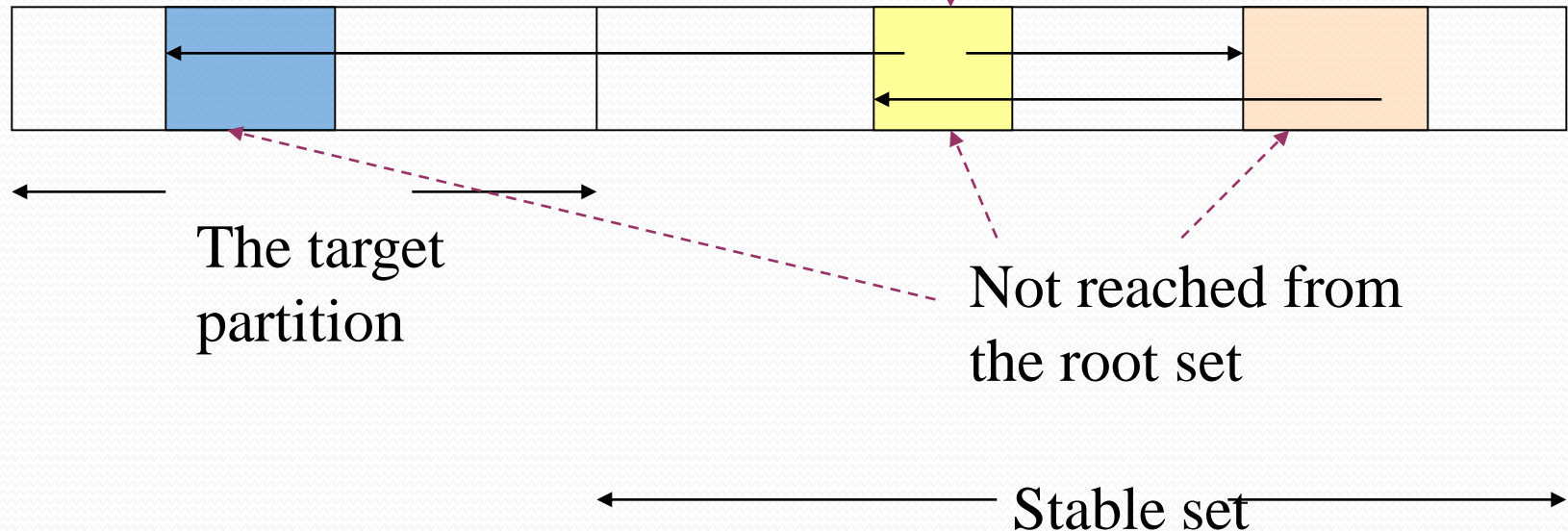
- We collect one part(ition) of the heap.
 - The *target* set.
- We maintain for each partition a *remembered* set of those objects outside the partition (the *stable* set) that refer to objects in the target set.
 - Write barriers can be used to maintain the remembered set.

Collecting a Partition

- To collect a part of the heap:
 1. Add the remembered set for that partition to the root set.
 2. Do a reachability analysis as before.
- Note the resulting **Scanned** set may include garbage.

Example: “Reachable” Garbage

In the remembered set



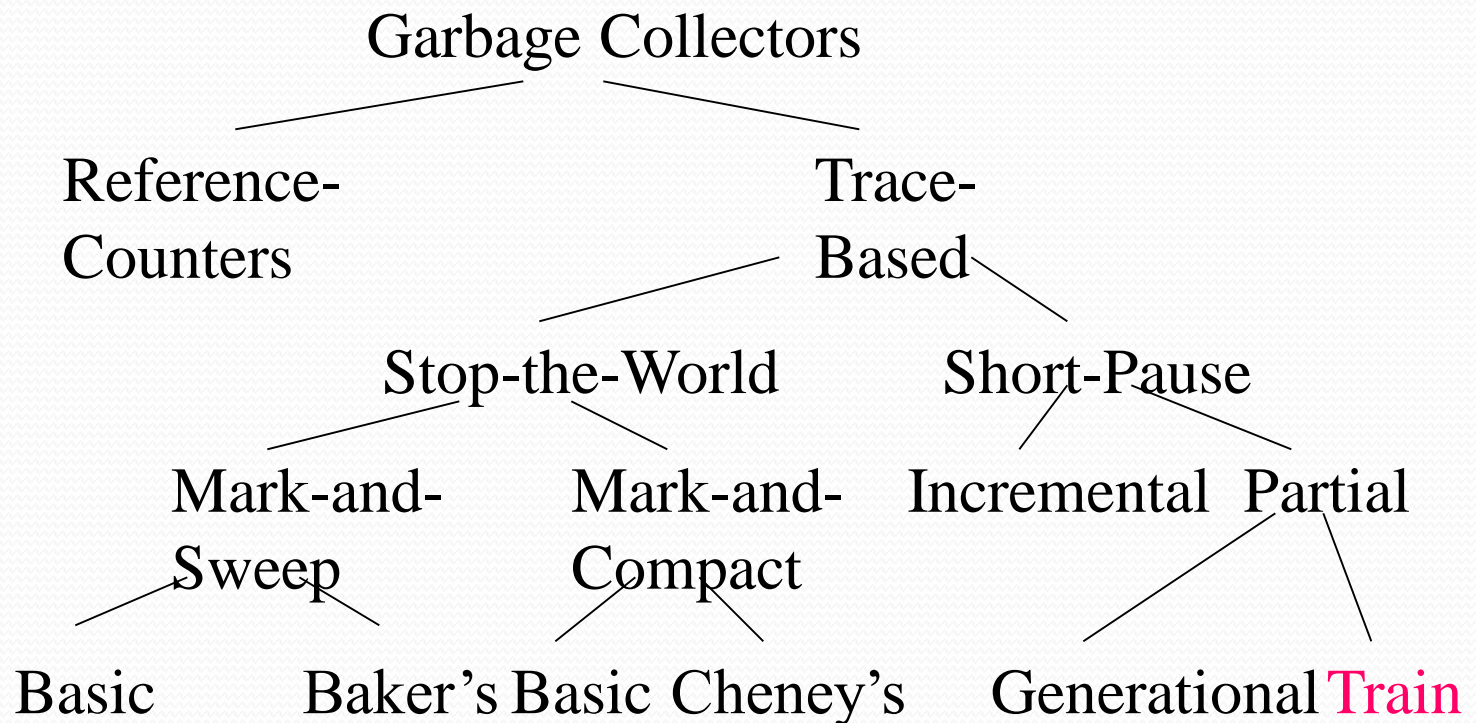
Generational Garbage Collection

- Divide the heap into partitions P_0, P_1, \dots
 - Each partition holds older objects than the one before it.
- Create new objects in P_0 , until it fills up.
- Garbage collect P_0 only, and move the reachable objects to P_1 .

Generational GC --- (2)

- When P_1 fills, garbage collect P_0 and P_1 , and put the reachable objects in P_2 .
- **In general:** When P_i fills, collect P_0, P_1, \dots, P_i and put the reachable objects in $P(i+1)$.

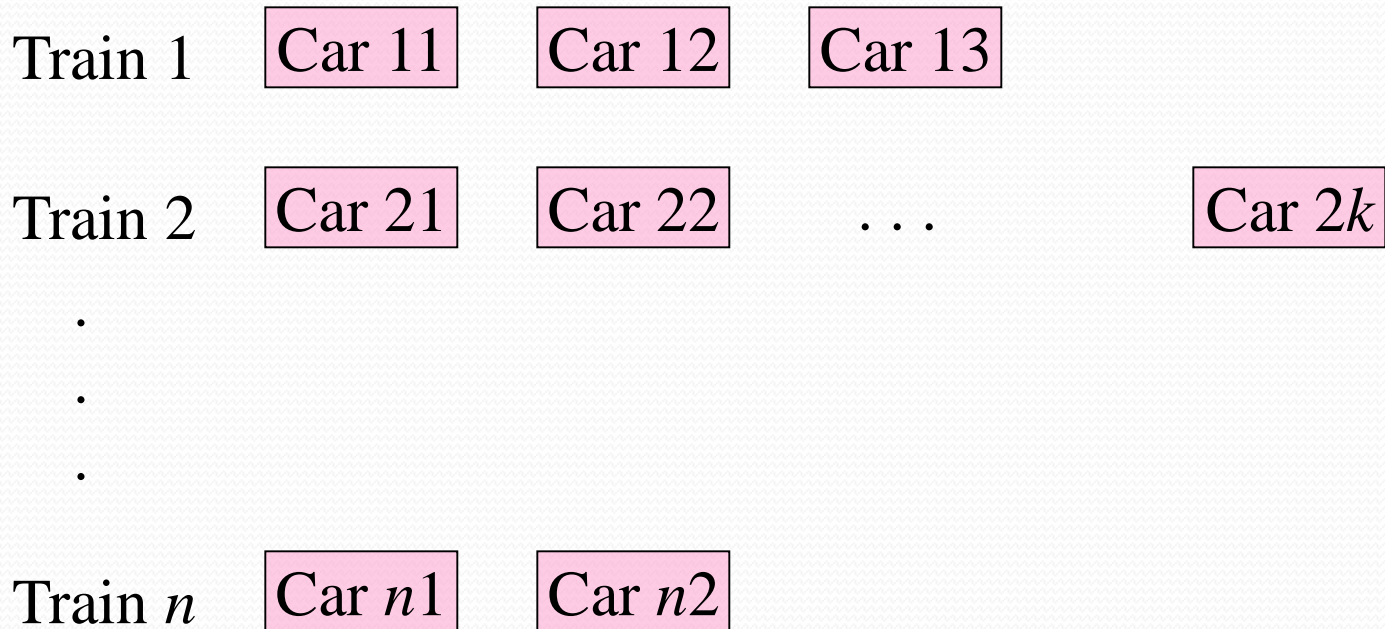
Taxonomy



The Train Algorithm

- Problem with generational GC:
 1. Occasional total collection (last partition).
 2. Long-lived objects move many times.
- Train algorithm useful for long-lived objects.
 - ◆ Replaces the higher-numbered partitions in generational GC.

Partitions = “Cars”



Organization of Heap

- There can be any number of trains, and each train can have any number of cars.
 - You need to decide on a policy that gives a reasonable number of each.
- New objects can be placed in last car of last train, or start a new car or even a new train.

Garbage-Collection Steps

1. Collect the first car of the first train.
2. Collect the entire first train if there are no references from the root set or other trains.
 - **Important:** this is how we find and eliminate large, cyclic garbage structures.

Remembered Sets

- Each car has a remembered set of references from later trains and later cars of the same train.
- **Important**: since we only collect first cars and trains, we never need to worry about “forward” references (to later trains or later cars of the same train).

Collecting the First Car of the First Train

- Do a partial collection as before, using every other car/train as the stable set.
- Move all **Reachable** objects of the first car somewhere else.
- Get rid of the car.

Moving Reachable Objects

- If object **o** has a reference from another train, pick one such train and move **o** to that train.
 - Same car as reference, if possible, else make new car.
- If references only from root set or first train, move **o** to another car of first train, or create new car.

Panic Mode

- **The problem:** it is possible that when collecting the first car, nothing is garbage.
- We then have to create a new car of the first train that is essentially the same as the old first car.

Panic Mode --- (2)

- If that happens, we go into *panic mode*, which requires that:
 1. If a reference to any object in the first train is rewritten, we make the new reference a “dummy” member of the root set.
 2. During GC, if we encounter a reference from the “root set,” we move the referenced object to another train.

Panic Mode --- (3)

- **Subtle point:** If references to the first train never mutate, eventually all reachable objects will be sucked out of the first train, leaving cyclic garbage.
- But perversely, the last reference to a first-train object could move around so it is never to the first car.