

### Instructor

David R. Hanson  
UCC 318  
621-6610  
Office Hours by appointment

### Required Texts

D. R. Hanson, *DECSys-10 C Programmer's Manual*, TR 83-12, Dept. of Computer Science, Univ. of Arizona, Aug. 1983.

### Optional Texts

D. R. Hanson, *Software Tools User's Manual*, TR 81-20, Dept. of Computer Science, Univ. of Arizona, Dec. 1981.

B. W. Kernighan and D. M. Ritchie, *The C Programming Language*, Prentice-Hall, 1978.

B. W. Kernighan and P. J. Plauger, *Software Tools*, Addison-Wesley, 1976.

~ R. M. Graham, *Principles of Systems Programming*, John Wiley & Sons, 1975.

X P. Calingaert, *Assemblers, Compilers, and Program Translation*, Computer Science Press, 1979.

D. W. Barron, *Assemblers and Loaders*, 3rd edition, Elsevier North-Holland, 1978.

✓ N. Wirth, *Algorithms + Data = Programs*, Prentice-Hall, 1976, chap. 5.

### Grading

1/2 programming assignments; 1/6 each of 3 exams (including the final). These weights are subject to minor variation.

### Attendance

Attendance is expected but will not be recorded. Students are, however, fully responsible for all material presented in class. Since much of the material is not available in textbooks, class attendance is strongly recommended.

### Rules

Exams will be scheduled in advance; unless prior arrangements are made, a grade of zero will be recorded for missed exams.

Students are encouraged to discuss programs in general terms and to help one another find bugs in existing programs. Writing code for use by another or using another's code is cheating, however.

Except in extraordinary circumstances, late programs are *not* accepted.

<i>date</i>	<i>topic</i>	<i>assignments</i>	<i>due</i>
Aug 23	introduction; compiling		
25	compiling: grammars	parser	
30	compiling: declarations		
Sep 1	compiling: procedures		
6	compiling: statements		
8	compiling: expressions		
13	compiling: expressions		
15	compiling: expressions		
20	code generation	code generator	parser
22	code generation		
27	code generation		
29	code generation		
Oct 4	exam		
6	linking	link - 450	
11	exam review/linking		
13	linking		code generator
18	linking		
20	no class		
25	linking		
27	linking		
Nov 1	linking	library searching	link
3	linking		
8	exam		
10	debugging	db 300	
15	exam review/linking		library searching
17	debugging		
22	debugging	cdb	
24	Thanksgiving		
29	debugging		db
Dec 1	debugging		
6	debugging		cdb
16	final exam (10:30 am)		

Name: Bill Mitchell

There are four problems on this exam. Each problem is worth 25 points. This is *open book* exam.

Please attempt to solve all problems in the space provided. Show your work in arriving at solutions; partial credit will be given. There is an extra page at the back of the exam. Be brief.

Some of the problems ask for DEC-10 instruction sequences. The following instruction examples provide all of the necessary details.

<i>instruction</i>	<i>effect</i>
move r,x	load the contents of the word at address 'x' into register 'r'
movem r,(r1)	store the contents of 'r' at the address contained in 'r1'
movem r,-5(14)	store the contents of 'r' at the address 5 plus the contents of 14
movei r,x	load the address 'x' into 'r'
movei r,\$4	load the address of the label '\$4' into 'r'
jrst \$4	jump to label '\$4'
jrst (r)	jump to the address given by the contents of 'r'
adjsp r,n	increment the stack pointer in 'r' by 'n' words
adjsp r,-n	decrement the stack pointer in 'r' by 'n' words
add r1,r2	add the contents of 'r2' to the contents of 'r1'

Good luck.

22  
20  
10  
20  

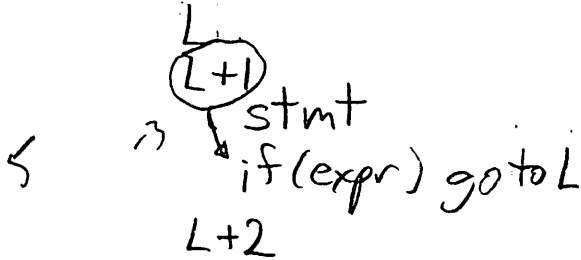
---

72

1. In addition to the **while** and **for** loops, C also has a **do** loop in which the test is at the bottom of the loop:

**do stmt while ( expr )**

(a) (8 pts) Give the 'code template' for the **do** statement, showing the correct placement of the necessary labels.



(b) (12 pts) Write the code necessary for parsing this statement, including the correct semantics (i.e. calls to *walk*, *deflabel*, *jump*, etc.). Make sure your code interacts properly with **break** and **continue** statements.

dostat()  
{

int lab;

lab = genlabel(3);

t = stok();

stmt(lab);

if (t != WHILE)

err("while' expected", "");

else

t = stok();

mustbe('(');

walk(cnode(expr()), lab, 0);

deflabel(lab+2)

mustbe(')');

}

(c) (5 pts) How does adding this statement affect the *first* and *follow* sets?

"do" is added to the First + Follow sets for stmt.



widen is called for all ints

2. C permits expressions like

$a = b < c;$

which sets  $a$  to 1 if  $b < c$  and to 0 otherwise. The type of the result for the comparison operators,  $\&\&$ ,  $\|$ , and unary  $!$  is integer. This feature can be implemented by 'converting' conditional expressions to integer expressions using the code template

```
if (e) goto L
c.cc = 0
goto L+1
L
c.cc = 1
L+1
```

↑  
attempt at  
a "big hint"

where  $c.cc$  is a global variable allocated elsewhere. Describe the changes necessary to implement this feature in our compiler. Be specific; give the relevant details concerning type checking and expression tree traversal, using code fragments or pseudo-code where necessary.

If widen gets <sup>actually, ~~the~~ as in value</sup> a T-COMP node, it adds a node, <sup>#</sup> later intercepted by walk that generates the code shown above.

good

In walk:

Case '#':

$r1 = \text{walk}(p \rightarrow e\_left, 0, 0);$  via additional nodes  
(best added in env)

use about  
template.

Generate code to check value of  $r1$  and produce a 0 or 1 value as result of exprcode in  $r$ .

This value is later used as an integer in any computations.

18

3. Generating code for the DEC-10 is simplified because it has many registers. Suppose, however, that it had only 3 registers. Describe how you would handle this situation, especially when you run out of registers. Your scheme should affect only the bottom half. Be specific about details concerning types, if any, and what changes to the generated code must be made to use your scheme. Keep your scheme simple, but be sure it works for recursive procedures:



In brief, use words in the procedure frame as temporary storage and use one of the three registers as a scratch register. More specifically, the procedure frame will be:



move r1 5, X  
move r6, (5)

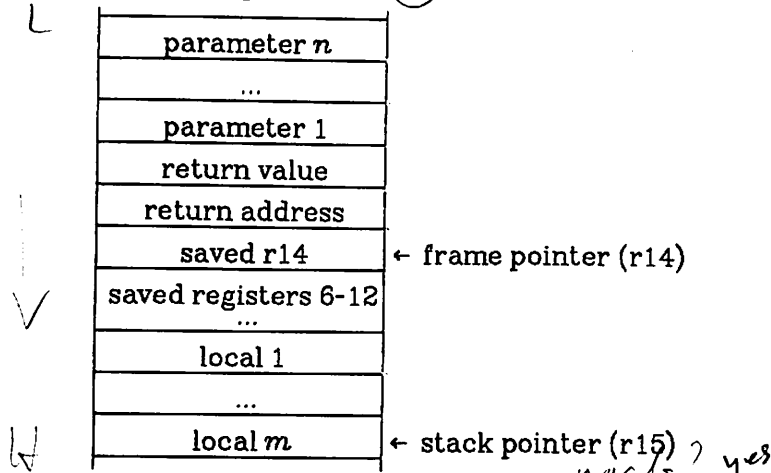
~~unnecessary~~  
how is this space allocated? -5  
just #?  
real

In the new scheme, ralloc manages two registers and some suitably large number of mem. registers. If a real register is available ralloc returns 1 or 2, else it returns a negative number corresponding to a memory register.

When the bottom half encounters a <sup>read reference</sup> neg. reg. number it generates code to get the right value out of the stack frame and puts it into the scratch register which is then used in a computation. Write references do the reverse, calculating the result in the scratch register + then storing it back in the frame.

(20)

4. Suppose the frame layout for procedures was as follows.



Note that the parameters are in the *reverse* order. Assuming that the DEC-10 does *not* have subroutine calling instructions, give appropriate instruction sequences for a call to an integer function *f* with *n* parameters and *m* locals, the entry sequence generated by *procbeg* for *f* and the exit sequence generated by *retcode* for *f*. As suggested by the diagram, your sequences must preserve the contents of registers 6-12, but may use registers 2-5 as desired. You do *not* have to write any C code; just give the relevant instruction sequences. If you are unsure of specific details of the DEC-10 instruction set, use English as necessary to make the intent of your sequences clear.

Calling sequence:

The value for each parm is calculated and put on the stack via: in correct order

```

move r3(sp)
adjsp 15,1
    
```

after parms are on stack

```

adjsp 15,2
movei 2,$1+1
movem 2,(15)
movei 3,f
$1: jrst (3)
    
```

Entry Sequence:

```

adjsp 15,7
movem 0,-7(15)
    7,-6(15)
    ...
movem 12,0(15)
adjsp (m)15
    
```

When do you not use reg. 14?  
-5

you can combine movem then adjsp

Return:

```

Assume ret val is in r1
movei 1,-(m+7+1)(15)
move 2,(1) # ret addr in 1
adjsp -(m+7+2+n),15
jrst (2)
    
```

assumes sub is called w/ correct # of args.

can be in caller.

return 14?

20

1. (a) The code template is

```
L      stmt
L+1   if (expr) goto L
L+2
```

(b) The case

```
case DO:
    dostat(loop = genlabel(3));
    break;
```

is added to *stmt*, where *loop* is the current loop handle. *dostat* is

```
dostat(lab)
int lab;
{
    t = gtok();
    deflabel(lab);
    stmt();
    deflabel(lab + 1);
    mustbe(WHILE);
    mustbe('(');
    walk(cnode(expr()), lab, 0);
    mustbe(')');
    deflabel(lab + 2);
}
```

(c) **do** affects the *first* and *follow* sets in just like **while** and **for**. **do** must be added to *first(stmt)* and to *follow(stmt)*.

2. 'Converting' a conditional expression to an integer expression can be done in *rvalue*:

```
struct node *rvalue(p)
struct node *p;
{
    ...
    if (p->e_type == T_COND)
        p = node('?', T_INT, p, NULL);
    return (p);
}
```

could use overloaded C++

When the ? operator is encountered in *walk*, the appropriate code is generated for the left operand, constructing the trees as necessary, and the result of computing the *rvalue* of *c.cc* is returned:

```
...
case '?':
    walk(p->e_left, lab = genlabel(2), 0);
    p1 = node(O_ID, T_INT|T_ADDR, lookup("c.cc"), NULL);
    walk(node(O_ASGN, T_INT, p1, constant("0", T_INT)), 0, 0);
    jump(lab + 1);
    deflabel(lab);
    walk(node(O_ASGN, T_INT, p1, constant("1", T_INT)), 0, 0);
    deflabel(lab + 1);
    r = walk(rvalue(p1), 0, 0);
    break;
```

3. A simple approach to dealing with a limited supply of registers is to provide a fixed number of temporaries and allocate them as if they were registers. When emitting an

instruction, it may be necessary to load temporary values into real registers and to store the previous values. Thus, *exprgen* should return an index into an array of resources, which may either registers or temporaries. This is necessary so that when a register is 'spilled' into a temporary the change affects all uses of that register. In order to work properly for recursive procedures, the temporaries must be located in the frame. Either a fixed number of temporary locations can be allocated in every frame (above the locals, for example), or only enough space for those temporaries actually used can be allocated. In the latter case, the size of the local frame is not known until the end of the procedure, which may require that the entry point be at the bottom of the procedure or that a constant containing the frame size be generated and referenced by the entry sequence.

4. The 'frame operations' are

```
procedure call:
    adjsp 15,n+3          ; allocate space for arguments & linkage info
    ...
    movem r,-i-2(15)     ; argument i
    ...
    movei 2,L            ; get return address
    jrst  f              ; transfer to procedure
L:    move  r,-2(15)     ; get return value
    adjsp 15,-n-3       ; remove arguments
```

```
procedure entry:
f:    movem 2,-1(15)     ; save return address
    movem 14,0(15)      ; save frame pointer
    move 14,15          ; establish new frame pointer
    adjsp 15,m+7        ; allocate space for locals & registers
    movem 6,1(14)       ; save register 6
    ...
    movem 12,7(14)      ; save register 12
```

```
procedure exit:
    movem r,-2(14)      ; set return value
    move 6,1(14)        ; restore register 6
    ...
    move 12,7(14)       ; restore register 12
    move 15,14          ; deallocate local portion of the frame
    move 14,0(15)       ; restore frame pointer
    move 2,-1(15)       ; get return address
    jrst (2)            ; return to caller
```

Name: Bill Mitchell

There are four problems on this exam. Each problem is worth 25 points. This is *closed book* exam.

Please attempt to solve all problems in the space provided. Show your work in arriving at solutions; partial credit will be given. There is an extra page at the back of the exam. Be brief.

Good luck.

Problem 4 asks about changes to the code generator, *cgen.c*. The routines in *cgen.c* and their current output are as follows.

<i>routine</i>	<i>output</i>
defconst(p)	.seg lit .def -s \$L lit offset + char, int, float constants: bit pattern char * constants: 0331100000001 lit offset + 9-bit bytes, 4/word
deflabel(lab)	.def -s \$lab segment offset +
exprcode(op,r1,r2,lab,p)	code for op
global(p)	.def p->s_name data offset +
jump(lab)	jrst \$lab
procbeg(p)	.seg code 0331100000000 lit offset + .seg lit p->s_name given in 9-bit bytes, 4/word .seg code .def p->s_name code offset + jsp 4,c.savr adjsp 15,size of locals (procbeg must also compute s_offset values for local variables)
procend(x)	no output
progbeg()	no output
progend(x)	.len code length .len data length .len lit length
retcode(r)	move 1,r jsp 4,c.retr
rfree(r)	no output free resource r

25  
25  
~~17~~  
19  
86

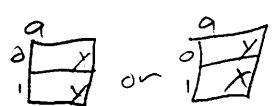
good job!

1. Named (and unnamed) common in Fortran permits separately compiled subroutines to access the same storage, e.g. in

```
subroutine A      subroutine B
common /a/ x      common /a/ y
...
end              ...
end              end
```

x in A and y in B refer to the same cell, which is resolved during linking.

(a) (10 points) Named common *cannot* be handled by simply making the common block name (a in the example above) a segment. Explain why not.

IF "a" were a segment + references to x were generated as "a 0 +" + refs to y were generated as "a 0 +" when loaded, x or y would lie at different locations because of relocation, eg: 

common.name.len

(b) (15 points) Can you devise a way to handle common using the *existing* facilities of link? Explain. Feel free to assume and state compiler and loader conventions to support your scheme, if that is necessary.

15  
10  
The basic problem is to avoid relocation. It <sup>almost</sup> seems that having the compiler use a separate segment for each common and then outputting .len w/o will do the trick. The problem w/ this is that the loader won't know how long the common blocks are. To solve this, modify the compiler to make common segment names of the form "common.name.len", and let the loader sort out everything. (Basically, the loader must find the longest length and use it.)

actual name of common  
generate refs to these segments for vars which lie in it.

Note that the arbitrary limitation on segment name length must be removed, if source-lang. limitations are to be avoided.

25

2. Suppose you have an object file, *table.o*, that contains routines and data for a general-purpose symbol table manager. *table.o* defines the symbols *lookup*, *install*, and *initialize*, which are the routines for accessing the table. You are building a system in which separate symbol tables are needed in each of two modules, *command.o* and *lex.o*. Describe how you could use *link* to construct the final object file, named *prog.out*. You may add new command line options to *link*, but you may not change any of the object files or *link* object code in any way.

```
link -S "lookup,install,initialize" command.o table.o -o command.o.o
link -S "lookup,install,initialize" lex.o table.o -o lex.o.o
link command.o.o lex.o.o -o prog.out
```

-S "idlist" causes the <sup>P?</sup>SUPPRESS bit to be set for each identifier in idlist. (Same as described in paper, I think.)

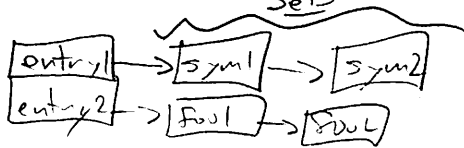


3. Suppose the entries in a library are *not* in topological order, but the index at the beginning of the library contains lines of the form

# symbol name  
- name1 name2

Lines of the first form indicate that *symbol* is defined in entry *name*; lines of the second form indicate that entry *name1* requires symbols defined in entry *name2*. Give *pseudo-code* for *lib(file)*, which completely searches the library in *file* in *one sequential pass* over *file*. Be sure to describe the data structure that results from reading the index; do *not* give the details of reading the index itself, however.

Reading the index produces <sup>①</sup> a topologically sorted list of symbol dependencies built from the "name1 name2" information, and <sup>②</sup> a table mapping entry-definition information, e.g.:



how does -name1 name2 lines give this list?  
(by {unrolling} the respective lists; idempotent)  
-5

which is called etab

Als:  $R - (DUS)$

for each entry in index {  
 if  $etab[entry] \cap \text{undefined symbols} \neq \emptyset$  then  
     pass1(entry)  
     ~~skip this entry~~ -3

{ Find first symbol in dlist in  $(R - (DUS))$  and add it + all ~~remaining~~ following symbols to  $R$ .

17

4. Typically, each entry in a library is compiled separately and then inserted into the library. Thus, for example, a collection of functions is usually kept one function per source file, and each one is compiled and inserted into the library individually. A more convenient approach is to have a compiler option that, given a source file containing only functions (no global data), causes the compiler to generate a library whose entries are the link object code for those functions. Describe what changes to the code generator (*cgen.c*) and to *link* are necessary to implement this strategy. Give enough detail so that a person familiar with *cgen.c* and link object code would understand how to make the necessary changes.

Overview: Have the compiler generate ~~code~~ code for each function into a separate file and then <sup>concatenate</sup> merge the files into a library.

More specifically,

Add a data structure <sup>ctab</sup> that maintains information about the routines called by a particular routine.

~~proceeds~~ <sup>proceeds</sup> causes an entry to be made in *ctab* for the current procedure and establish a global ptr to the entry. When walk encounters a CALL node, it adds (but not redundantly) the ~~routine~~ routine being called to the current *ctab* list. When *proceed* is called, it outputs the index entry by scanning *ctab* + traversing the ~~call~~ call lists. Then, *ctab* is scanned again <sup>+</sup> each temp file is appended to the library.

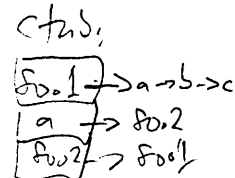
Note that it is a simple matter to produce a top. sorted library.

BTW - No link mods at all. (!)

```

So.1()
{
  a()
  b()
  c();
}
a()
{
  So.2()
}
So.2()
{
  So.1();
}

```



unnecessary for index.

details you forgot:  
1) links -> 3  
2) next location -> 3

19

1. (a) Segments alone cannot be used for labeled common because `link` concatenates segments of the same name from different files and relocates references to segments. In the example given, if `a` is made a segment, linking the two files would result in `a` being two words in length and `y` being incorrectly defined as `'a 1 +'`. A common block is more like a segment that may be defined in many files but is 'allocated' only once. That is, subsequent occurrences of a common block should simply refer to the first one.

(b) Clearly, `link` itself cannot handle common. It is, however, possible to establish compiler and loader conventions so that the compiler can produce segments that are handled as usual by `link` but are recognized by the loader as common blocks and treated accordingly. Perhaps the simplest scheme is to represent a common block `x` as a segment named `comm_x_n`, where `n` is the length of the common block, and to specify `0` on the associated `len` command. Specifying a length of `0` avoids relocation. The loader can obtain the correct length from the segment name. In example in the problem, the output from the common declaration in `A` would be

```
.seg comm_a_1
.def x comm_a_1 0 +
.len comm_a_1 0
```

and the output from `B` would be

```
.seg comm_a_1
.def y comm_a_1 0 +
.len comm_a_1 0
```

could pass length via  
unrecog command, but  
wouldn't know where length  
will come out

Another, similar, approach is to represent common block `x` as a segment named `comm_sub_x`, where `sub` is the subroutine name in which `x` is referenced. In the example in the problem, subroutine `A` would define a segment named `comm_A_a` and `B` would define a segment named `comm_B_a`. Using different names keeps `link` from concatenating the segments. The loader, however, can treat segment names of the form `comm_sub_x` properly by removing the `comm_sub` prefix. Thus, `comm_A_a` and `comm_B_a` will be assigned the same base address. Note that their lengths could also be checked (they should be equal).

2. The solution to this problem is an application of 'link-time binding' described on page 361 of the `link` paper. A copy of `table.o` can be linked with `command.o` and with `lex.o`, providing the symbols `initialize`, `lookup`, and `install` are suppressed, i.e. defined with the `-s` option. The resulting two object files can then be linked to form `prog.out`. Since we cannot alter the object files, the command line option

```
-s symbol,symbol,...
```

which sets the `SUPPRESS` bit for the given symbols is added. Then the object files for `command.o` and `table.o` and for `lex.o` and `table.o` are constructed as by the commands.

```
link command.o table.o -s initialize,lookup,install -o command.out
link lex.o table.o -s initialize,lookup,install -o lex.out
```

The final object file is constructed by linking `command.out` and `lex.out`:

```
link lex.out command.out -o prog.out
```

3. The index lines `'- name1 name2'` can be used to augment the entry structure with a list of other entries that are needed if the entry of interest is linked. The relevant C declarations are

```
struct entry {
    char *e_name;
    int e_needed;
    struct enode *e_family;
    struct entry *e_link;
}
```

If `name1` must be loaded, then `name2` must be loaded

```
struct enode {  
    struct entry *f_entry;  
    struct enode *f_link;  
}
```

*e\_family* is the head of a list of pointers to entry structures for the other entries in the same 'family', i.e. those needed if this entry is linked. This list is constructed as index lines of the form given above are read. *lib(file)* can then be written as follows.

```
lib(file)  
char *file;  
{  
    struct symbol *p;  
    struct entry *q;  
    struct enode *f;  
  
    read index and build entry structure;  
    for each (symbol p in (L∩R)-D) {  
        p->s_entry->e_needed++;  
        for (f = p->s_entry->e_family; f; f = f->f_link)  
            if (q = f->f_entry)  
                q->e_needed++;  
    }  
    for each (entry in the library) {  
        q = associated entry structure;  
        if (q->e_needed)  
            pass1(this entry);  
        else  
            skip this entry;  
    }  
}
```

4. The code generator can generate libraries by emitting code for each function and for the index into temporary files and then concatenating the temporary files at the end of compilation. In doing so, the compiler must insure that each temporary file is complete, e.g. that each includes the necessary literals. Specifically, the major modifications to the routines in *cgen.c* are as follows.

*progbeg* opens a temporary file for the index and initializes a list of structures that describe the temporary files.

*procbeg* opens a new temporary and adds a node to the list initialized by *progbeg*. It also writes an index entry for the procedure to the index temporary file.

*procbeg* must insure that the output is complete. To do so, it executes

```
foreach(defconst, CONST, 0);
```

to include the necessary literals, and emits the appropriate *len* commands. It also must *reset* the location counters to 0, since each library entry is a separate object file.

*procbeg* copies the temporary files, beginning with the index, to the standard output, prefacing each one with the appropriate archive header, and deletes them.

*defconst* should set the *p->s\_offset* field to 0 so that future references to the constant in other functions will cause its inclusion in the object file. (See *cgen.c*). *cgen.c* [1055,21] includes the 60 or so lines needed to implement this feature, which is activated by the *-l* option to the compiler.

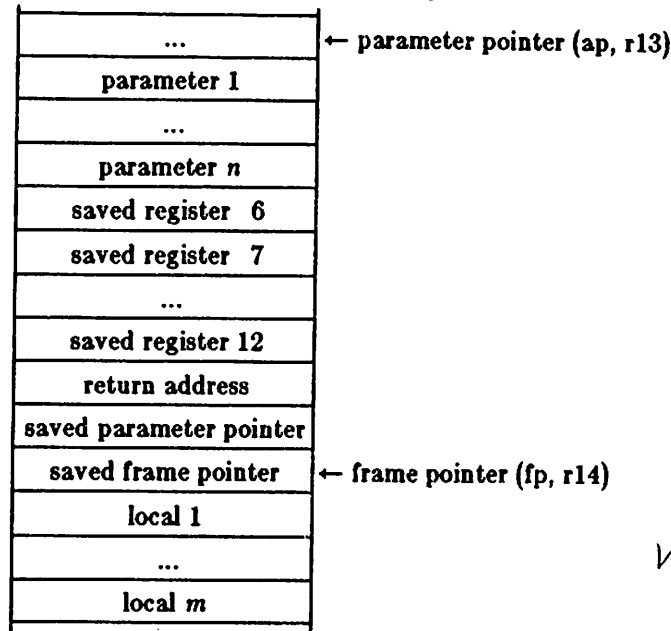
~~Global data has one entry~~  
Library entry for each global data item  
Var name lengths

Name: Bill Mitchell

There are four problems on this exam. Each problem is worth 25 points. This is *closed book* exam.  
Please attempt to solve all problems in the space provided. Show your work in arriving at solutions; partial credit will be given. There is an extra page at the back of the exam. Be brief.  
Good luck.

**Useful Information**

The layout of an activation record for a procedure with  $n$  parameters and  $m$  locals is as follows.



25  
25  
25  
5  
—  
80

*nick work!*

Parameters and locals are addressed by constant offsets from the 'parameter pointer' (register 13) and the 'frame pointer' (register 14), respectively.

*findframe*( $n$ ) positions a global frame pointer, *cursor* (a *frame* structure), to the  $n$ th frame from the top of the stack. It returns 0 if the positioning fails and 1 otherwise.

Possibly useful definitions are as follows.

```

struct frame {
    struct symbol *f_func;    /* frame variables */
    int f_depth;             /* function */
    int f_nargs;             /* depth from the top */
    int f_ap;                /* number of arguments */
    int f_fp;                /* argument pointer */
    int f_sp;                /* frame (local) pointer */
};
struct breakpt {
    ADDRESS b_addr;         /* breakpoint information */
    int b_inst;             /* breakpoint address */
    int b_jumps[2];        /* instruction at breakpoint */
    struct symbol *b_func;  /* jumps back into main code */
    int b_epoint;          /* function in which breakpoint is set */
};
    
```

```
struct state {                                /* machine execution state */
    ADDRESS x_pc;                             /* location counter */
    int x_regs[16];                           /* registers */
    int x_level;                              /* call level */
    int x_trace;                              /* trace counter */
};
struct symbol {                               /* symbol table entries */
    char *s_name;                             /* entry name */
    int s_flags;                              /* flags (see below) */
    int s_type;                               /* data type (see below) */
    int s_size;                              /* size of object in bits */
    int s_offset;                            /* activation frame offset in words */
    struct symbol *s_link;                   /* link to next table entry */
#ifdef CDB
    int *s_points[1];                        /* execution points for functions */
#endif
};
struct node {                                 /* expression nodes */
    int e_op;                                 /* operator */
    int e_type;                              /* type of result */
    struct node *e_left;                    /* left operand */
    struct node *e_right;                  /* right operand */
    int e_result;                           /* execution result */
    int e_count;                            /* reference count */
};

/* symbol table entry flags */
#define DEF    0001                          /* symbol is defined */
#define REF    0002                          /* symbol is referenced */
#define GLOBAL 0004                          /* symbol is a global */
#define LOCAL  0010                          /* symbol is a local */
#define PARAM  0020                          /* symbol is a formal parameter */
#define CONST  0040                          /* symbol is a constant */
#define KEYWORD 0100                         /* symbol is a keyword */
#define ARRAY  0200                          /* symbol is an array */
```

1. Suppose the breakpoint command was extended to be

`b name n { cdb commands... }`

where the braces enclose a sequence of `cdb` commands, separated by semicolons, to be executed when the breakpoint is encountered. If the `{...}` specification, is omitted, a `c` command is executed and commands are accepted from the user, as usual. For example,

`b strcmp 3 { -;C;j }`

would set a breakpoint at the 4th execution point in `strcmp` and execute the `-`, `C`, and `j` commands upon encountering the breakpoint. Describe what changes or additions you would make to `cdb` to implement this extension. Be specific.

Add something to the `breakpt` structure that allows a string to be associated w/ a breakpoint. When `{-}` is specified, `b_cmd` is pointed at a copy of the string `sum` the brackets, otherwise, `b_cmd = NULL`.

Next, restructure the command loop of `db` so that the user can type strings like `-jC;j` at the prompt. This is fairly trivial,

just split the line at `j`'s (watch out for quoted strings!) and make a queue of those lines.

Input lines are added to the queue and `cdb` takes lines from the queue (as if the user had typed them in) and executes them.

With the new command parsing in hand, modify the `j` case (after the resume) to check the `b_cmd` field, if `NULL`, just do a `c`; otherwise, point the input line pointer at the `b_cmd` string and jump to the point just after where commands are read in. Thus, it looks like the user typed in the commands associated w/ the breakpoint.

i.e.  
char \*b\_cmd  
↑  
allocate

of the  
breakpoint  
hit

✓/

2. It is useful to be able to call user-defined C functions during the evaluation of expressions typed to cdb. For example, typing

```
p lookup("int", &root)
```

f(x),

at a breakpoint in main in test6.c would print the index of the entry for int. During execution of a user-defined function in response to a cdb expression, breakpoints are not set. Describe the changes you would make to cdb to implement this feature. Assume user-defined functions must be executed at full speed, i.e. you cannot interpret the code. Be specific about the modules in which changes are necessary and the details concerning stack usage.

Add two global variables;

```
int argno;
int fargs[NPARAM];
```

↙ max value for # of parameters, the 10 seems to limit this to about 2.

In eval:

```
argno = 0;
```

In expcode(op, r1, r2, p)

Case O\_ARG:

```
fargs[argno++] = results[r1]; // value;
break;
```

Case O\_CALL;

```
int (*f)()
```

f = <sup>(int)</sup> p → f-func → s-offset;

```
rval = (*f)(fargs[0], fargs[1], ... fargs[NPARAM-1])
```

```
results[r] = rval
```

allocated

```
break;
```

↖ increment argno.

Problems; runs on cdb stack; could do switch, but this seems like a reasonable solution. Also, could craft stack by hand to avoid (fargs[0] ...), but again, has bad win/lose ratio. Also limited to # of fargs, but since the 10 doesn't seem to like more than 7, that doesn't seem like a problem (Hand crafting stack avoids max fargs problem.) ~~Over.~~ Over.

25/30



This doesn't allow for  $f(x(1), y(2))$  ← it's very close, however.

but I can't figure it out. Perhaps something ~~of~~ w/ allocating parameter values in the results array and maintaining a list of ~~of~~ the elements w/ the values and then ~~of~~ filling in arrays at the point of call.

3. Suppose `edb` had a `k` command, which caused the next expression to be executed as if a breakpoint was set at the next execution point. Successive `k` commands could then be used to 'single step' the program.

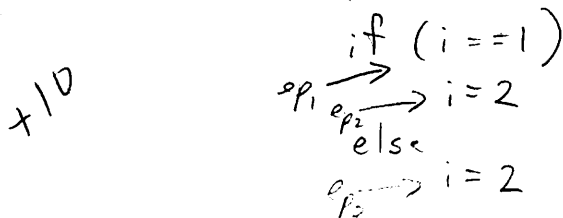
(a) (10 points) Consider the following implementation of the `k` command.

"If execution is currently halted at execution point  $n$  in function  $f$ , set a breakpoint at execution point  $n+1$ . Add a `b_flags` field to the breakpoint structure and set a flag indicating the breakpoint at  $n+1$  is a 'single-step' breakpoint. The `k` command is then identical to the `j` command, except that 'single-step' breakpoints are removed after they are encountered."

This implementation does *not* work; explain why.

⊗ Execution points are non-sequentially reached, because the flow of control is not necessarily sequential.

Consider:



If  $i == 0$ , the flow is 1, 3;  $ep_2$  is never reached. (Hence, a `k` at  $ep_1$  might go for a good way.)

(b) (15 points) Ignoring the complications implied by the `return` statement, describe what modifications to the implementation outlined in part (a) are necessary to correctly implement the `k` command. Be specific about any changes you make to the compiler's output or to `edb`. Do not describe any changes to the compiler itself; concentrate on the changes in its output.

The compiler needs to tell the debugger about execution points ( $ep_i$ ) where the code might cause  $ep_{i+1}$  to not be reached "next", and for such  $ep_i$ , which  $ep_j$  might be next. For example, in an `if` stmt, the first  $ep_i$  in either arm might follow the clause. In a `while` stmt, the body might be skipped. The `for` stmt is also a problem, and `break` + `continue` are also trouble.

So, modify the compiler so that for each  $ep_i$  it outputs the address of the first instruction and follows the address w/ negative addresses of instructions that begin execution points that might be reached next.

↓ Over

25

inst. address

For example:

ep<sub>0</sub>, 000 → i = f();

ep<sub>1</sub>, 010 → if (i == 1)

ep<sub>2</sub>, 014 → i = 2;

else

ep<sub>3</sub>, 017 → i = 3;

ep<sub>4</sub>, 021 j = i \* 2

ep array:

ep <sub>0</sub> →	1
ep <sub>1</sub> →	010
	-014
	-017
ep <sub>2</sub> →	014
	-021
ep <sub>3</sub> →	017
ep <sub>4</sub> →	021

No neg #'s if next ep is sequential.

The k command ~~will~~ finds the current ep + determines addresses of all the possible next eps. It makes a list of them, sets ~~to~~ breakpoints at all ~~of~~ of them. ~~It~~ It does a switch, a breakpoint is hit, it cleans up the points using the list + returns to the cdb command level.

4. C permits 'register' declarations, e.g.

```
register int n; only, looks
```

*sym table*  
*output only*

which instruct the compiler to use machine registers (instead of allocating space in the frame) for the indicated local variable. Suppose our compiler used registers 6-12 as necessary for such variables, and set the REGISTER flag in the symbol table entry field, *s\_flags*, to indicate that *s\_offset* contained a register number instead of a frame offset. Describe the steps necessary to compute the *lvalue* of a register variable during the evaluation of an expression, i.e. the steps in the O\_ID case in *exprcode*. Be careful!

O\_ID;

if (ISREGISTER(p))

lvalue = &(up->x\_regs[p->s\_offset]);

*ptr to user state*

*ptr to symbol*

This seems to be enough. If the value is required, it is gotten from the saved register value. If it is modified, the saved value is written over and subsequently restored by resume.

*for frame 0,  
not for  
frame i, i > 0  
-20*

*15/*

1. Breakpoint commands can be handled by adding a *b\_cmd* field to the breakpoint structure that points to the command string given in braces. This field is set on a *b* command, and, upon encountering a breakpoint, its contents are appended to an input buffer, substituting newlines for semicolons. The main command loop is modified to take commands from the command buffer until it is empty, and then resume accepting commands from the terminal. This approach preserves the current method for handling errors and permits the breakpoint commands to set and remove breakpoints correctly.
2. User-defined functions can be called by simply calling them from *cdb*, i.e. using *cdb*'s stack. This is possible because the locations of the arguments and locals, given by the frame offsets, are independent of which stack is used, providing the arguments are pushed onto the right stack. The values of the actual arguments can be saved in an array during evaluation in *exprcode*:

```
union words args[MAXARGS];
int argptr = 0;
...
case O_ARG:
    args[argptr++] = wlp->w_int;
    break;
```

Calling the function amounts to passing the correct number of arguments and getting the address of the function:

```
union word *ap;
int (*f)();
...
case O_CALL:
    ap = &args[argptr-r2]; /* r2 = argument count */
    f = p->e_left->e_left->s_offset;
    switch (r2) {
        case 0: wp->w_int = (*f)(); break;
        case 1: wp->w_int = (*f)(*ap); break;
        case 2: wp->w_int = (*f)(*ap, *(ap+1)); break;
        ...
        case MAXARGS:
            wp->w_int = (*f)(*ap, *(ap+1), ..., *(ap+MAXARGS-1));
            break;
    }
    argptr -= r2;
    break;
```

This code takes advantage of the architectural property that all scalar types occupy 1 word of storage.

3. (a) The proposed implementation of the *k* command does not work because it assumes that the order in which execution points are encountered during *execution* is the same as during *compilation*, which is incorrect. Consider, for example, the following fragment and its execution points.

```
while ( i < 10)
    ↑0
    f(i++);
    ↑1
g();
↑2
```

Stepping through the loop with the *k* command implemented as described will only step through the first iteration of the loop because the execution point following point 1 is number 2. During *execution*, however, control flows to the execution points in the sequence 0, 1, 0, 1, ..., 0, 2.

- (b) The flaw in the implementation strategy is that the notion of "successor" execution point is incorrect. This can be corrected by having the compiler generate the list of execution points and their *execution-time* successors. For example, the *e\_points* array could be organized so that it contains the address of the

corresponding execution point followed by a list of the addresses of its successor execution points, terminated by -1. (Currently, this list will have 0, 1, or 2 entries). The entire list is terminated by a 0 execution point address, as before. For the code fragment given above, the output for the *e\_points* array might be as follows.

```
code 056 +
code 065 +
code 075 +
-1
code 065 +
code 056 +
-1
code 074 +
-1
0
```

As shown, execution point 0 is at code 056 + and its execution-time successors are points 1 and 2, which are at code 065 + and code 075 +, respectively. Likewise, point 1 is at code 065 + and its successor is point 0, which is at code 056 +, and point 2 is at code 074 + and it has no successors. Given this output, the strategy described in part (a) can be generalized to set 'single-step' breakpoints at all of the *successors* of the current breakpoint. This approach can also be extended to handle **return** statements and flowing off the end of a function.

4. The basic question of this problem is "where is register *r* for frame *i*", where frames are numbered beginning at 0 from the top of the stack. The answer depends on the value of *i*. If *i* is 0, register *r* is stored in the *x\_regs* fields of the user's program state. Because registers are saved by the callee during the function entry sequence, if *i* is non-zero, register *r* is stored in the frame for the function called by the function associated with frame *i*, i.e. register *r* is stored at offset *fp-15+r* in frame *i-1*. The appropriate code for the *O\_ID* case in *exprcode* is as follows.

```
case O_ID:
    q = (struct symbol *) p->e_left;
    if (q->s_flags&GLOBAL)
        wp->w_addr = q->s_offset;
    else if (q->s_flags&REGISTER) {
        if ((i = cfrp->f_depth) == 0)
            wp->w_addr = &up->x_regs[q->s_offset];
        else {
            findframe(i - 1);
            wp->w_addr = cursor.f_fp - 15 + q->s_offset;
            findframe(i);
        }
    }
    ...
    break;
```

C Parser

Design and implement a recursive descent parser for the subset of C defined by the grammar given below and the attached syntax diagrams. Divide your parser into modules defined by the following files.

- name* *purpose*
- cc.h* definitions
- ezpr.c* expression parsing - 530
- egen.c* symbolic code generation
- clerc.c* lexical analysis
- cmain.c* main program
- cpars.c* declaration parsing - 360
- csat.c* statement parsing - 180
- csym.c* symbol table management

Both the source and object files (.rel files) are available on [1055,21] for those files indicated by an asterisk. You are free to copy and examine the source code as desired, but you should load your program with the object files.

Syntax

Your parser must accept a subset of C defined by the following extended BNF grammar in which alternatives are separated by vertical bars, i.e. *alb* stands for *a* or *b*, optionally by square brackets, i.e. [*a*] stands for *a*, repetition is indicated by curly braces, i.e. {*a*} stands for *a*{*a*}{*a*}{...}, and parentheses are used for grouping, i.e. (*alb*)*c* stands for *alb**c*. Identifiers for which there is no production and operators denote terminal symbols. When braces, brackets, vertical bars, or parentheses are used as terminal symbols, they are enclosed in quotes to distinguish such usage from their grammatical use.

```

prog : { decl ; | func }
decl : type decl { , decl }
decl : ptr id [ '[' [ con ] ']' | '(' ')' ]
type : char | float | int
ptr : { * } - do stms # $ for use in symbol table
func : [ type ptr ] id '(' [ id { , id } ] ')' { decl ; } { stmt ; }
stmt : expr ;
      | if '(' expr ')' stmt [ else stmt ]
      | while '(' expr ')' stmt
      | for '(' [ expr ] : [ expr ] ')' stmt
      | return [ expr ] ;
      | break ;
      | continue ;
      | '{' stmt { '}'
      | ;
expr : expr [ binop ] = expr
      | expr ( relop | binop ) expr
      | term
  
```

*Assume untyped fens*

```

term : term ( ++ | -- )
      | ( ! | * | & | ~ | - | + | ++ | -- ) term
      | term '[' expr ']'
      | term '(' [ expr { , expr } ] ')'
      | id
      | con
      | fcon
      | ccon
      | scon
      | '(' expr ')'
  
```

```

relop : == | != | < | <= | > | >= | < | > | && | ||
  
```

```

binop : ~ | & | ! | ? | < | > | < | > | + | - | * | / | %
  
```

Syntax diagrams corresponding to each of the productions in the above grammar are attached and should be used to guide the implementation of your parser. The corresponding *first* and *follow* sets for the above grammar are as follows.

```

first(prog) = first(decl) ∪ first(func) ∪ { ε EOF }
            = { char int float id EOF }

first(decl) = first(type)
            = { char int float }

first(ptr) = { * ε }

first(func) = first(type) ∪ { id }
            = { char int float id }

first(stmt) = first(expr) ∪ { if while for return break continue '{' ; }
            = { ! * & ~ ++ -- id con fcon ccon scon }
            return break continue '{' ; }

first(expr) = first(term)
            = { ! * & ~ ++ -- id con fcon ccon scon }

first(term) = { ! * & ~ ++ -- id con fcon ccon scon }

follow(prog) = { EOF }

follow(decl) = { ; }

follow(decl) = { , } ∪ follow(decl)

follow(type) = first(decl) ∪ first(ptr) - { ε } ∪ { id }

follow(ptr) = { id }
  
```

```

follow(func) = first(prog) ∪ follow(prog)
              = { char int float id EOF }
follow(stmt) = first(stmt) ∪ { } else {
              = { ! * & ~ ~ + + - - id con fcon ccon soon ( if while for
                return break continue '{ ; ' } else }
follow(expr) = { ; } = binop relop ] . {
follow(term) = follow(expr) ∪ { + + - - [ ( {
              = { ; } = binop relop ] . + + - - [ ( {
    
```

### Code Generation

*cgem1.c* contains a symbolic code generator that must be loaded with your program to obtain output. The code generation interface consists of a set of functions that are called at appropriate points during parsing as follows.

*z = progbeg(); progend(x)* should be called before any parsing is initiated, and *progend* should be called after all parsing is completed, passing it the integer returned by *progbeg*.

*deflabel(tab)* should be called at the point of definition of label number *tab*. Label numbers should be generated sequentially by the parser beginning at 1.

*global(p)* should be called with the symbol table pointer for each global identifier, except functions, as they are declared.

*defconst(p)* After parsing all of the input *defconst* should be called with the symbol table pointer for each constant referenced in the input.

*z = procbeg(p); procend(x)* *procbeg* should be called when parsing for a procedure body is initiated with the symbol table pointer for the procedure. *procend* should be called after parsing a procedure body, but before the local variables are expunged from the symbol table, passing it the integer returned by *procbeg*.

*jump(tab)* should be called to generate a jump to label number *tab*.

*r = exprcode(op, r1, r2, tab, p)* should be called for each node in an expression tree in the correct evaluation order. *op* is the operator (see *cc.h*), *r1* and *r2* are the results from calling *exprcode* for the descendants of the node, *tab* is the associated label (used only for conditional operators), and *p* is a pointer to the node. *exprcode* should return an integer describing the result of executing the operator, such as a temporary number. For the operator *O\_CALL*, *r2* should be the number of arguments. (See *cgem1.c*).

*retcode(r)* should be called to generate code for a *return* statement. It is passed the result of *exprcode* for the root of the expression tree for the return expression, or 0 if the return expression is omitted.

### Semantics

Your parser must implement the correct semantics for C. The following code templates define the semantics of statements. In the following, *L* stands for a label number. For loops, *L* labels the beginning of the loop body, *L+1* is the 'continue' label, and *L+2* is the 'break' label.

```

if ( expr ) stmt1 else stmt2
    if ( ~expr ) goto L
    stmt1
    goto L+1
    stmt2
    L
    L+1
    
```

```

while ( expr ) stmt
    L
    L+1
    if ( ~expr ) goto L+2
    stmt
    goto L
    L+2
    
```

```

for ( expr1 ; expr2 ; expr3 ) stmt
    expr1
    if ( ~expr2 ) goto L+2
    stmt
    expr3
    goto L
    L
    L+1
    L+2
    
```

```

break
    goto L+2(L is current loop label)
    
```

```

continue
    goto L+1(L is current loop label)
    
```

Type checking and the order of evaluation are the important semantic issues for expressions. The following tables summarize the type behavior of the operators. In all cases, characters are 'widened' to integers, if necessary. Consult the *C Reference Manual* at the back of *The C Programming Language* for further details.

expression	equivalent expression
<i>e1[e2]</i>	<i>*(e1 + e2)</i>
<i>e1 op e2</i>	<i>e1 = e1 op e2</i> ( <i>e1</i> evaluated once)
<i>++e</i>	<i>e ++ = 1</i>
<i>--e</i>	<i>e -- = 1</i>
<i>e++</i>	<i>(t = e, e ++ = 1, t)</i> ( <i>e</i> evaluated once)
<i>e--</i>	<i>(t = e, e -- = 1, t)</i> ( <i>e</i> evaluated once)
<i>e</i>	<i>e != 0</i> ( <i>e</i> used in conditional)



expression	result
!conditional	conditional
*(ptr T)	lvalue T
&(lvalue T)	ptr T
-int	int
-float	float
~int	int
(ptr T) ( +   - ) int	ptr T
T(...)	T
int ( +   -   *   / ) int	int
(int   float) ( +   -   *   / ) (int   float)	float
int % int	int
(int   float) ( ==   !=   < =   > =   <   > ) (int   float)	conditional
(ptr T) ( ==   != ) int	conditional
(ptr T) ( ==   !=   < =   > =   <   > ) (ptr T)	conditional
int ( <<   >>   &   '   ~ ) int	int
conditional ( &&      ) conditional	conditional
lvalue T = T	T
lvalue int = float	int
lvalue float = int	float
lvalue (ptr T) = int	ptr T

Evaluating expressions in the correct order amounts to making the correct sequence of calls to *exprcode*. For binary operators, including assignment, the left operand must be evaluated first. Evaluating conditional operators ( ! && || == etc.) is similar except that evaluation causes transfer of control instead of computing a value. In addition, 'short circuit' evaluation is used: the expression is evaluated only as far as is necessary to determine the truth value. Conditional expressions are always evaluated with a 'true' or 'false' label (but not both) to which control should be transferred. Assume *walk(p, flab, flab)* traverses the tree rooted at *p* with true and false labels *flab* and *flab*, respectively. The following table shows the correct sequence of recursive traversals depending on *p*'s operator, *flab*, and *flab*; *left* and *right* denote the left and right operands of *p*, respectively, and *L* denotes a generated label.

op, flab, flab	sequence
, lab, 0	walk(left, lab, 0) walk(right, lab, 0)
, 0, lab	walk(left, L, 0) walk(right, 0, lab)
&&, lab, 0	L: walk(left, 0, L) walk(right, lab, 0)
&&, 0, lab	L: walk(left, 0, lab) walk(right, 0, lab)
<, lab, 0	r1 = walk(left, 0, 0) r2 = walk(right, 0, 0) exprcode(<, r1, r2, lab, p)
<, 0, lab	r1 = walk(left, 0, 0)

	r2 = walk(right, 0, 0) exprcode(> =, r1, r2, lab, p)
!, flab, flab	walk(left, flab, flab)
binaryop, 0, 0	r1 = walk(left, 0, 0) r2 = walk(right, 0, 0) r = exprcode(binaryop, r1, r2, 0, p)
unaryop, 0, 0	r1 = walk(left, 0, 0) r = exprcode(unaryop, r1, 0, 0, p)

As shown, evaluation of the comparison operators requires both the comparisons and their complements.

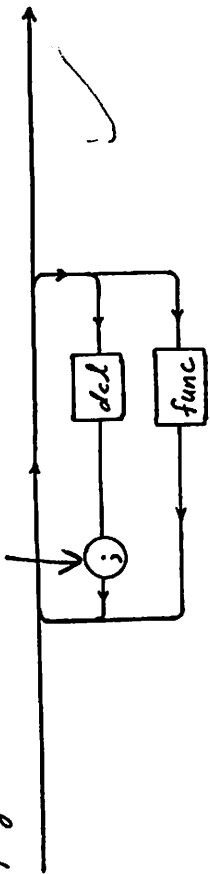
*c0.exe[1055,21]* is a working version of the parser. In addition, the files *test?.c[1055,21]* contain test programs of increasing complexity; using these files as input, your parser should produce output that is *identical* to that produced by *c0*.

This is a moderately large program; the following table gives the line counts for my version of each of modules in this program.

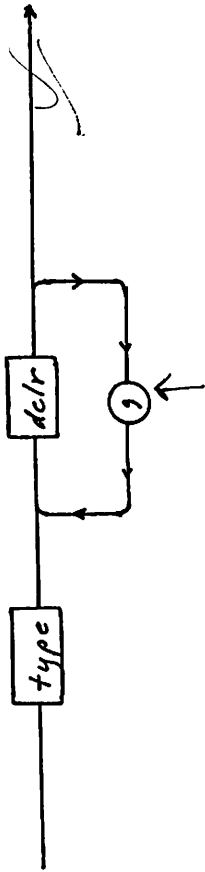
521	<i>cepr.c</i>
137	<i>cgen1.c</i>
251	<i>clcx.c</i>
32	<i>cmain.c</i>
360	<i>cpars.c</i>
173	<i>cstat.c</i>
189	<i>csym.c</i>
1663	<i>total</i>

For further information on recursive descent parsing, see Chap. 5 of N. Wirth, *Algorithms + Data = Programs*, Prentice-Hall, 1976. Consult B. W. Kernighan and D. M. Ritchie, *The C Programming Language*, Prentice-Hall, 1978 for details concerning C. Further information (alot of it!) will be given in class; don't miss it.

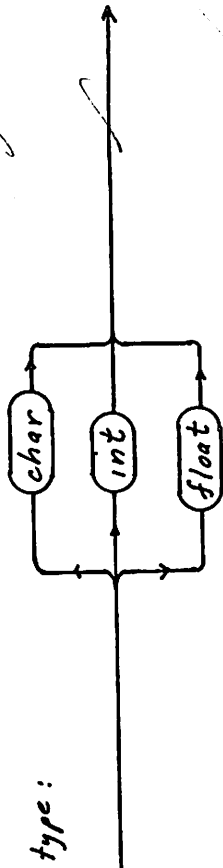
prog:



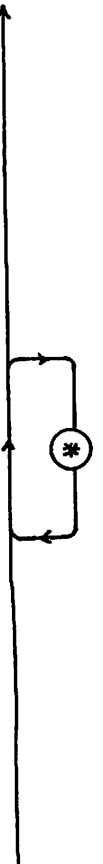
decl:



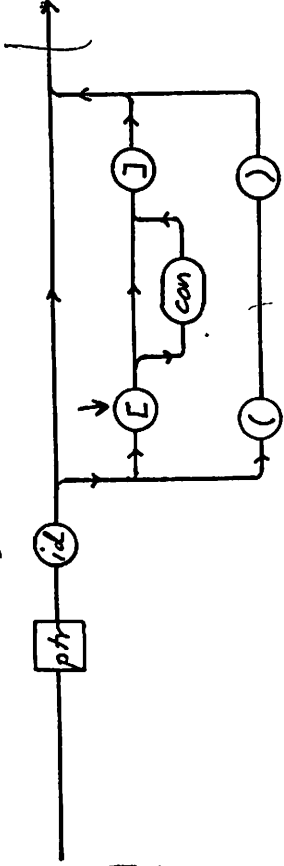
type:



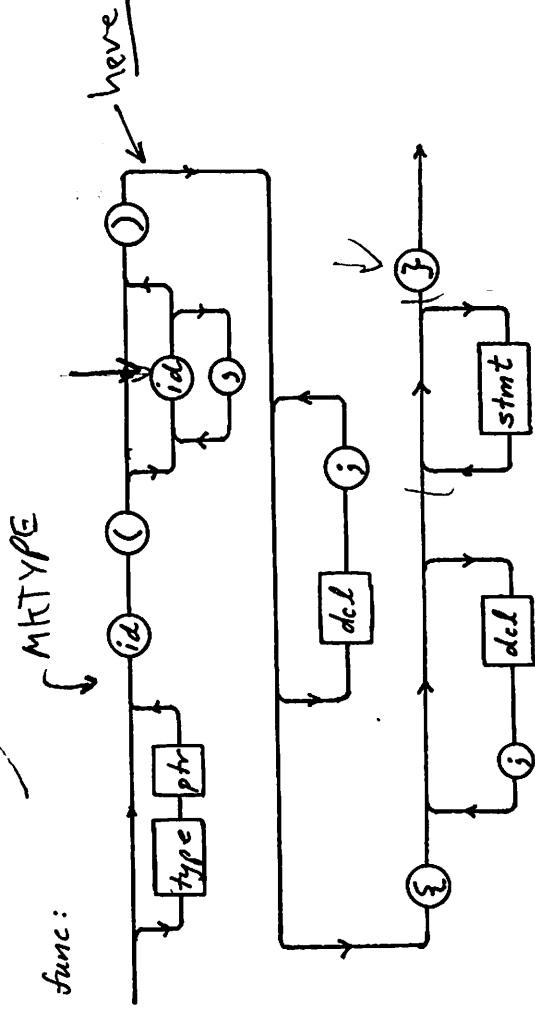
ptr:



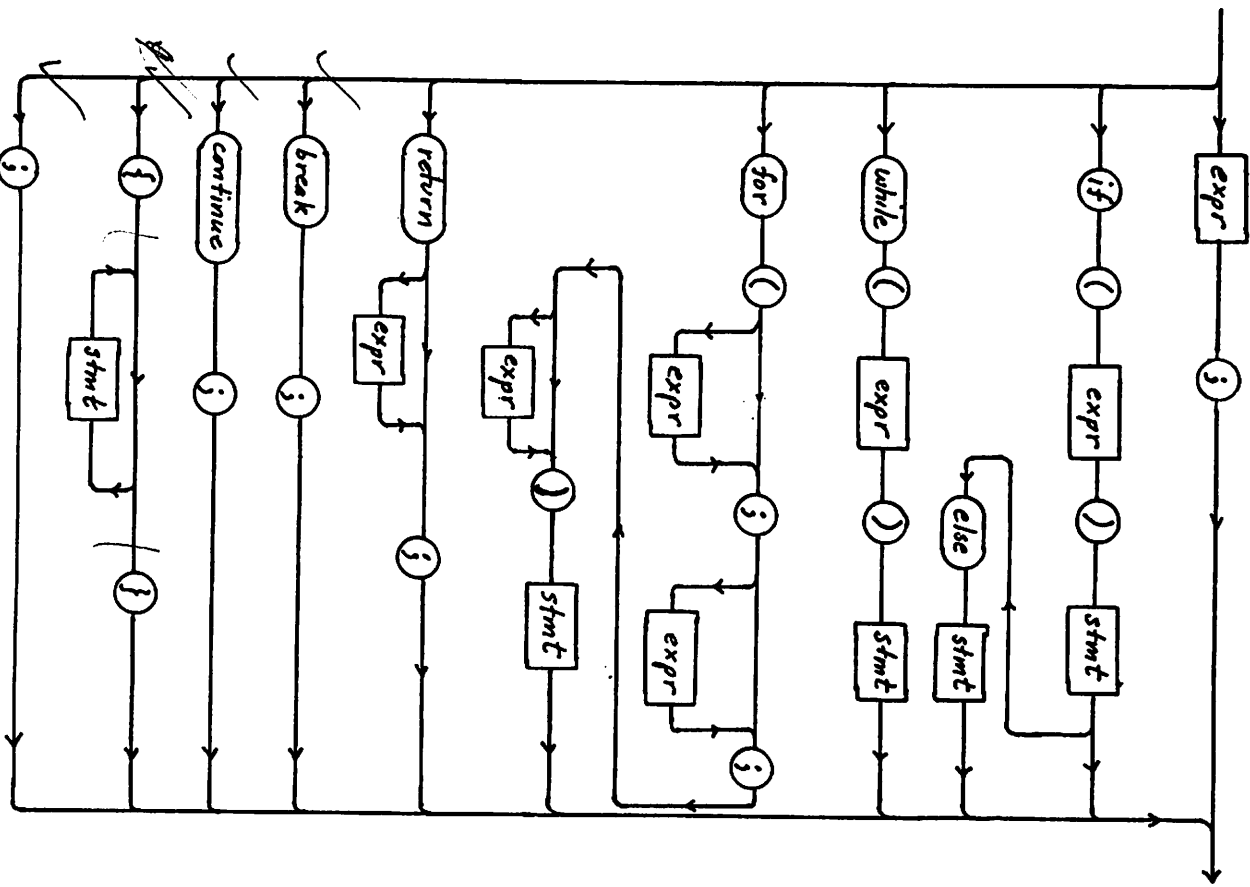
de/r:



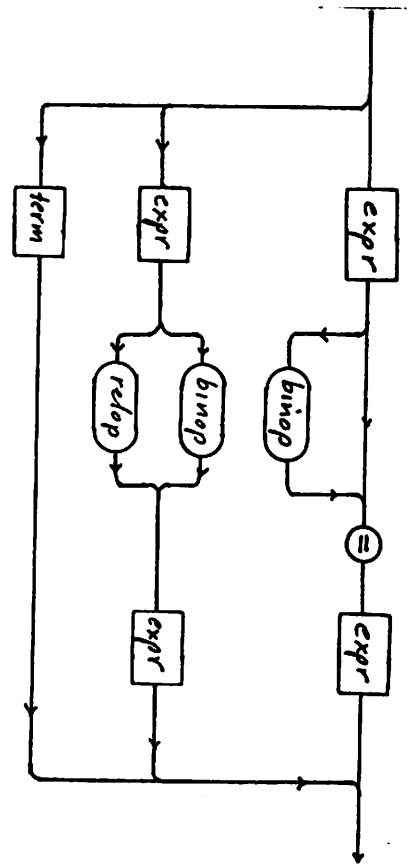
func:

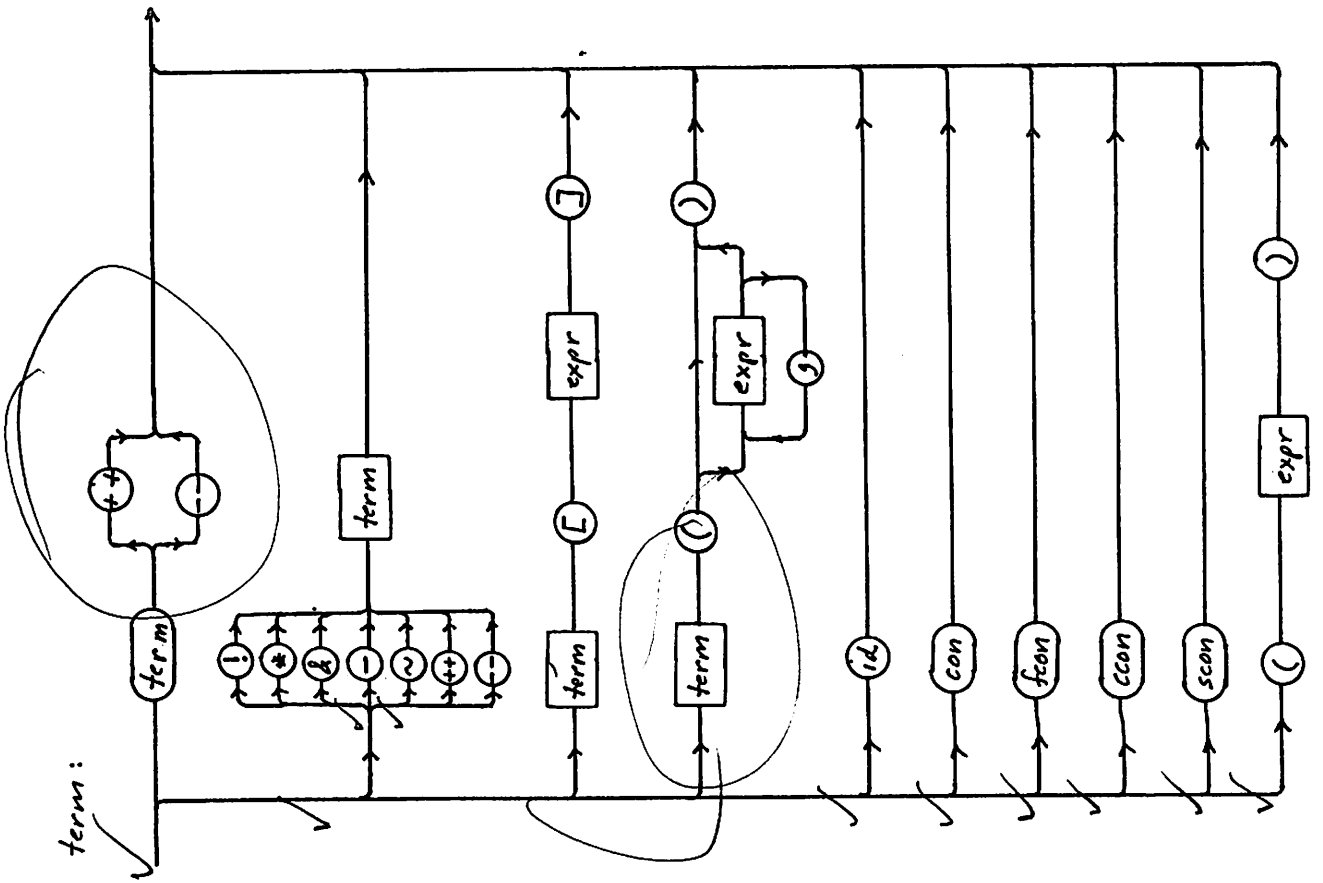


stmt:



expr:





### C Code Generation

The object of this assignment is to replace the routines in *cgen1.c* by routines that generate DEC-10 instructions in **link** object code format. Templates of the expected output for each of the routines and for expressions are given below. As before, the relevant files are as follows.

<i>name</i>	<i>purpose</i>
* <i>cc.h</i>	definitions
* <i>opcode.h</i>	instruction definitions
<i>cexpr.c</i>	expression parsing
<i>cgen.c</i>	DEC-10 code generation
* <i>clax.c</i>	lexical analysis
* <i>cmain.c</i>	main program
<i>cparse.c</i>	declaration parsing
<i>cstat.c</i>	statement parsing
* <i>csym.c</i>	symbol table management

23,5,6

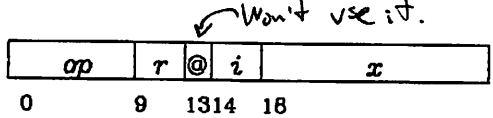
275 can edit cgen1.c

Both the source and object files (*.rel* files) are available on [1055,21] for those files indicated by an asterisk. Use your versions of *cparse.c*, *cstat.c*, and *cexpr.c* and place your code generator in *cgen.c*.

The general assembly language form of DEC-10 instructions is

*op* *r*,@*x*(*i*)

where *op* is a 9-bit instruction code, *r* is a 4-bit register number, @ is a 1-bit indirection indicator, *x* is an 18-bit address expression or small constant, and *i* is a 4-bit index register number. Any of *r*, @, *x*, or *i*, can be omitted. These fields appear in the 36-bit word as follows.



Object code consists of *commands*, which are introduced by a period in column 1, and *object text*, which contains symbolic 36-bit polish suffix expressions whose values occupy some memory cell. Object code is defined by the following grammar.

*file* : { *cmd* | *expr* } names segmented

*cmd* : .def [ -s ] *id* *expr*  
 | .seg *id*  
 | .len *id* *con*

*expr* : *expr* *expr* ( + | - | < | > ) | *id* | *con*

Code is emitted into *segments*, which are blocks of object text that are loaded in contiguous blocks of memory. There are three segments: **code**, which contains executable program code, **data**, which contains global data, and **lit**, which contains constants (e.g. program literals). **seg** commands are used to start or to resume emitting code into the named segment. **def** commands are used to define symbols, e.g.

```
.def rows data 56 +
.def -s $4 code 100 +
```

define the address of global variable **rows** to be the 57th word in the **data** segment and the address of label **\$4** to be the 101st word in the **code** segment. The **-s** indicates that **\$4** is a temporary symbol. **len** commands give the length, in words, of each segment.

The DEC-10 has 16 registers, which are to be used as follows.

Size fields in bits, but allocated words

register	use
0	not used
1	function return values
2-5	scratch, used in entry/exit sequence
6-12	scratch, saved across calls
13	parameter pointer
14	frame pointer
15	stack pointer

Expressions should be evaluated using registers 6-12.

Global variables are addressed by simply referring to their names, e.g. loading the address of the global variable `rows` into register 6 is accomplished by `movei 6,rows` in assembly language and by

```
0201300000000 rows +
```

in `link` object code. Constants are addressed by generating a label for them and placing the actual constant in the `lit` segment, e.g. loading the constant 45 into register 6 is accomplished by

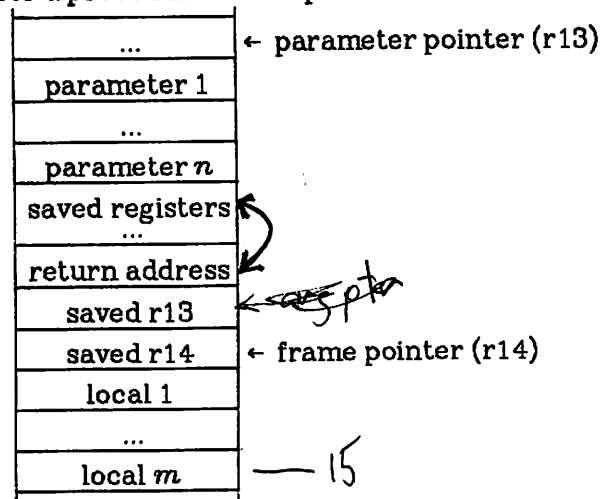
```
move 6,$15
...
$15: exp 45
```

in assembly language and by

```
0200300000000 $15 +
...
.seg lit
.def -s $15 lit 1 +
45
```

in `link` object code.

Addressing local variables and parameters depends on the run-time representation of procedure activation records. Locals and parameters are located in the current activation record or 'frame' for the procedure in which they are declared. The layout of an activation record for a procedure with  $n$  parameters and  $m$  locals is as follows.



Parameters and locals are addressed by constant offsets from registers 13 (the 'parameter pointer') and 14 (the 'frame pointer'), respectively. These offsets should be computed by `procbeg` stored in the `s_offset` fields of the symbol table entries for relevant variables. As shown, the offset for the  $k$ th parameter is  $k$ . For example, if `a` is the third

parameter, loading the address of a into register 6 is accomplished by

**movei 6,3(13)**

in assembly language and by

**020131500003**

in **link** object code. Offsets for locals begin at 1 and increase according to the width of each local.

The output for each routine in the bottom half is as follows. For explanatory purposes, Assembly language output is given in some cases for DEC-10 instructions; **link** object text should be emitted for the actual output, however.

<i>routine</i>	<i>output</i>
✓ defconst(p)	<pre>                     .seg lit                     def \$L lit offset +                     char, int, float constants:                         bit pattern                     char * constants:                     0331100000001 lit offset +                     9-bit bytes, 4/word                 </pre> <p><i>Defining \$L as offset</i>  <i>code on lit</i></p>
✓ deflabel(lab)	<pre>                     .def -s \$lab segment offset +                 </pre> <p><i>code?</i></p>
exprcode(op,r1,r2,lab,p)	see below
✓ global(p)	<pre>                     .def p-&gt; s_name data offset +                 </pre>
✓ jump(lab)	<pre>                     jrst \$lab                 </pre> <p><i>\$</i></p>
✓ procbeg(p)	<pre>                     .seg code                     0331100000000 lit offset +                     .seg lit                     p-&gt; s_name given in 9-bit bytes, 4/word                     .seg code                     .def p-&gt; s_name code offset +                     jsp 4,c.savr                     adjsp 15,size of locals                     (procbeg must also compute s_offset                     values for local variables)                 </pre> <p><i>2, 3, 5, 6</i>  <i>int done if size is 0</i></p>
procend(x)	no output
progbeg()	no output
✓ progend(x)	<pre>                     .len code length                     .len data length                     .len lit length                 </pre>
✓ retcode(r)	<pre>                     move 1,r                     jsp 4,c.retr                 </pre>
✓ rfree(r)	<pre>                     no output                     free resource r                 </pre>

In generating code for expressions, each operator computes its result and leaves it in a register (6..12) for use as an operand by other operators. In the following table, unless otherwise specified, the operands of each operator are given by *r1* and *r2* and the result is given in *r*, and *r1* and *r2* should be freed. Generated code is given for each operator and combination of operand types. Omitted type information indicates the default code for all types not explicitly mentioned.

<i>operator</i>	<i>operand type / scope</i>	<i>code</i>
<del>O_ID</del>	global	movei r,s_name
<del>O_ID</del>	local	movei r,s_offset(14)
<del>O_ID</del>	parameter	movei r,s_offset(13) if <i>id</i> is a char, also do hrli r,0001100 else if <i>p</i> -> <i>e_type</i> is char *, also do hrli r,0331100
<del>O_CON</del>	char	movei r,constant
<del>O_CON</del>	int <= 18 bits	movei r,constant
<del>O_CON</del>		move r,\$L (L in lit segment)
<del>O_CVI</del>	float	fix r,r1
<del>O_CVI</del>	char - no-op	fix r,r1
<del>O_CVF</del>		fitr r,r1
<del>O_NEG</del>		movn r,r1
<del>O_BCOM</del>		setcm r,r1
<del>O_INDIR</del>	lvalue char	ldb r,r1
<del>O_INDIR</del>		move r,(r1)
<del>O_ARG</del>		push 15,r1
O_CALL		movni 5,r2 (r2 = number of args) pushj 15,p->e_left->e_left->s_name move r,1
<del>O_ASGN</del>	lvalue char, char	dpb r2,r1
<del>O_ASGN</del>		movem r2,(r1) (result in r2)
O_LT	char *, char *	push 15,r1 push 15,r2 pushj 15,c.scop adjsp 15,-2 jumpl 0,lab
O_LT		camge r1,r2 jrst lab

*fix r,r1*  
*only have result ptrs to lvalues, up to rvalues would mess this up*



✓ O_ADD	char *, int	adjbp r2,r1 (result in r2)
✓ O_ADD	int, int	add r1,r2 (result in r1)
✓ O_ADD	float, float	fadr r1,r2 (result in r1)
✓ O_SUB	char *, int	movn r2,r2 adjbp r2,r1 (result in r2)
✓ O_SUB	int, int	sub r1,r2 (result in r1)
✓ O_SUB	float, float	fsbr r1,r2 (result in r1)
✓ O_MUL	int, int	imul r1,r2 (result in r1)
✓ O_MUL	float, float	fmpr r1,r2 (result in r1)
✓ O_DIV	int, int	idiv r1,r2 (r1+1 must be free) (result in r1)
✓ O_DIV	float, float	fdvr r1,r2 (result in r1)
✓ O_MOD	int, int	idiv r1,r2 (r1+1 must be free) (result in r1+1, free r1)
✓ O_BXOR		xor r1,r2 (result in r1)
✓ O_BOR		ior r1,r2 (result in r1)
✓ O_BAND		and r1,r2 (result in r1)
✓ O_LSH		lsh r1,(r2) (result in r1)
✓ O_RSH		movn r2,r2 lsh r1,(r2) (result in r1)

15  
 9-GE → cam!  
 17 13

The remaining comparison instructions are as follows.

operator	jump?	cam?
O_LT	jumpl	camge
O_LE	jumple	camg
O_EQ	jumpe	camn
O_NE	jumpn	came
O_GE	jumpge	caml
O_GT	jumpg	camle

c1.exe[1055,21] is a working version of the compiler. When run with a -d option, c1

prints instructions in assembly language format; without the `-d`, `link` object text is printed. In addition, the files `test[2356].c[1055,21]` are meaningful programs, for which your compiler should be capable of producing executable code. To execute the test programs, the `link` object text must be converted to a DEC-10 relocatable file. This is accomplished by running `link` (which is on the Computer Science library) and `casm[1055,21]` and loading the resulting assembly language file and the skeleton runtime library `rtlib.rel[1055,21]` with `ccom`; for example, to compile and execute `test2.c` use the commands

```
c1 < test2.c | link | casm > test2.mac  
ccom test2.mac rtlib.rel[1055,21]  
a
```

My version of `cgen.c` is ~375 lines long.

163

### Linking

Write the two-pass link editor **link** for the object code described in the paper, 'A Machine Independent Linker'. **link** is invoked by the command

**link** [ *file* | **-o** *ofile* | **-u** ]...

Input is taken from *files*, in the order given. Standard input is read in the absence of *files* or if the file '-' is given. Output is written to *ofile*; if *ofile* is omitted, the output is written to the standard output. The **-u** option causes a list of undefined symbols to be printed on standard error. *only at end?*

Object text consists of command lines and polish suffix expressions. The following grammar defines the format of the input files.

```

file      : { cmd | expr }

cmd       : .def [ -s ] id expr
           | .seg id
           | .len id con
           | .org con

expr      : expr expr ( + | - | < | > ) | id | con

```

*loader units!*

Identifiers are sequences of letters, digits, percent signs, underscores, or periods, beginning with any of these except a digit. Constants are decimal or octal integers; a leading 0 indicates an octal constant. The semantics of each command and of expressions are described in the paper. Unrecognized commands should be treated as comments and copied to the output. *cmds?*

The version of **link** described in the paper is available in the Computer Science library and is invoked by the command given above. *c1.exe*[1055,21] is a working version of the C compiler that emits **link** object code in 36-bit units. Giving the **-18** option causes *c1* to emit object code in 18-bit units. Compile the programs, *test?.c*[1055,21] to produce object code for testing your linker. Additional test data will be provided later.

Use the following files in constructing your program.

<i>name</i>	<i>purpose</i>
* <i>link.h</i>	definitions
<i>lexpr.c</i>	expression evaluation
<i>lmain.c</i>	main program
<i>lpass1.c</i>	pass one
<i>lpass2.c</i>	pass two
* <i>lsym.c</i>	symbol table

Both the source and object files (*.rel* files) are available on [1055,21] for those files indicated by an asterisk.

The approximate lengths (in lines) of the files comprising my version of **link** are

152	<i>lexpr.c</i>	<i>mylink.h</i>
76	<i>lmain.c</i>	
141	<i>lpass1.c</i>	
114	<i>lpass2.c</i>	
103	<del><i>lsym.c</i></del>	
586	<i>total</i>	

*+ typedef Int long*

### Library Searching

Extend `link` to search a library for undefined references. This revised version of `link` is invoked by the command

```
link [ file | -o ofile | -l lfile | -u ]...
```

As before, input is taken from *files*, in the order given; standard input is read in the absence of *files* or if the file '-' is given. Output is written to *ofile*; if *ofile* is omitted, the output is written to the standard output.

During pass one, files are read in the order given. The `-l` option indicates that the library *lfile* should be searched in an attempt to resolve any undefined references. The search must be made immediately upon encountering the `-l` option using the set of undefined references at that point.

Libraries are simply archive files (see `ar` and *Software Tools*). Each entry in an archive begins with a header of the form

```
#-h- name size year date time
```

The *name* is the file name of the entry and *size* is its size in characters (not including the header). The remaining fields indicate the date and time of insertion into the archive. The first entry in the archive, named 'index', is an index that contains lines of the form

```
# symbol name
```

where *name* is the name of the entry in which *symbol* is defined. Note that the presence of the index makes it unnecessary to keep the library topologically sorted. Portions of the library may need to be read twice, however. Your program should read the library sequentially constructing an appropriate data structure to represent the library contents. If, after reading the library once, it still contains entries that define undefined but referenced symbols, use random access (see *seek* and *stell* in the C Programmer's Manual) to link the necessary library entries. There are several test libraries, named *lib?.lib*, and short test C programs, named *l?.c*, in [1055,21].

You must use your other `link` modules the revised version of *lsym.c*; place the library searching code in *lib.c*. My version of *lib.c* is about 180 lines long.

entry = library entry  
= library member

② l2.c + lib1  
uses - lxl.c + lib3  
code + lib3  
l1.c + lib1  
+ lib2  
l3.c + lib3  
l4.c + lib4  
+ lib3

mylink.h  
lexpr.c — not modified  
lib.c  
lmain.c  
lpass1.c  
lpass2.c

Three Weeks

### Low-Level Interactive Debugging

Write **db**, a low-level interactive debugger that allows examination, modification and execution of an executing program. The relocatable binary file for **db** is loaded with the program to be debugged using **ccom** and execution begins in **db**. For example, the command

```
ccom test2.rel rtlb.rel db.rel
```

loads **db** with the test program and creates the executable file *a.exe*. Execution of *a.exe* begins by calling the function `_db` in *db.rel* instead of `main` in *test2.rel*.

**db** issues the prompt "*db*>" and accepts the following edit-like commands.

<i>addr</i>	<b>b</b>	set a breakpoint at <i>addr</i>
<i>addr</i>	<b>c/x</b>	change word at <i>addr</i> to <i>x</i>
<i>addr</i>	<b>d</b>	remove the breakpoint at <i>addr</i>
	<b>J</b>	begin or resume execution of user's code
<i>addr</i>	<b>p fmt</b>	print word at <i>addr</i>
<i>addr</i>	<b>= fmt</b>	print <i>addr</i>
<i>addr</i>	<b>- fmt</b>	print word at <i>-addr</i>
<i>addr</i>	<b>newline</b>	print word at <i>++addr</i>

The *addr* argument, which is an octal address, is optional; if it is omitted, the most recently specified *addr* is used.

**p**, **=**, and **-** commands can be followed by a *sprintf* format string that is used in printing the relevant value. This format is remembered and reused in subsequent **p**, **-**, **newline**, and **=** commands until it is replaced. Initialize it to "`0%o\n`". A simplified version of *sprintf* in *sprintf.c[1055,21]* should be included in your program. In addition to the usual formats, it supports the "`%i`" format for printing words as instructions, assuming the array *opc* is initialized properly (see *casem.c[1055,21]*). The "newline" command increments *addr* and prints the word at the resulting address; the **-** decrements *addr* and prints the word at the resulting address; thus repeated **newline** and **-** commands can be used to step through memory. The output of **p**, **-**, and **newline** commands should be preceded by the associated address (in octal).

The **c** command changes the word at *addr* to *x*, where *x* is a signed decimal or octal constant. If *x* is omitted, the word is set to 0.

**b** and **d** commands are used to set and removed breakpoints at *addr*, respectively. **db** should support 5 breakpoints. The **J** command begins execution of the user's program or resumes its execution after stopping at a breakpoint.

Further details concerning the operation of **db** can be obtained by loading the relocatable binary file for my version, which is in *db.rel[1055,21]*, and experimenting with the test programs in *test[2356].rel[1055,21]*.

*db.h[1055,21]* contains the relevant definitions for **db**. *db.c[1055,21]* contains a skeleton version of **db** that can be used as a starting point; place your debugger in a file named *db.c*. You should "include" *sprintf.c[1055,21]*, which is a version of *sprintf*, and *swtch.c[1055,21]*, which includes the context switching routines, in your version (see *db.c*). Your version of **db** must define the following two routines. *get(addr)* which returns the contents of the word at *adr*, and *put(addr,x)*, which replaces the contents of the word at *adr* with *x* and returns the replaced value. Both of these routines must check the legality of *adr* to avoid addressing errors. In addition, *put* must clear the 'user write protect bit' if *adr* is a legal address in the high segment. This is accomplished by *\_setuwp(bit)*, which sets the write bit to *bit* and returns the previous value so that it can be reset after replacing the word at *adr*. *\_setuwp* is in *swtch.c[1055,21]*.

**J** commands are executed by calling *\_resume(xp,&dbstate,brks)*, where *xp* is a pointer to the state of the user's program (a *struct state \**, see *db.h*), *dbstate* is the structure that holds the state of **db** during the execution of the user's program, and *brks* is an array of pointers to breakpoint structures (see *db.h*). If *\_resume* returns, a breakpoint was encountered. Set the "most recently specified *addr*" to the address of the breakpoint and accept **db** commands.

My version of **db** is 140 lines long.

### High-Level Interactive Debugging

Write `cdb`, a high-level interactive debugger that allows examination, modification and execution of an executing C program. The relocatable binary files comprising `cdb` are loaded with the program to be debugged using `ccom` and execution begins in begins by calling the function `_db` in `cdb.rel` instead of `main`.

`cdb` issues the prompt "`cdb>`" and accepts the following commands.

<code>b</code>		print a list of breakpoints
<code>b name n</code>		set a breakpoint at execution point <code>n</code> in <code>name</code>
<code>c</code>		print the call to the focus
<code>C</code>		print the call to the focus including locals
<code>d</code>		delete all breakpoints
<code>d name n</code>		delete the breakpoint at execution point <code>n</code> in <code>name</code>
<code>J</code>		begin or resume execution of user's program
<code>p expr</code>		evaluate C expression <code>expr</code> and print the result
<code>s</code>		print the symbols for the focus
<code>S</code>		print the symbols for the focus and all globals
<code>T n</code>		set the user's <code>_trace</code> to <code>n</code>
<code>-</code>		move the focus "down" one activation
<code>+</code>		move the focus "up" one activation
<code>0-9! - n</code>		move the focus to the <code>n</code> th activation

Unrecognized commands are treated as `p` commands.

Debugging focuses on a "currently active procedure", or *focus*. Initially, there is no focus; after encountering a breakpoint, the function containing the breakpoint, whose activation record is on the top of the stack, becomes the focus. The runtime symbol table is organized so that the symbols for the focus and global symbols are available, in that order.

`-` and `+` commands change the focus by moving it to the function associated with the next activation record "down" or "up" the stack, where "up" means towards the top of the stack. Thus, for example, a series of `-` commands can be used to walk down the stack. A number, `n`, can be given as a command to set the focus to the function associated with the `n`th activation record from the top of the stack. Thus, `0` sets the focus to the function whose activation is at the top of the stack.

The `c` and `C` commands print the call to the focus. For example, if the focus is set to an activation of `queens`, `c` prints

```
queens(c=6)
```

`C` prints the call and lists the names and values of the locals:

```
queens(c=6)
  r=7
```

Whenever the focus is changed, by a command or by encountering a breakpoint, a `c` command is executed.

The `s` and `S` commands print lists of currently available symbols. The `s` command prints the name of the focus and its local symbols. For example, for the eight queens program, `s` prints

```
int queens()
  int r
  int c
```

`S` print all of the global symbols, including the locals for the focus:

```
int up[]
int down[]
int rows[]
int x[]
int main()
int queens()
```

```

    int r
    int c
    int print()
  
```

**p** commands evaluate C expressions and print the results. Variables appearing in expressions are interpreted according to C scope rules assuming the current function is given by the focus. Variables can be changed by assignment expressions.

**b** and **d** commands are used to set and remove breakpoints. Breakpoint locations are given as "execution points" within a function. An execution point is the address of the first instruction in the instruction sequence for an expression. Expressions are numbered sequentially beginning at 0 in the order of appearance in the function. The compiler emits a list of execution points immediately following the runtime symbol table entry for functions (see *db.h*[1055,21]). The list is terminated by a 0. **cdb** should support 5 breakpoints. The **J** command begins execution of the user's program or resumes its execution after stopping at a breakpoint.

Further details concerning the operation of **cdb** can be obtained by loading the relocatable binary file for my version, which is in *cdb.rel*[1055,21], and experimenting with the test programs in *test*[2356].*rel*[1055,21].

The relevant files for **cdb** are as follows.

<i>name</i>	<i>purpose</i>
* <i>db.h</i>	<b>db/cdb</b> definitions
* <i>cc.h</i>	compiler definitions
<i>cdb.c</i>	main program
* <i>eval.c</i>	expression evaluation
* <i>cezpr.c</i>	expression parsing
* <i>clz.c</i>	lexical analysis
* <i>csym.c</i>	symbol table management
* <i>swtch.c</i>	context switching primitives

Both the source and object files (*.rel* files) are available on [1055,21] for those files indicated by an asterisk. *cdb.c*[1055,21] contains a skeleton version of **cdb** that can be used as a starting point; place your debugger in a file named *cdb.c*.

The symbol table routine *lookup(str)* searches for the identifier given in *str* by first examining the locals and parameters for the function whose symbol table entry is pointed to by *\_focus* following by the globals. Thus, maintaining the focus is accomplished by setting *\_focus* correctly. Likewise, *foreach* calls its argument function for each global symbol and each local symbol in the focus. *ptype(p)* prints symbol table entry **p** in the format shown above; it should be used in the **s** and **S** commands.

*result* † *\*eval(str,frp)* (in *eval.c*) parses and evaluates the expression in *str* using the activation record represented by the *frame* structure pointed to by *frp* for accessing locals and parameters. *eval* results a pointer to a *result* structure. *frame* and *result* are defined in *db.h*. *eval* should be used for the **p** command.

*findframe(n)* (in *cdb.c*) positions a global frame pointer, *cursor*, to the *n*th frame from the top of the stack. It returns 0 if the positioning fails and 1 otherwise. *findframe* should be used for the **-**, **+**, and **number** commands.

As in **db**, **j** commands are executed by calling *\_resume(up,&dbstate,brks)*, where *up* is a pointer to the state of the user's program, *dbstate* is the structure that holds the state of **cdb** during the execution of the user's program, and *brks* is an array of pointers to breakpoint structures (see *db.h*). If *\_resume* returns, a breakpoint was encountered. Issue a message giving the execution point of the breakpoint, set the focus (*cursor* and *\_focus*) to the function associated with the frame at the top of the stack, execute a **c** command, and accept **cdb** commands.

My version of that portion of **cdb** not provided in *cdb.c* is 275 lines long.

Bill Mitchell

Q&C 327

Computer Science 453

Fall 1983

Class Notes, Set 1

August 25, 1983



# Software Layers

"applications" software  
e.g. software tools, editors,  
DBMS, etc.

CS  
430

"basic" systems software  
e.g. compilers, linkers, editors,  
performance monitoring, etc.  
i.e. essential tools

CS  
453

operating system  
system calls : { i/o services, etc.  
                  { processes

CS  
452

"bare" machine

## Essential Systems Software

- in presence of virtual machine presented by the hardware + operating system, i.e. instructions + system calls.
- essential tools are those necessary to construct other tools, i.e. translators + associated systems software.

- Editors for program preparation,

- CS 453 {
- • Compilers for language translations;
  - • Linkers for combining program modules,
  - • Debuggers for finding program errors
  - • Monitors for understanding program behavior + cost
  - Command language interpreters (e.g. a shell) for program invocation.

[Note absence of assembler?]

## Editors

- program preparation
- file manipulation/creation
- provide model of general user interface
- examples:
  - line-oriented, e.g. ed (UNIX), sos (Dec10)
  - display-oriented, e.g. Emacs, vi, edt, Rand editor, etc.
  - syntax-directed editors for editing program text, other "objects" Maintain correct format at all times; speed program preparation; bypass some compilation tasks. e.g. Cornell Program Synthesizer, Z, Emacs (sorta), sds.
  - display-editor interface to everything, i.e. editing as the computing metaphor.

## Compilers

- for systems programming language, i.e. language for writing operating system and other tools, e.g. BCPL, C, Pascal (sorta), Mesa (Xerox), Bliss, etc.
- must produce reasonable code, and be relatively simple; more importantly, must produce correct code!

### For CS 453 :

- write a recursive-descent compiler for a subset of C in (full) C, on the Dec-10.
- 1-pass, emitting link object code (will also write the linker!).
- machine-independent parser ("top half"); machine-dependent code generator ("bottom half").
- classical design.

## Linkers

- machine-independent link editor
- links textual object language
  - machine-dependent content
  - machine-independent form
- permits "iterative linking" i.e. output can be used as input.
- requires simple, machine-dependent loader
- includes library searching  
(libraries are archive files)
- replaces traditional use of an assembler as compiler's "last pass".

## Debuggers

- low-level debugging, i.e. at the instruction level
- use to built high-level debugger, i.e. at the C statement level.
- features
  - editor style interface
  - breakpoints
  - "single stepping"
  - procedure tracing
- requires modification of Compiler.
- programmer-defined debugging functions.

## Monitoring

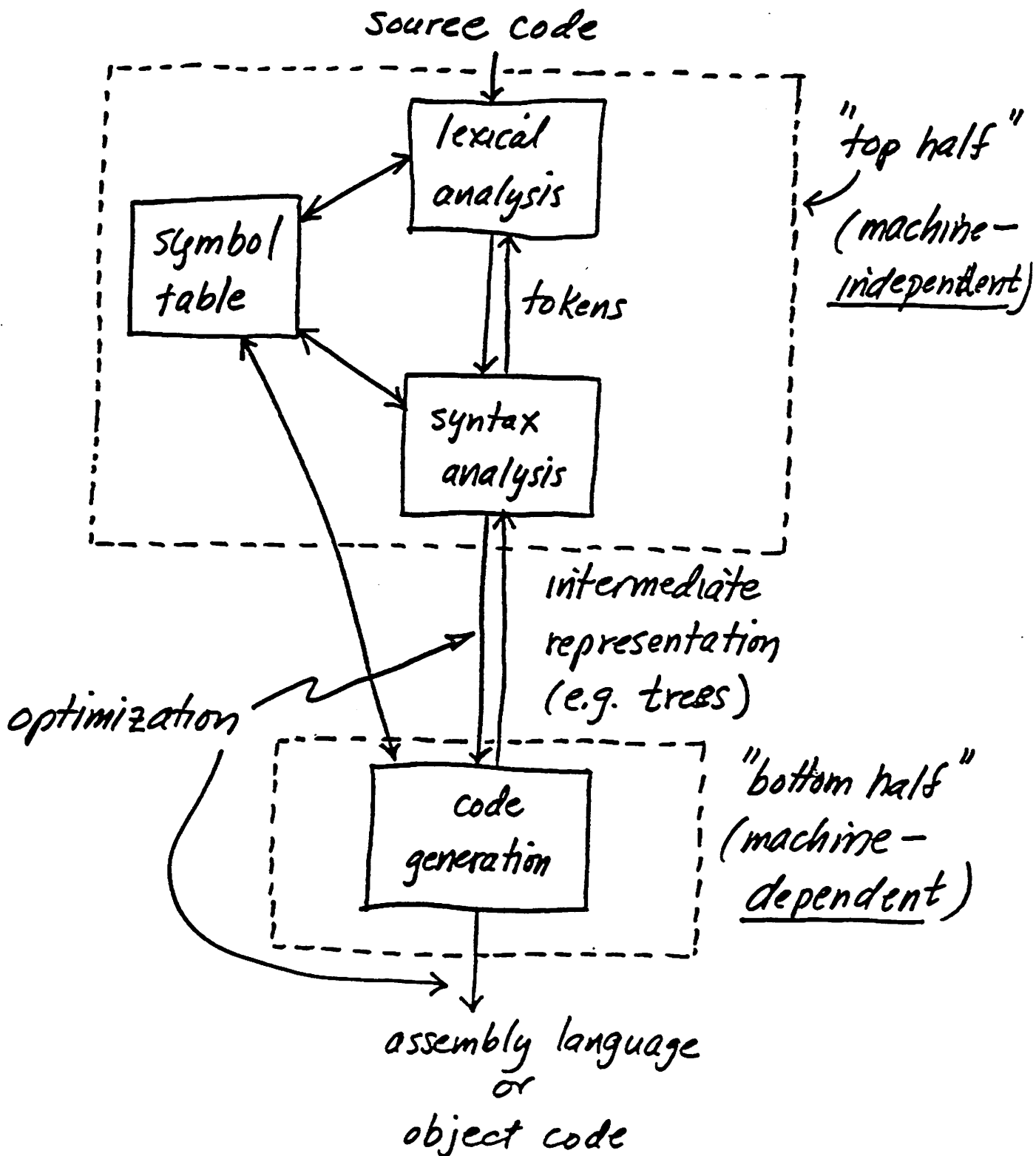
- statement level performance information
- execution frequencies  $\left\{ \begin{array}{l} \text{statements} \\ \text{procedures} \\ \text{~~program~~} \end{array} \right.$
- time histogram; sampling techniques.
- data presentation techniques; program "slices"; call graphs.

## Command Languages

- UNIX shell; shell programming
- display shells; "point & execute"
- interactiveness vs. "batch"  
(the UNIX shell is batch-oriented!).
- language design / binding issues /  
data type issues.
- unifying programming & command  
languages; EZ and related systems.



# Compiling: Classical Design



## Syntax Analysis

[ cparse.c  
cexpr.c  
cstat.c ]

- read tokens
- parse token stream according to a grammar that defines the language.
- perform semantic processing, eg. type checking, execution sequencing
- produce intermediate representation (IR), and interface with code generator, e.g. generate calls to emit jumps and define labels.
- error handling (syntactic & semantics errors).
- machine-independent processing

## Lexical Analysis [cllex.c]

- read source code ,  
return "tokens" , e.g.
- single chars : 'c'  $\Rightarrow$  'c'
- operators : '>='  $\Rightarrow$  GE  
↖  
symbolic constant
- constants:
  - 65  $\Rightarrow$  CON , "65"
  - 5.7  $\Rightarrow$  FCON , "5.7"
  - 'a'  $\Rightarrow$  CCON , "97"
  - "hi there"  $\Rightarrow$  SCON , "hi there"
- identifiers : puts  $\Rightarrow$  ID , puts
- keywords : else  $\Rightarrow$  ELSE
- also return "associated" value, e.g.  
symbol table pointer for constant,  
identifier string,  
value of integer constant.

## Code Generation

[ cgen1.c  
cgen.c ]

- read IR, e.g. abstract syntax trees, misc. directives ("generate a jump instruction").
- manage machine resources, e.g. registers
- assign other machine-dependent parameters, e.g. procedure frame layout.
- emit output, e.g. assembly language or object code

## Symbol Table [csym.c]

- maintain semantic information on identifiers, constants, e.g. type, address, size, etc.
- maintain multiple instances of same identifier, e.g. global function  $x$ , local integer  $x$ .
- maintain single-copy of all strings. two or more occurrences of " $x$ " result in entry of only one instance of " $x$ ", e.g. like unique software tool.
  - keep all unique strings in the input!
  - permits fast string comparison
  - permits "constant pooling"

[see csym.c, function stlookup]

## Grammars

grammar: precise specification of syntax

language: set (possibly  $\infty$ ) of sentences  
of terminal symbols (e.g. "a", sb).

$$L = L(T, N, P, S)$$

$T$  = vocabulary of terminal symbols

$N$  = set of non-terminal symbols  
("grammatical categories")

$P$  = set of productions  
("syntactic rules")

$S$  = start symbol from  $N$ .

Example: grammar given in Backus-Naur Form:

$E$	:	$E + E$	}	productions
$E$	:	$E * E$		
$E$	:	$( E )$		
$E$	:	$- E$		
$E$	:	$id$		

$$E : E + E \mid E * E \mid (E) \mid -E \mid id$$

"or" ↗

$$L = L(\{+ * ( ) - id EOF\},$$

$$\{E\},$$

$$\{E : E + E, E : E * E, E : (E), E : -E, E : id\},$$

$$E).$$

Derivation: repeated replacements (according to productions) to yield sentence from start symbol, e.g.

$$\bar{E} \Rightarrow -E$$

$$\Rightarrow -(E)$$

$$\Rightarrow -(E + E)$$

$$\Rightarrow -(id + E)$$

$$\Rightarrow -(id + id).$$

( this is a  
leftmost derivation )

Sentence recognition ("parsing"): reconstruct the derivation given the sentence.

## Top-down Parsing ("Recursive Descent")

- using the next input symbol and current derivation to properly "guess" the next derivation step.
- implement using a recursive procedure for each non-terminal.

### Pitfalls

1. left recursion, e.g.  $E : E + E$   
causes parser to loop.
2. backtracking, e.g.  $A : xB \mid xC$   
might have to "back-up";  
which production to apply?
3. ambiguity:  $E : E + E \mid E * E$   
"precedence" of + and \*?
4. error handling: what happens when we get "stuck"?



## Solutions

3. ambiguity: add more productions to enforce precedence, associativity, etc.  
e.g.

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

becomes

$$E : E + T \mid T$$

$$T : T * F \mid F$$

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

1. left recursion: replace  $A : A\alpha \mid \beta$  with  $A : \beta A'$  and  $A' : \alpha A' \mid \epsilon$  (i.e. right recursion).

" $\beta$  followed by 0 or more  $\alpha$ 's"

$$\begin{cases} E : TE' \\ E' : +TE' \mid \epsilon \end{cases}$$

$$\begin{cases} T : FT' \\ T' : *FT' \mid \epsilon \end{cases}$$

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

## Solutions, cont'd

2. backtracking: in  $A : \alpha \mid \beta$ , must be able to determine which alternative,  $\alpha$  or  $\beta$ , to choose based on current input symbol.

Let:  $\text{first}(\alpha)$  = terminals that begin sentences derived from  $\alpha$ .

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

Computing  $\text{first}(x)$ : continue until no change.

a) if  $x$  is a terminal,  $\text{first}(x) = \{x\}$

b) if  $x : a\alpha$ , where  $a$  is a terminal, add  $a$  to  $\text{first}(x)$

c) if  $x : \epsilon$ , add  $\epsilon$  to  $\text{first}(x)$

d) if  $x : \underbrace{Y_1 Y_2 \dots Y_{i-1}}_{\neq \epsilon} Y_i \dots Y_k$ , add  $-\{\epsilon\}$

$\text{first}(Y_1) \cup \overset{\neq \epsilon}{\Rightarrow} \text{first}(Y_2) \dots$  to  $\text{first}(x)$

and add  $\text{first}(Y_i)$  to  $\text{first}(x)$ .

i.e. can "look thru" via  $\epsilon$ -productions

$Y_1 \dots Y_{i-1}$ .

Let  $\text{follow}(A)$  = terminals that can follow non-terminal  $A$ .

Computing follow : repeat until no change

- if  $S$  is the start symbol, add EOF to  $\text{follow}(S)$ .
- if  $A : \alpha B \beta$  and  $\beta \neq \epsilon$ , add  $\text{first}(\beta) - \{\epsilon\}$  to  $\text{follow}(B)$ .
- if  $A : \alpha B$  or  $A : \alpha B \beta$  and  $\epsilon \in \text{first}(\beta)$ , add  $\text{follow}(A)$  to  $\text{follow}(B)$ .

Example :

<u>rule</u>	<u>first</u>	<u>follow</u>
$E : TE'$	{ ( id }	{ ) EOF }
$E' : +TE'   \epsilon$	{ + $\epsilon$ }	{ ) EOF }
$T : FT'$	{ ( id }	{ + ) EOF }
$T' : *FT'   \epsilon$	{ * $\epsilon$ }	{ + ) EOF }
$F : (E)   id$	{ ( id }	{ * + ) EOF }

(see Asgn #1 handout)

## Extended BNF

- recursion and  $\epsilon$ -productions are troublesome
- extend BNF to include explicit repetition and optional constructs
- use syntax diagrams.

### Extensions

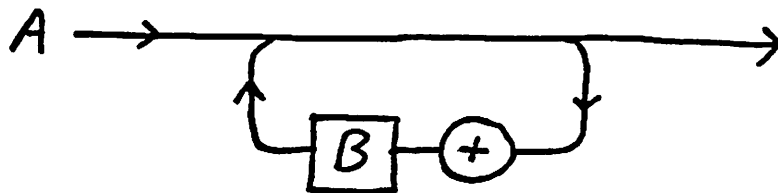
$[a]$  stands for  $a | \epsilon$   
 $\{a\}$  " "  $\epsilon | a | aa | \dots$   
 $(a|b)c$  " "  $ac | ab$

### Example

$E' : \epsilon | +T | +T+T | \dots \Rightarrow \{+T\}$   
 $E : T \{+T\}$   
 $T : F \{*F\}$   
 $F : (E) | id$

## first & follow sets

- first - computed by initial branches in diagram
  - includes  $\epsilon$  if diagram can be traversed without encountering  $\square$  or  $\circ$ .
- follow - computed from examining where  $\square$  is used.



$$\text{first}(A) = \{ + \epsilon \}$$

follow(B) includes follow(A)  $\cup$  {+}

Want :

1. every fork to be selectable by looking only at next i.e.  $\text{first}(\alpha) \cap \text{first}(\beta) = \emptyset$
2. if any graph can be traversed without reading input, label its exit path with follow set (which affects  $\neq \epsilon$ ).

# Translating EBNF $\rightarrow$ Diagrams

1. terminal  $x$

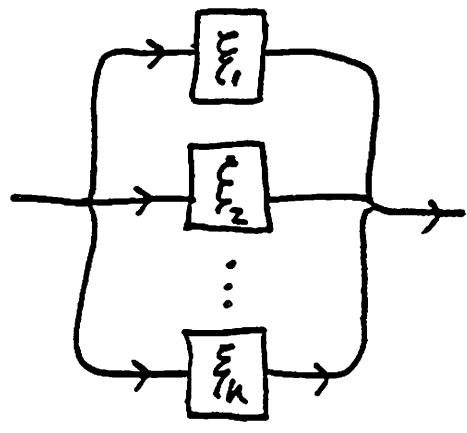


2. non-terminal  $A$

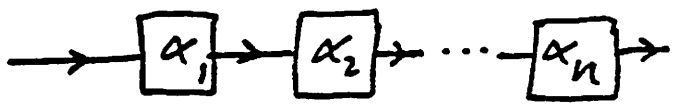


3.  $A : \xi_1 | \xi_2 | \dots | \xi_n$

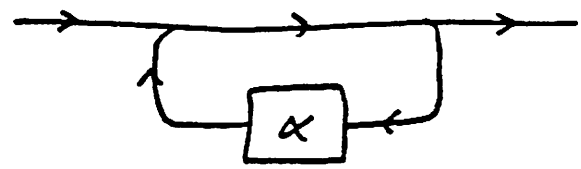
where  $\xi_k$  is  
syntax diagram for  $\xi_k$



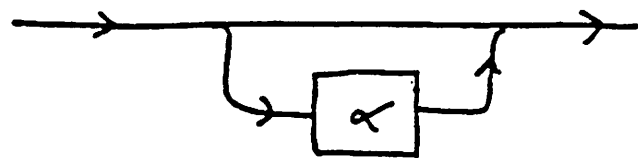
4.  $A : \alpha_1 \alpha_2 \dots \alpha_n$



5.  $\{ \alpha \}$



6.  $[ \alpha ]$



## Translating Diagrams $\rightarrow$ C Code

input tokens:  $t = \text{gtok}();$   
 $\uparrow$  "current token"

1. reduce graphs to few manageable graphs via substitution.
2. for each graph  $S$ , write a procedure  $T(S)$  that "traverses" the graph guided by the input.

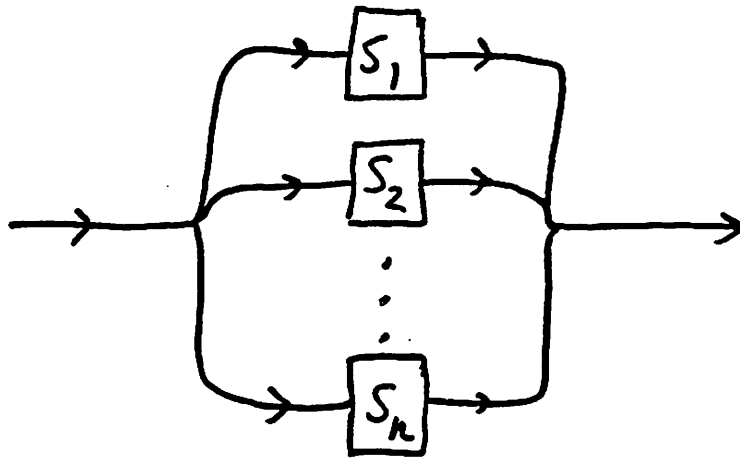
Use following templates:



```

{ T(S1);
  T(S2);
  ⋮
  T(Sn); }
  
```

i.e. non-terminals  $\Rightarrow$  procedure calls

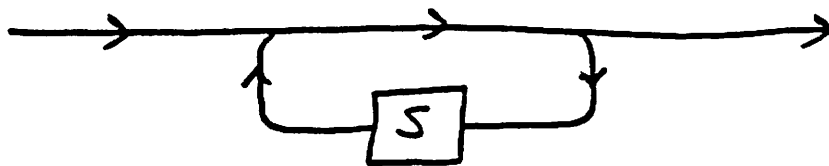


```

if (  $t \in \text{first}(S_1)$  )  $T(S_1)$ ;
else if (  $t \in \text{first}(S_2)$  )  $T(S_2)$ ;
  ⋮
else if (  $t \in \text{first}(S_n)$  )  $T(S_n)$ ;
else err( "syntax error" );

```

(or use equivalent switch statement).



```

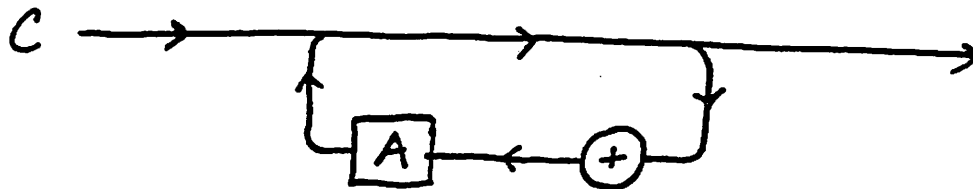
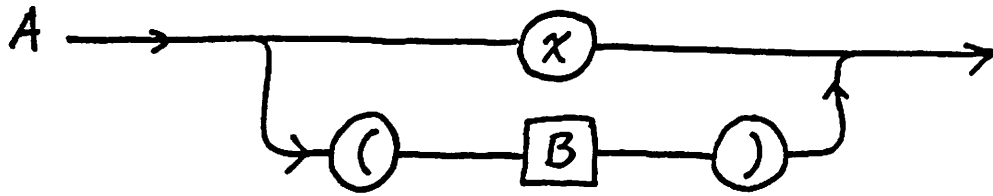
while (  $t \in \text{first}(S)$  )
   $T(S)$ ;

```

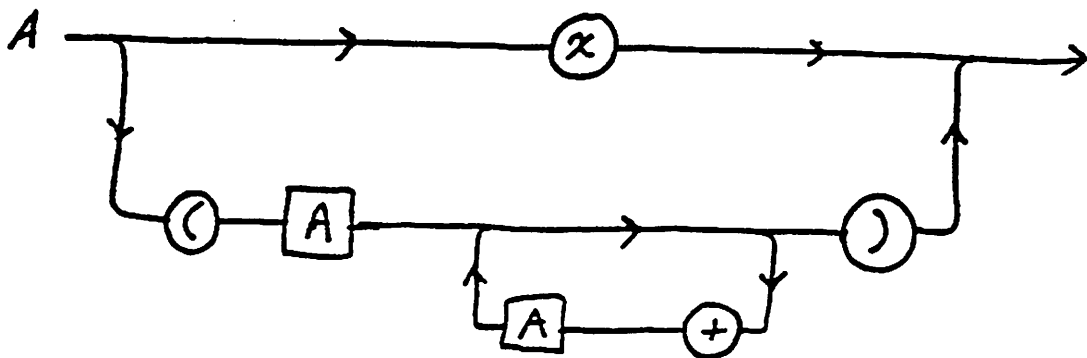


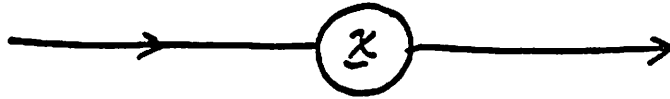
Example :  $A : x | '( B )'$   
 $B : AC$   
 $C : \{ + A \}$

sentences are  $x, (x), (x+x), ((x)), \dots$



Could draw one diagram for A:





```
if ( t == x )
    t = getToken();
```

```
else err( "missing x", "" );
```

Example : syntax diagram for A

```
int t;
main()
{
    t = getToken();
    A();
    if ( t != EOF )
        err( "syntax error" );
}
```

(after some refinement)

A()  
{

if (t == 'x')

t = getToken();

else if (t == '(') {

t = getToken();

A();

while (t == '+') {

t = getToken();

A();

}

if (t == ')')

t = getToken();

else err("missing ");

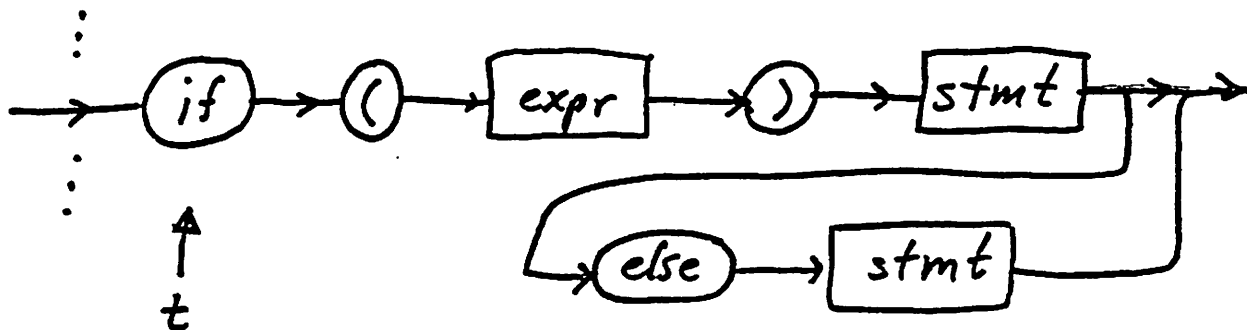
}

else err("syntax error");

}

Example: C if statement

stmt:



stmt()

{

switch (t) {

⋮

case IF:

if stat ();

break;

⋮

default: err("unrecognized statement")

}

}

```

ifstat()          /* without semantics */
{                (5 more lines!)

    t = gtok();
    mustbe ('(');
    expr();
    mustbe (')');
    stmt();

    if (t == ELSE) {
        t = gtok();
        stmt();
    }
}

```

```

mustbe(c)
{
    static char missing[] = "missing X";
    if (t == c)
        t = gtok();
    else {
        missing[8] = c;
        err(missing, "");
    }
}

```

## Compiler Organization

(see grammar)

cparse.c { prog  
dcl  
dclr  
type  
ptr  
func

cstat.c { stmt

cexpr.c { expr  
term

### Advice

1. write a procedure (or several) for each of the above non-terminals.  
Don't take short cuts!
2. resolve ambiguities locally; i.e. ignore them, then revise code to fix them.



- grammar accepts invalid function decls; use semantic checks
  - globals: install as usual, but don't turn on DEF (establishes type for forward references).
  - parameters: illegal.
  - locals: install as global, as above.

### Strategy for writing dclr & dcl:

```

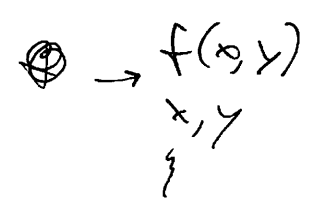
dclr(type, f)
  int type; /* type of this list */
  struct symbol *(*f)();
  {
    ... gather up one declaration
      item ... call appropriate routine
    (*f)(name, type, size, ...)
  }

```

symbol

$f$  is a function; one for each of global, parameter, local, which installs symbol & performs semantic checks.



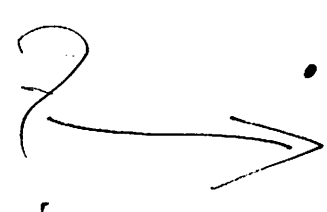


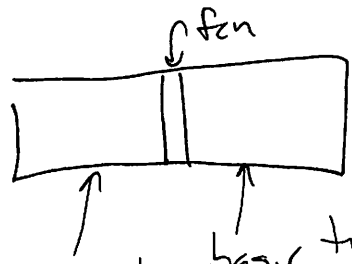
```
dcl(type, f)
int type;
struct symbol * (*f)();
{
  while (t == ID || t == '*') {
    dclr(type, f);
    if (t == ',')
      t = gotoK();
  }
}
```

```
⋮
/* parse locals */
while (t == CHAR || t == INT || t == FLOAT) {
  dcl(type(), declLocal);
  mustbe(',');
}
```

installs local identifiers

- Note:
- must always check for erroneous redeclaration, or for missing or extraneous declarations for parameters.
  - need to call global(p) for global variables (will allocate storage).





# Symbol table entries

Be sure to use the macros

```

struct symbol {
    char *s_name; /* name */
    int s_flags; /* scope & other flags */
    int s_type;
    int s_size; /* size in bits */
    int s_offset; /* for use by bottom half,
                  arg numbers, etc. */
    struct symbol *s_link;
};

```

flags: DEF REF  
 GLOBAL · LOCAL PARAM CONST  
 KEYWORD

## Examples

<u>input</u>	<u>s_type</u>	<u>s_size</u>
int a;	int	36
char s[8];	char *	72
float *y[10];	float **	360

~~any other~~

## Compiling Procedures

- straightforward traversal of syntax diagrams (once ambiguity is resolved).

func()

{

install function name as global;  
 parse parameter list, setting  
 s-offset of each parameter  
 to parameter number;  
 parse parameter declarations;  
 check for semantic errors in params.;  
 mustbe('{');

parse locals;

x = proc beg (symbol table entry  
 for function name);

stmt list();

retcode(0); /\* gratis return \*/

mustbe('}');

proc end(x);

expunge();

} order is important!

}

## Compiling Statements

- make calls to `deflabel`, `jump`, and expression code generator at appropriate points in syntax diagram.
- use three contiguous labels for loops, e.g.  $L$ ,  $L+1$  ("continue" label),  $L+2$  ("break" label).
- suggested organization:
  - `stmtlist`
  - `stmt`
  - `ifstat`
  - `whilostat`
  - `forstat`
  - (others-inline code)
- write a label generator; e.g. `genlabel(n)` generates  $n$  labels, returns first one.  
Uses integers for labels first call - ret 1
- generate code for expressions via `walk(p, tlab, flab)`; where
  - $p$  - pointer to root of expression tree
  - $tlab$  - "true" label
  - $flab$  - "false" label

stmt()

```

{
  switch (t) {
    case first(expr):
      walk(expr(), 0, 0);
      mustbe(';');
      break;
    case IF:
      ifstat(genlabel(2));
      break;
    case WHILE:
      whilestat(genlabel(3));
      break;
    :
    default:
      err("unrecognized stmt", "");
  }
  if (t != follow(stmt)) {
    err("illegal stmt termination", "");
    skip to t < follow(stmt);
  }
}

```

Version of expr that eats only id + ret id node

$\text{if}(\text{expr}) \text{stmt}_1 \text{ else } \text{stmt}_2$

$\text{if}(\sim \text{expr}) \text{ go to } L$

$\text{stmt}_1$

$\text{go to } L+1$

$L: \text{stmt}_2$

$L+1:$

} omit if no else clause

"false" label

$\text{ifstat}(\text{lab})$

$\text{int lab};$

{

$t = \text{g tok}();$

$\text{mustbe}('(');$

$\text{walk}(\text{expr}, 0, \text{lab});$

$\text{mustbe}(')');$

$\text{stmt}();$

$\text{if}(t == \text{ELSE}) \{$

$t = \text{g tok}();$

$\text{jump}(\text{lab} + 1);$

$\text{deflabel}(\text{lab});$

$\text{stmt}();$

}

$\text{else } \text{deflabel}(\text{lab});$

$\text{deflabel}(\text{lab} + 1);$

}

## Loops

while (expr) stmt

```

L
L+1  if (~expr) goto L+2
      stmt
      goto L

```

→ L+2  
L+2 is always break label if  $ctrl$  is; ~~add~~ no expr, etc.

for (expr<sub>1</sub>; expr<sub>2</sub>; expr<sub>3</sub>) stmt

optional!

```

      expr1
L      if (~expr2) goto L+2
      stmt
L+1    expr3
      goto L
L+2

```

---

break

```

      goto L+2
continue
      goto L+1

```

} L is current  
loop "handle"  
error is there is  
no current loop!

Example

```

for(i = 1; i < 10; i++) {
    sum += x[i];
    if(sum > 100)
        break;
}

```

---

```

i = 1
L1:  if (i >= 10) go to L3
      sum += x[i]

      if (sum <= 100) go to L4
      go to L3
L4: L5: L2:
      i++
      go to L1
L3:

```

---

L1, L2, L3 - for loop labels  
L4, L5 - if labels



## Error Handling

- lexical errors ( " .... )
- syntactic errors
  - extra symbols
  - missing symbols
  - use of incorrect symbol
  - transposition
- semantic errors
  - undeclared identifier — uses "?" and int type + uses it on demand
  - break/continue outside a loop
  - duplicate declarations
  - missing declarations (e.g. parameter)
  - illegal array size
  - type conflict
- strategy: continue parsing, skipping input if necessary; must make forward progress.
- maintain correct parsing state according to first/follow sets.
- do something "reasonable" for semantic errors, e.g. enter undeclared identifiers, assume int for type errors.
- avoid "cascading"

## Errors

- missing symbols, e.g. ')' - assume it was there.

```

mustbe(c)
int c;
{
    if (t == c)
        t = getk();
    else {
        err("missing c", "");
    }
}

```

do not advance input →

- $t \notin \text{first}(X)$  - skip input until  
 $t \in \text{first}(X)$  or  $t \in \text{follow}(X)$   
or  $t == \text{EOF}$ .

alternative: skip one input symbol  
 (useful in expression parsing).

- $t \notin \text{follow}(X)$  - skip input until  
 $t \in \text{follow}(X)$  or  $t == \text{EOF}$
- use more restricted sets in some cases,  
 e.g. stmts.

## Semantic Errors

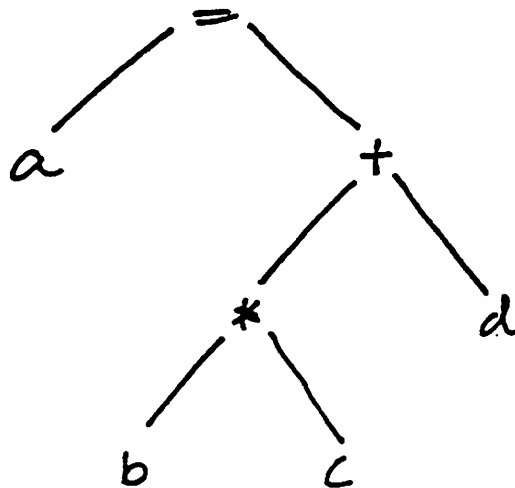
- maintain correct expression trees, i.e.
  - operators must have proper number of arguments
  - assume type information in presence of type errors, e.g. T\_INT.
- install undeclared identifiers as globals with type int.
- ignore duplicate declarations (but issue msg!)
- assume 1 for illegal array sizes
- ignore extraneous break/continues.
- bottom line: make sure code generator gets "reasonable" parameters, even if code is incorrect.

## Compiling Expressions

- parse expression syntax.
- construct trees to represent expressions.
- perform semantic checks (i.e. type checking) as trees are built.
- traverse trees calling expr code to generate code, in correct order. (i.e. postorder traversal).

Example: int a, b, c, d;

a = b \* c + d;



a b c d + =

## Expression Nodes

```

struct node {
    int e_op;    /* operator O_xxx */
    int e_type; /* type of result */
    struct node *e_left;
    struct node *e_right; /* descendants */
    int e_result;
}
int e_count; /* ref count */

```

### Operators:

- |            |                                     |
|------------|-------------------------------------|
| • O_ASSIGN | } operators expected by<br>exprcode |
| • O_NEG    |                                     |
| • ⋮        |                                     |
| • O_CON    |                                     |
- may use other operator codes to guide traversal, e.g. ||

if (a < 10 || b == 6) ...

↑  
never passed to exprcode,  
but exists in tree.

addr val  
Lvalues / Rvalues

- Rvalue: value associated with a name.
- Lvalue: location associated with a name.

$a = b$   
 ↗ need a's lvalue      ↖ need b's rvalue

- must distinguish between pointers (which are rvalues) and lvalues, e.g.

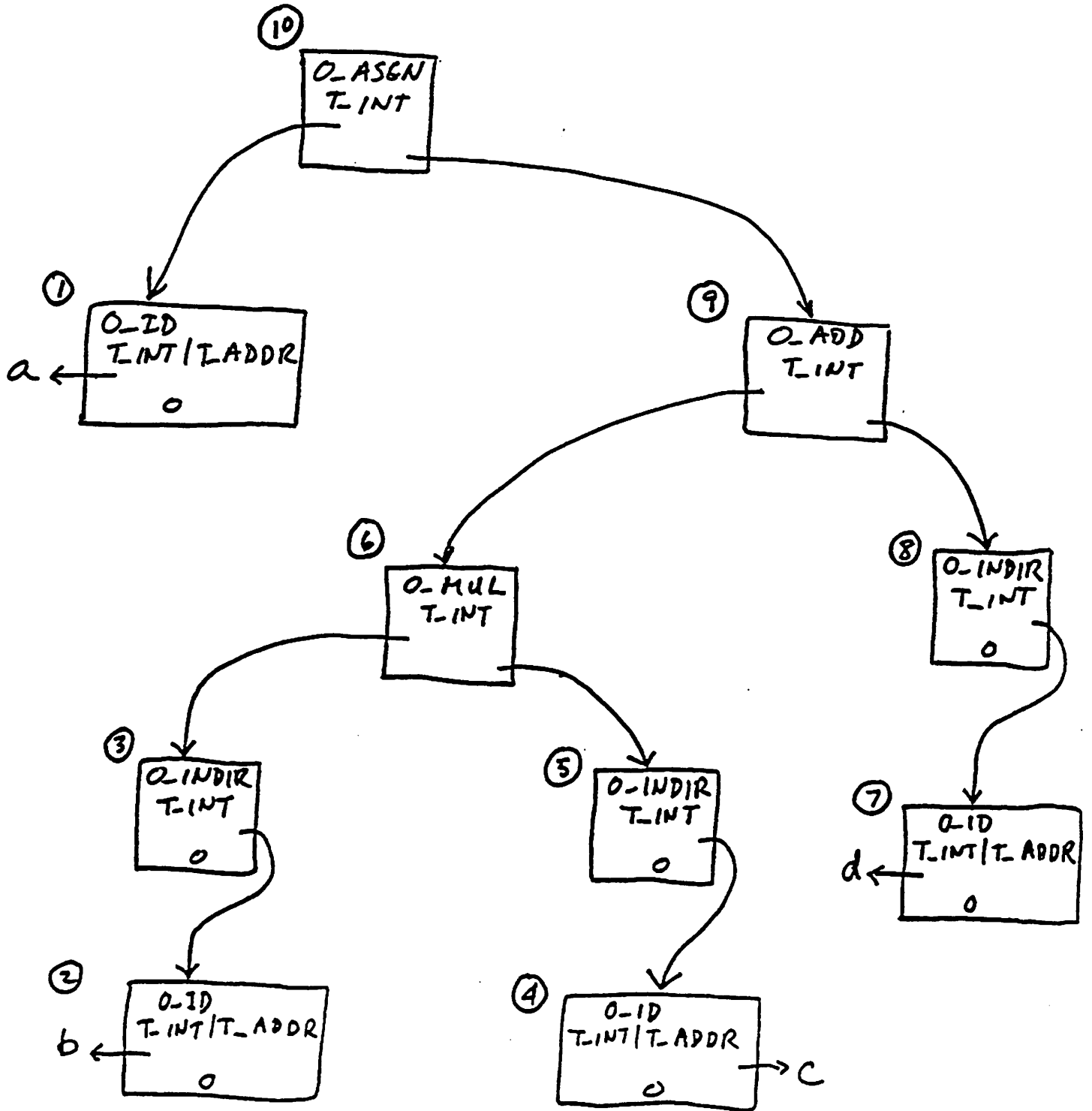
int a, \*p;

↗ takes rvalue & rets it as lvalue  
 ↗ need rvalue  
 ↖ returns lvalue: the location referred to by the value of p.

\*p = a;

- use T\_ADDR type flag to denote lvalue.
- use O\_INDIR (unary \*) to "convert" lvalue → rvalue.

$$a = b * c + d$$



• (i) denotes order of calls to exprcode.

$$a = b * c + d$$

$r1 = \text{exprcode}(O\_ID, 0, 0, \textcircled{1});$

$T1 = \&a$  ( $r1$  is assigned  $T$ )

$r2 = \text{exprcode}(O\_ID, 0, 0, \textcircled{2});$

$T2 = \&b$

$r3 = \text{exprcode}(O\_INDIR, r2, 0, \textcircled{3});$

$T3 = *T2$

$r4 = \text{exprcode}(O\_ID, 0, 0, \textcircled{4});$

$T4 = \&c$

$r5 = \text{exprcode}(O\_INDIR, r4, 0, \textcircled{5});$

$T5 = *T4$

$r6 = \text{exprcode}(O\_MUL, r3, r5, \textcircled{6});$

$T6 = T3 * T5$

$r7 = \text{exprcode}(O\_ID, 0, 0, \textcircled{7});$

$T7 = \&d$

$r8 = \text{exprcode}(O\_INDIR, r7, 0, \textcircled{8});$

$T8 = *T7$

$r9 = \text{exprcode}(O\_ADD, r6, r8, \textcircled{9});$

$T9 = T6 + T8$

$r10 = \text{exprcode}(O\_ASGN, r1, r9, \textcircled{10});$

$*T1 = T9$



## Parsing Expressions

- consider subset of C expressions:

$expr : expr '=' expr$   
 $| expr ('+' | '-' | '*' | '/') expr$   
 $| term$

$term : '++' term$   
 $| term '++'$   
 $| id$

$++ ++ ++ ++ ++ ++$   
 $+ \rightarrow + ++$   
 $+ \rightarrow + ++ ++$   
 $id id ++ ++$

- grammar is a problem:
  - left recursion
  - ambiguous (no precedence/assoc.)
- revised grammar:

$expr : expr1 \{ = expr \}$

$expr1 : expr2 \{ (+ | -) expr2 \}$

$expr2 : expr3 \{ (* | /) expr3 \}$

$expr3 : term$

$term : ( ++ term | id ) \{ ++ \}$

gives right assoc.

```
struct node * expr()
```

```

{
  struct node * p;
  p = expr1();
  while ( t == '=' ) {
    t = getToken();
    p = enode( O_ASSIGN, p, expr() );
  }
  return (p);
}

```

check follow set

No point in checking follow set in other routines.

performs type-checking, allocates nodes

```
struct node * expr1()
```

```

{
  int op;
  struct node * p;
  p = expr2();
  while ( t == '+' || t == '-' ) {
    op = t == '+' ? O_ADD : O_SUB;
    t = getToken();
    p = enode( op, p, expr2() );
  }
  return (p);
}

```

if (c, i, e) (c)? i:e

Y \* Z + P;

```
struct node *expr2()  
{  
    int op;  
    struct node *p;  
    p = expr3();  
    while (t == '*' || t == '/') {  
        op = t == '*' ? O_MUL : O_DIV;  
        t = getok();  
        p = enode(op, p, expr3());  
    }  
    return (p);  
}
```

```
struct node *expr3()  
{  
    return (term());  
}
```

```

struct node *term()
{
    struct node *p;
    if (t == INCR) {
        t = gtok();
        p = enode(—, term(), 0);
    }
    else if (t == ID) {
        p = enode(0-ID, lookup(lexval.l-str), 0);
        t = gtok();
    }
    else err(—);
    while (t == INCR) {
        p = enode(—, p, 0);
        t = gtok();
    }
    return (p);
}

```

?  
 ↓  
 ?  
 ↑  
 add ident

?  
 {

- idea: parameterize  $expr_1$ ,  $expr_2$  sequence to handle arbitrary number of precedence levels.

```

int optab[][3] = {
    { '+', '-', 0 },
    { '*', '/', 0 },
    { 0 }
};

```

tokens for each precedence level  
 assign doesn't fit model  
 Dealig w/ lower level ops.

if run off end, call term

```

struct node *expr()
{
    struct node *p;
    p = expr1(0);
    while (t == '=' || binop) {
        t = getk();
        p = enode(O_ASSIGN, p, expr());
    }
    return (p);
}

```

"This is sort of the only neat part of the program"

struct node \*expr1(n) ← ~~part of~~ ~~optab~~

int n;

{

int op;

struct node \*p;

if (\*optab[n] == 0)

return (term());

if ~~of~~ level, lowest call term, p = expr1(n+1); - expr1 call's expr

while (t < optab[n]) {

op = "correct operator"

t = strtok();

p = enode(op, p, expr1(n+1));

have aug. assignment;

return (p);

}

p: (struct node \*) - alloc (stack (avoid null); sym-~~table~~ alloc

while (op = inset(t, optab[n])) {

t = strtok();

p = enode(op, p, expr1(n+1));

...

Bill Mitchell

Computer Science 453

Fall 1983

Class Notes, Set 2

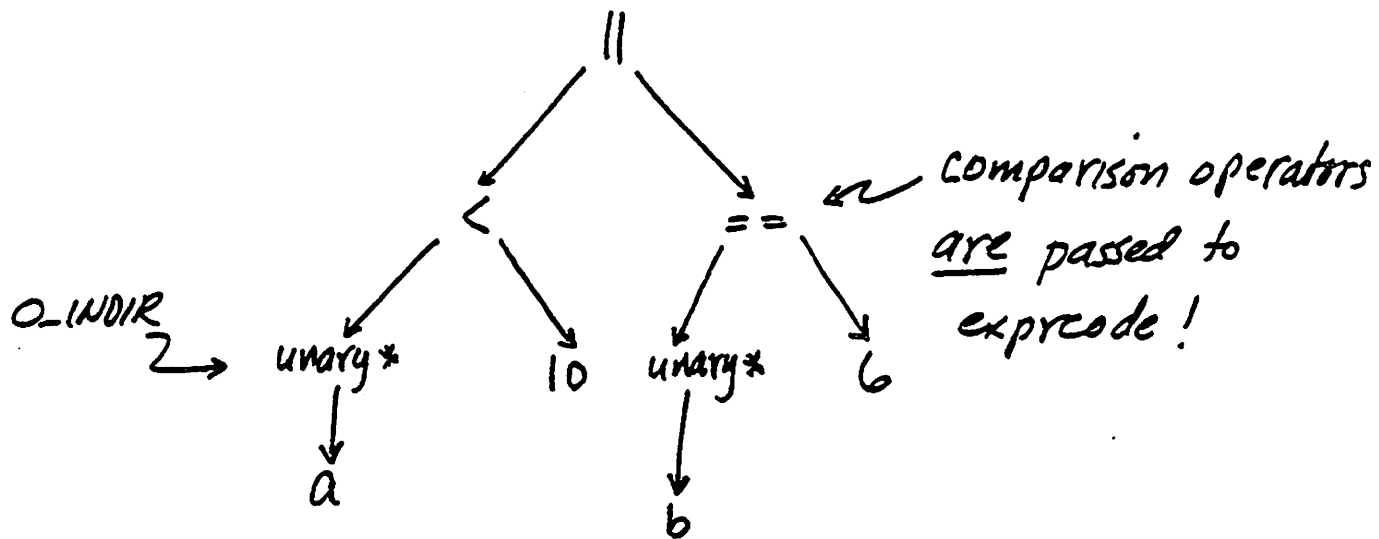
August 29, 1983

Start of notes

## Conditional Expressions

- unary  $!$ ,  $||$ ,  $\&\&$  alter flow of control;  
do not produce values.
- included in trees to guide traversal;  
not passed to exprcode.

Example:  $a < 10 \ || \ b == 6$



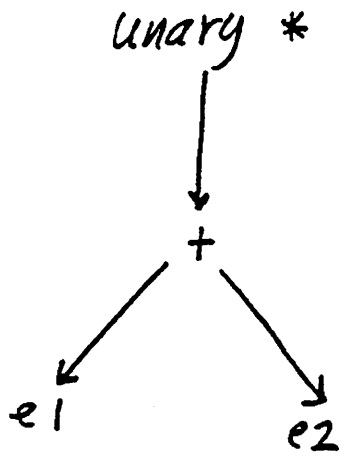
- general technique: use additional operator  
as necessary to guide traversal or to  
simplify code generation, e.g.  $x[e]$



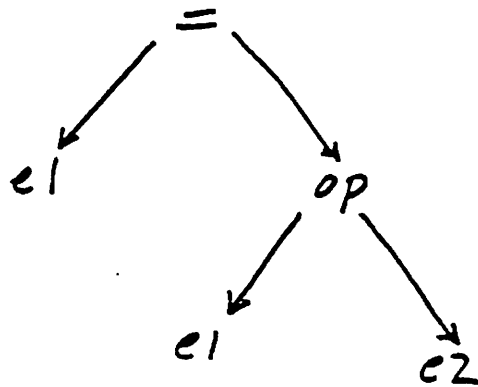
$e1[e2]$

$\therefore p = \text{enode}('[' , p, \text{expr}());$

$\text{enode}$  returns  $*(e1 + e2)$  :



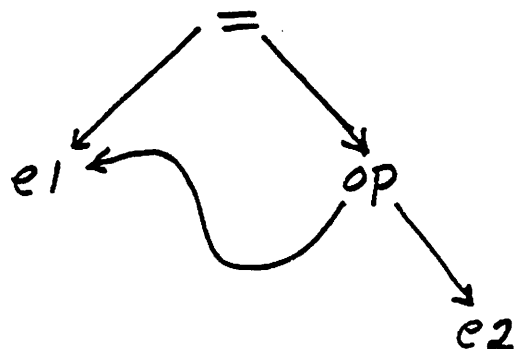
$e1 \text{ op} = e2$



Examples:  $a += 5;$

$*f() += a;$  ← oops!  $f()$  is called twice!

$\therefore$  must evaluate  $e1$  only once.

$$e1 \ op = \ e2$$


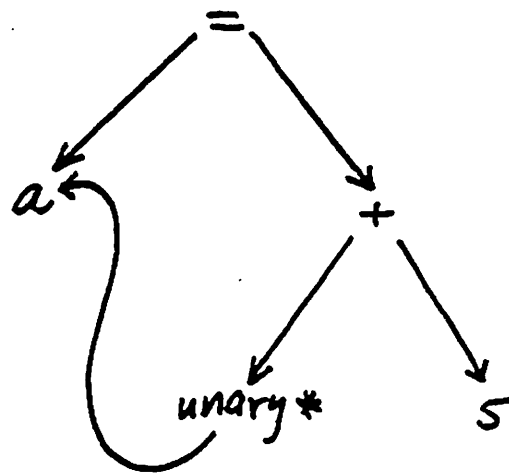
- need to "remember" result of expr code for 1<sup>st</sup> evaluation of  $e1$ , and re-use it for 2<sup>nd</sup> evaluation.
- store "result" for each node in `e_result` field
- `walk(...)` returns `e_result` field instead of traversing the subtree
- technique simplifies type-checking, e.g. can compile  $e1 += e2$  as  $e1 = e1 + e2$ ; special type-checking for `+=` is unnecessary.

res #

$e_1 \text{ op} = e_2$

- left operand of op must be an rvalue;  $e_1$  must be an lvalue.

Example:  $a += 5;$



•  $++e \iff e += 1$

•  $--e \iff e -= 1$

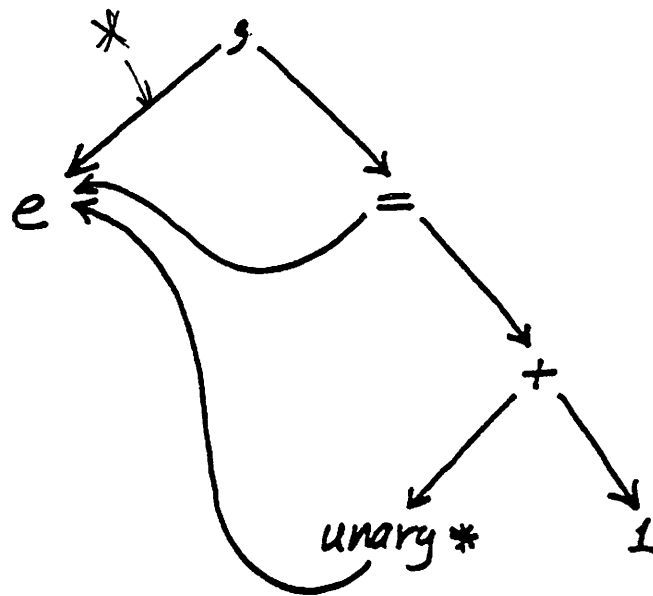
•  $e++ \iff (t' = e, e += 1, t')$

•  $e-- \iff (t' = e, e -= 1, t')$

need to fetch value, perform assignment, return value.

e++

- introduce new operator that evaluates both operands, but returns result of left operand:



- e-- is similar.
- 

Tree Traversal

$r = \text{walk}(p, \text{tlab}, \text{flab})$

- traverses tree rooted at  $p$ , calling `exprcode` as necessary.
- `tlab` ("tree" label), `flab` ("false" label) are used for conditional expressions.
- returns result of calling `exprcode`.

```

int walk(p, tlab, flab)
struct node *p;
int tlab, flab;

```

```
{
```

```
    int r, r1, r2, op;
```

```
    if (p == NULL)
```

```
        return (0);
```

```
    if (r = p->e.result)
```

```
        return (r);
```

```
    r = r1 = r2 = 0;
```

```
    switch (op = p->e.op) {
```

```
        :
```

```
        case O_ADD:
```

```
            :
```

```
            (set r)
```

```
            break;
```

```
        default:
```

```
            fprintf(stderr, "compiler error");
```

```
            exit(1);
```

```
    }
```

```
    p->e.result = r;
```

```
    return (r);
```

```
}
```

result already  
calculated?

case O\_ID: case O\_CON:

r = exprcode(op, 0, 0, 0, p);  
break;

op is  
+ comes from p

case O\_ADD: case O\_SUB: case O\_MUL:

case O\_DIV: case O\_MOD: case O\_LSH:

case O\_RSH: case O\_BXOR: case O\_BOR:

case O\_BAND:

case O\_NEG: case O\_INDIR:

case O\_BCOM: case O\_CVI: case O\_CVF:

case O\_ASGN:

r1 = walk(p → e\_left, 0, 0);

r2 = walk(p → e\_right, 0, 0);

r = exprcode(op, r1, r2, 0, p);

break;

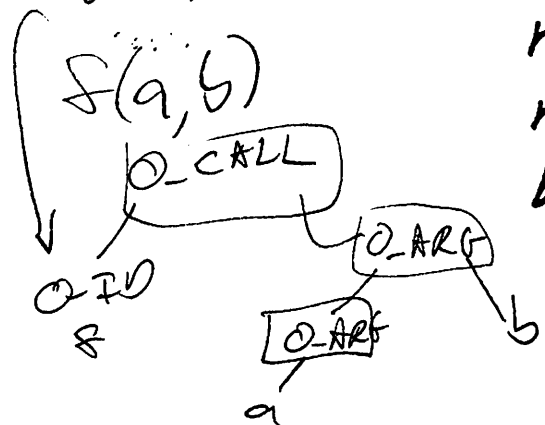
case O\_CALL:

r2 = walk(p → e\_right, 0, 0);

r = exprcode(op, 0, r2, 0, p);

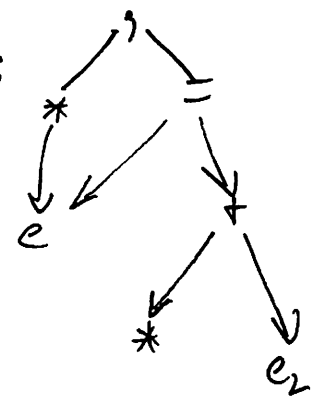
break;

full ~~ac~~ could  
be expr.



case ',': <sup>only for e++ e++=</sup>

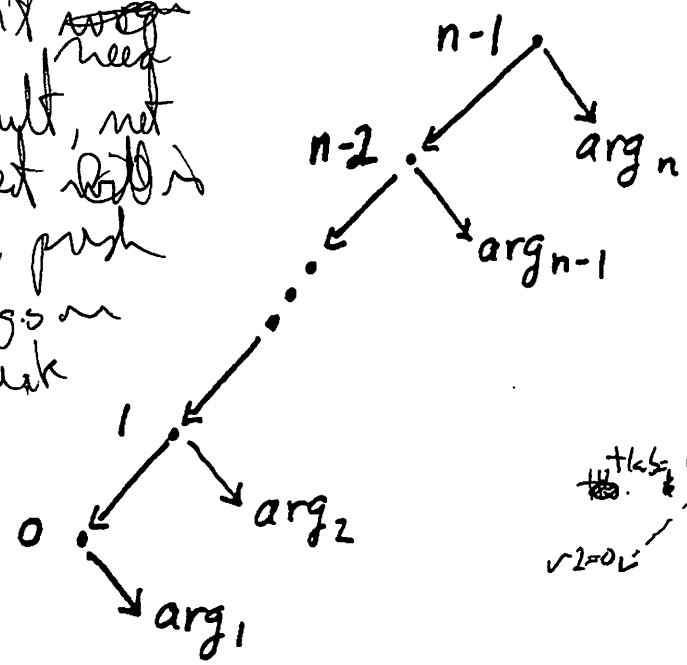
```
r = walk(p -> e_left, 0, 0);
walk(p -> e_right, 0, 0);
break;
```



case O\_ARG:

```
r2 = walk(p -> e_left, ++tlab, 0);
r1 = walk(p -> e_right, 0, 0);
exprcode(op, r1, r = r2 + 1, 0, p);
break;
```

don't need result, not effect, push args on stack



p =

## Short-Circuit Evaluation

- evaluate conditional only as far as is necessary to determine truth value.

```
if (a < 10 || b == 6)
    a = 1;
```

semantically equivalent to

```
if (a < 10) goto L3
if (b != 6) goto L1
```

L3: a = 1

L1:L2:

- obtain short-circuit evaluation by correct traversal of trees for conditional expressions.

```
ifstat(lab)
int lab;
{
    ...walk(expr(), o, lab);
    ...
}
```



~~void~~ never called w/ both tlab & flab

- must consider  $\parallel$  && ! comparisons with both "true" and "false" labels.
- $walk(e1 \parallel e2, tlab, 0)$  "fall thru on false"

want:  $\begin{matrix} \text{if}(e1) \text{ goto } tlab \\ \text{if}(e2) \text{ goto } tlab \\ \downarrow \\ \text{"false"} \end{matrix}$  if

use:  $walk(p \rightarrow e\_left, tlab, 0);$   
 $walk(p \rightarrow e\_right, tlab, 0);$

- $walk(e1 \parallel e2, 0, flab)$  "fall thru on true"

want:  $\begin{matrix} \text{if}(e1) \text{ goto } L \\ \text{if}(\sim e2) \text{ goto } flab \\ L: \downarrow \\ \text{"true"} \end{matrix}$  while

use:  $walk(p \rightarrow e\_left, tlab = genlabel(1), 0);$   
 $walk(p \rightarrow e\_right, 0, flab);$   
 $deflabel(tlab);$

- $walk(e_1 \ \&\& \ e_2, \ tlab, \ 0)$  "fall thru on false"

want:

```

if (~e1) goto L
if (e2) goto tlab
L: "false"

```

use:

```

walk(p → e-left, 0, flab = genlabel(1));
walk(p → e-right, tlab, 0);
deflabel(flab);

```

- $walk(e_1 \ \&\& \ e_2, \ 0, \ flab)$  "fall thru on true"

want:

```

if (~e1) goto flab
if (~e2) goto flab
↓
"true"

```

or, and,  
not, never  
separate use:  
to upgrade

```

walk(p → e-left, 0, flab);
walk(p → e-right, 0, flab);

```

- $walk(!e, \ flab, \ flab)$

use:

```

walk(p → e-left, flab, tlab)

```

$f(a) \rightarrow !f(a \neq 0)$   
 code(walk)  $\rightarrow$  not a copy of walk  
 new f(i)  $\rightarrow$  add i to walk  
 @  $\rightarrow$  add i to walk

use:  
 $r1 = \text{walk}(p \rightarrow e\text{-left}, 0, \text{flag});$   
 $r2 = \text{walk}(p \rightarrow e\text{-right}, 0, \text{flag});$   
 $\text{exprcode}(0\text{-E}, r1, r2, \text{flag}, p);$

$\rightarrow$  complement operator

want:  $\text{if}(e1 > e2) \text{ goto flag}$

•  $\text{walk}(e1 < e2, 0, \text{flag})$  "fall thru on true"  
 used in code generator

use:  
 $r1 = \text{walk}(p \rightarrow e\text{-left}, \text{flag}, 0);$   
 $r2 = \text{walk}(p \rightarrow e\text{-right}, \text{flag}, 0);$   
 $\text{exprcode}(0\text{-LT}, r1, r2, \text{flag}, p);$

want:  $\text{if}(e1 < e2) \text{ goto flag}$

•  $\text{walk}(e1 < e2, \text{flag}, 0)$  "fall thru on false"

# Storage Management

Do this  
[cb]!

- must avoid using storage (for nodes) proportional to program size.
- want to reach "steady state".
- maintain free list of available nodes as in csym.c for symbol table entries.
- add nodes (after traversal) to free list:

```
int walk(p, tlab, flab)
struct node *p;
int tlab, flab;
{
```

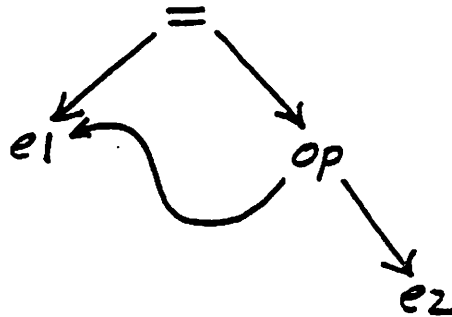
```
    :
    switch (op = p->e-op) {
        :
    }
```

```
    freenode(p->e-left);
    freenode(p->e-right);
    p->e-result = r;
    return (r);
}
```

o-ID, o-con  
need special  
handling  
don't call  
freenode  
for symbol ptrs }

if (op == O\_ID || op == O\_CON) { }

- problem: what about  $e1 \text{ op} = e2$  ?



- must avoid putting  $e1$  on free list twice.
- maintain reference count in  $e\_count$  field;   
  $freenode$  decrements  $e\_count$ , frees node if  $e\_count == 0$ .   
  $pointer$  points to node
- $e\_count$  is number of references to the node; computed during tree construction

bump  
count of  
l & r  
trees

```

freenode(p)
struct node *p;
{
    if (p && --p->e_count == 0) {
        p->e_left = freelist
        freelist = p;
    }
}
  
```

- useful during code generation for resource management.

enode does  
semantic work  
expr does syntax

## Type Checking

- $p = \text{enode}(op, \text{left}, \text{right})$ 
  - check types of left & right
  - allocate node - and others if necessary
  - fill in e-type field
  - return pointer to new node
- enode embodies expression semantics
- may need to add "conversion nodes",  
e.g.  $0\_CVF$ ,  $0\_CVI$
- may need to "convert" lvalues to rvalues.
- may need to "widen" chars to ints.
- need to do "type algebra", e.g. if  $p$  is  $*T$ , then  $*p$  is a  $T$ ; if  $a$  is an int, then  $\&a$  is a  $*T$ ; etc.
- may need to add " $!= 0$ " to expressions to yield conditional expressions.

```

struct node *enode(op, left, right)
int op;
struct node *left, *right;
{
    int ty;    /* type of result */

    ty = T_VOID;
    switch (op) {
        case Q_ID:
            :

        case AND: case OR:
            :

        default:
            fprintf(stderr, "compiler error");
            exit(1);
    }
    return (node(op, ty, left, right))
}

```

where node allocates a struct node and initializes its fields as given; also initializes e\_result (0), e\_count (0); and increments reference counts of left & right nodes (watch out!) id + const

```
struct symbol *q;
```

```
:
```

```
case O_ID:
```

```
q = (struct symbol *) left;
```

```
ty = q->s-type;
```

```
if((q->s.flags & ARRAY) == 0)
```

```
    ty |= T_ADDR; param
```

```
break;
```

"this node is an lvalue"

```
case O_CON:
```

```
q = (struct symbol *) left;
```

```
ty = q->s-type;
```

- O\_ID & O\_CON could be handled directly by calling node in term(), since there is no type checking.
- same for O\_ARG; others that do not require checking, such as 's'.



case QNEG:

```

    left = widen(left);
    ty = left → e-type;
    if (ty != T_INT && ty != T_FLOAT)
        err("type conflict for unary -", "");
    break;
    :

```

→ = T\_INT; →

```
struct node *widen(p)
```

```
struct node *p;
```

```
{
```

```
    p = rvalue(p);
```

```
    if (p → e-type == T_CHAR)
```

```
        p = node(Q_CVI, T_INT, p, NULL);
```

```
    return (p);
```

```
}
```

```
struct node *rvalue(p)
```

```
struct node *p;
```

```
{
```

```
    if (ISADDR(p → e-type))
```

```
        p = node(Q_INDIR, p → e-type & ~T_ADDR,
```

```
                p, NULL);
```

```
    return (p);
```

```
}
```

Case O\_ADD: Case O\_SUB: Case O\_MUL: Case O\_DIV:

```

left = widen(left);
right = widen(right);
tl = left->e-type;
tr = right->e-type;
if (tl == T_INT && tr == T_INT) {
    ty = T_INT;
    break;
}

```

hint: make this mess a function!

```

if (tl == T_FLOAT && tr == T_FLOAT) {
    ty = T_FLOAT;
    break;
}
if (tl == T_INT && tr == T_FLOAT) {
    ty = T_FLOAT;
    left = node(O_CVF, T_FLOAT, left,
                NULL);
    break;
}

```

```

if (tl == T_FLOAT && tr == T_INT) {
    ty = T_FLOAT;
    right = ...
    break;
}

```

p++ → p+1  
tr



```
if (ISPTR(tl) && tr == T_INT &&
    (op == O_ADD || op == O_SUB)) {
```

```
    ty = tl;
    break;
```

```
err("type conflict for ", opnames[op]);
```

```
    ty = T_INT;
    break;
```

→ avoids further errors.

case O\_LSH: case O\_RSH: case BXOR:

case O BOR: case O\_BAND: case O\_MOD:

```
left = widen(left);
```

```
right = widen(right);
```

```
if (left->e-type != T_INT ||
    right->e-type != T_INT)
```

```
    err("type conflict for ", opnames[op],
```

```
    ty = T_INT;
```

```
    break;
```

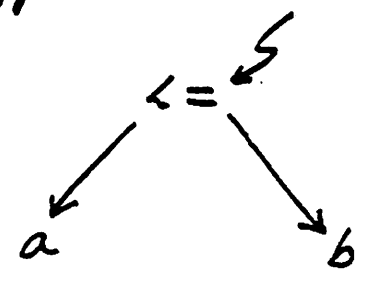
case O\_BCOM:

as above, omitting right tree



# Conditional Expressions

- have "type" T\_COND

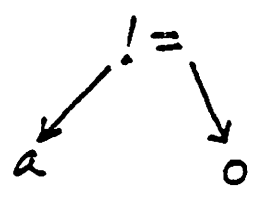


- "conversion" of expression to T\_COND involves adding comparison with 0, with associated type-checking => use enode.

if (a) ...

is equivalent to

if (a != 0) ...



```

struct node *cnode(p)
struct node *p;
{

```

```

  if (p->e_type != T_COND)

```

```

    p = enode(O_NE, p, node(O_CON,
      T_INT, constant("0", T_INT),
      NULL));

```

use enode  
to get type  
checking.

```

  return (p);
}

```

enode cont'd

case AND: case OR:

left = cnode(left);

right = cnode(right);

ty = T\_COND;

break;

case O\_LT: case O\_LE: case O\_EQ:

case O\_NE: case O\_GE: case O\_GT:

ty = T\_COND;

same as O\_ADD, etc.

for T\_INT, T\_FLOAT,

but don't set ty

? { if (ISPTR(tl) && tr == T\_INT &&  
(op == O\_EQ || op == O\_NE))

break;

if (ISPTR(tl) && tl == tr)

break;

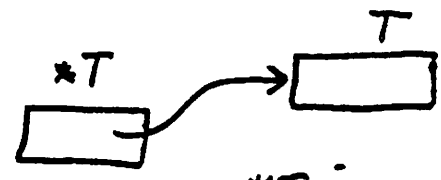
err("type conflict for ", opnames[op]);

break;

enode cont'd

Type Algebra:

case O\_INDIR:



\*T is an lvalue.

```

left = rvalue(left);
if ((n = ISPTR(left->e-type)) == 0)
    err("type conflict for unary *", "");
ty = MKTYPE(n-1, BASETYPE(left->e-type));
left->e-type = ty | T_ADDR;
return (left);

```

- Note: node for O\_INDIR constructed by rvalue, if T\_ADDR set.
- assume you want lvalue, let rvalue handle other case.

case UNARY '&':

{ define UNARY to be something to distinguish operators, e.g. 010001 }

```

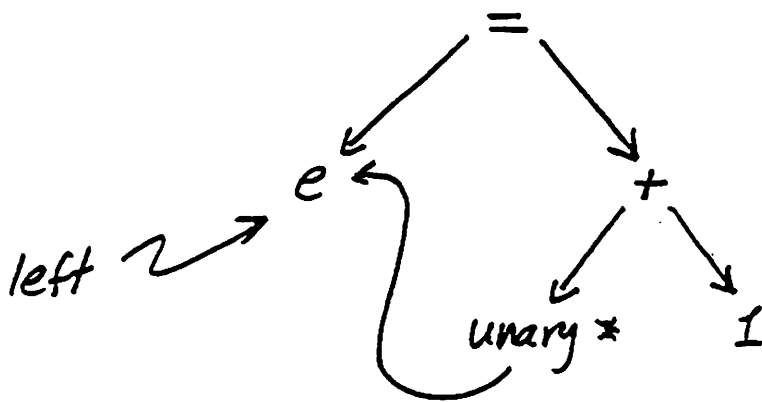
if (!ISADDR(left->e-type))
    err("type conflict for unary &", "");
ty = MKPTR MKPTR(1 + ISPTR(left->e-type),
    left->e-type & ~ T_ADDR);
left->e-type = ty;
return (left);

```

only "real" effect of unary & !!

$$++e \equiv e += 1$$

(see previous slide #49)



(likewise for --e)

CASE UNARY INCR:

```
left = lvalue (left);
```

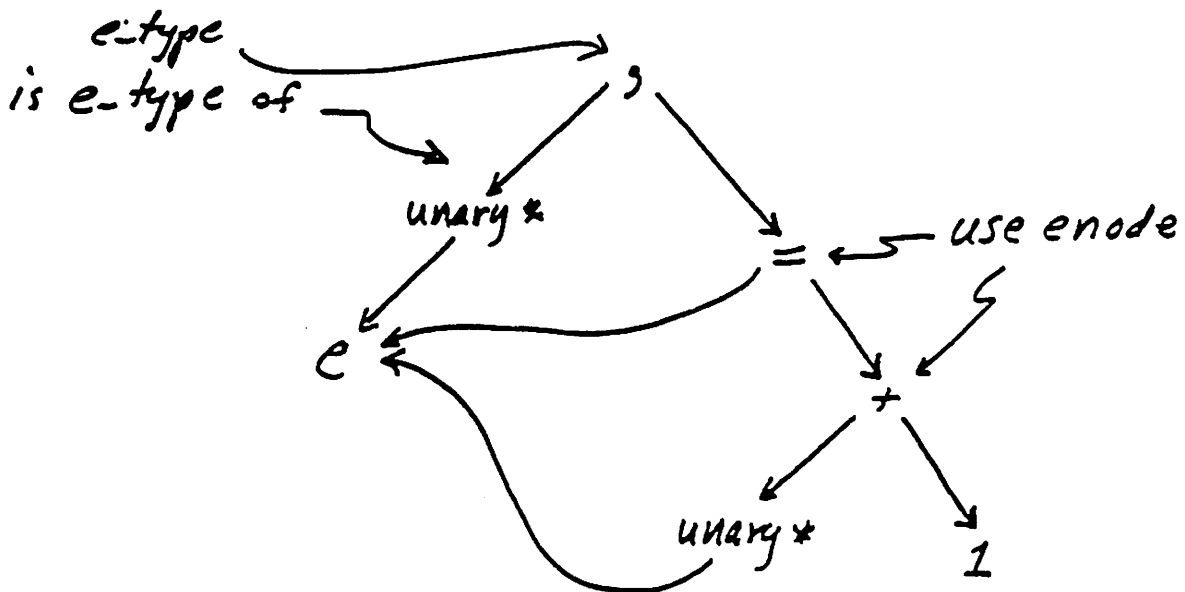
```
right = node(O_CON, T_INT,
             constant("1", T_INT), NULL);
```

enode adds unary \* ↗

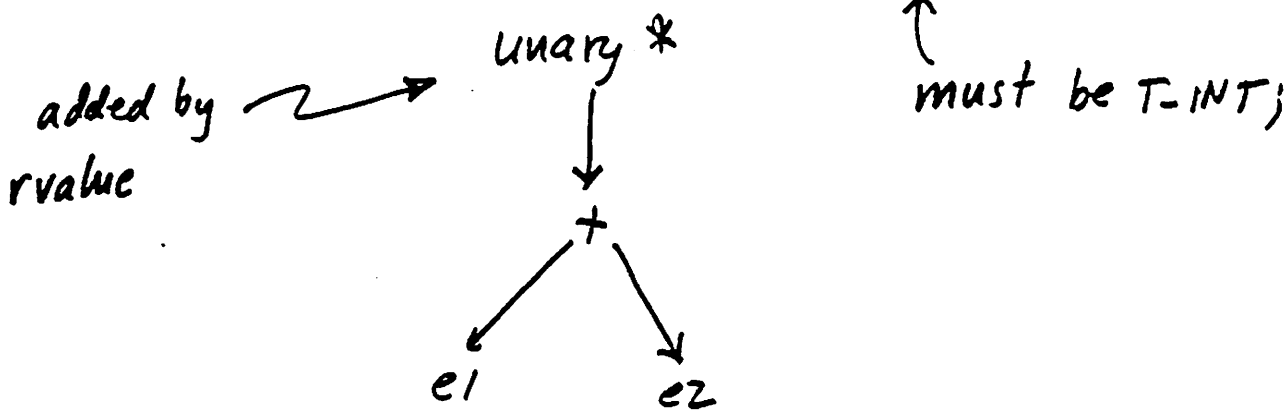
```
right = enode(O_ADD, left, right);
```

```
return(enode(O_ASSIGN, left, right));
```

$$e++ \equiv (t' = e, e += 1, t')$$



$$e1[e2] \equiv *(e1 + e2)$$



case '[':

left = rvalue(left);

right = widen(right);

if (right->e\_type != T\_INT)

err("type conflict for subscript", "");

if ((n = ISPTR(left->e\_type)) == 0)

err("type conflict for array name", "");

ty = MKTYPE(n-1, BASETYPE(left->e\_type));

return (node(O\_ADD, ty | T\_ADDR, left, right));

↑  
don't need enode because  
dirty work is done.

like  
Unary \*



$$\underline{e_1 = e_2}$$

lvalue T = T

lvalue int = float (convert to int)

lvalue float = int (convert to float)

lvalue (ptr T) = int

↑ result type

case O\_ASSIGN:

left = lvalue(left);

right = widen(right)

ty = left → e\_type & ~T\_ADDR;

if (ty == T\_CHAR && right → e\_type == T\_INT)

;

else if (ty == T\_INT && right → e\_type == T\_FLOAT)

right = node(O\_CVI, T\_INT, right, NULL);

else if (ty == T\_FLOAT && right → e\_type == T\_INT)

right = node(O\_CVF, T\_FLOAT, right, NULL);

else if (ISPTR(ty) && right → e\_type == T\_INT)

;

else if (ty != right → e\_type)

err("type conflict for = ", "");

break;

Only call  
O\_ASSIGN  
for  
source ref  
of "\*" & "&".

$$f(\dots)$$

$$f(a, b, c)$$

Case O\_CALL:

```

if (!ISFUNC(left → e-type)) {
    err("function expected", "");
    semantic error recovery → left = node(O-ID, T-INT,
                                         lookup("?" ), NULL);
}
ty = left → e-type & ~ (T-FUNC | T-ADDR);
break;

```

- enter "?" in symbol table during initialization
- use "?" whenever an ID is needed but you don't have one!
- give code generator well-formed, but maybe semantically incorrect, trees.

Computer Science 453

Fall 1983

Class Notes, Set 3

September 19, 1983

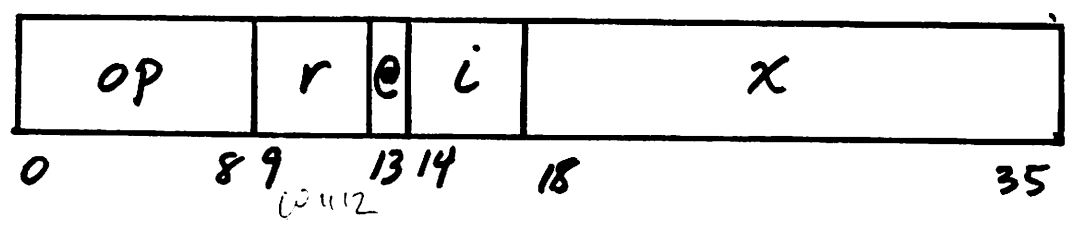
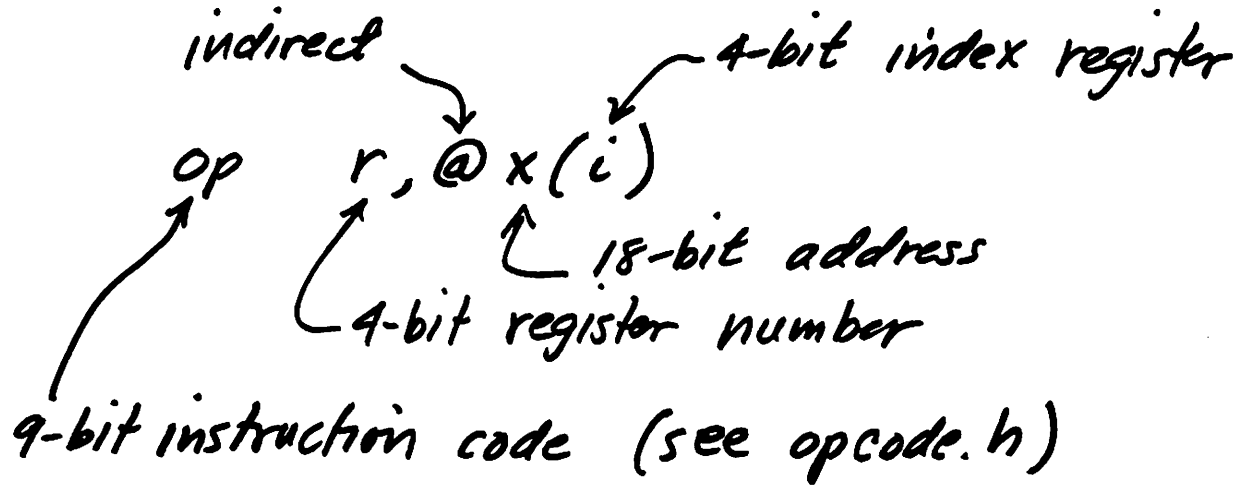
## Code Generation

- replace bottom half routines by set for DEC-10
- generate link object code
- generate poor, but correct code
- manage machine resources, e.g. registers
- design run-time representation for programs, e.g. memory allocation for code, data; activation record layout for procedures.

## Machine Model

- 2-address machine w/16 registers
- operators leave results in registers
- $2^{18}$  words (36-bit) of memory;  
0-15 address the registers
- three logical segments
  - "code" - executable program code
  - "data" - global data
  - "lit" - constants
- code + lit are loaded into read-only memory.
- stack (in the data segment) for procedure activation records ("frames").
- fixed-width (36-bits) instructions

# DEC-10 Instruction Set

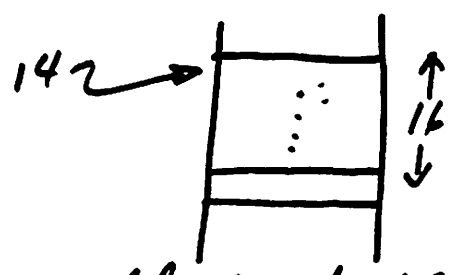


Examples: (assembly language) [see opcode.h]

```

move 2, x ; load value of x
add 6, 7 ; add register 7 to 6
movem 13, 16(14) ; store register 13,

```



```

addi 6, 40 ; add 40 to reg. 6

```

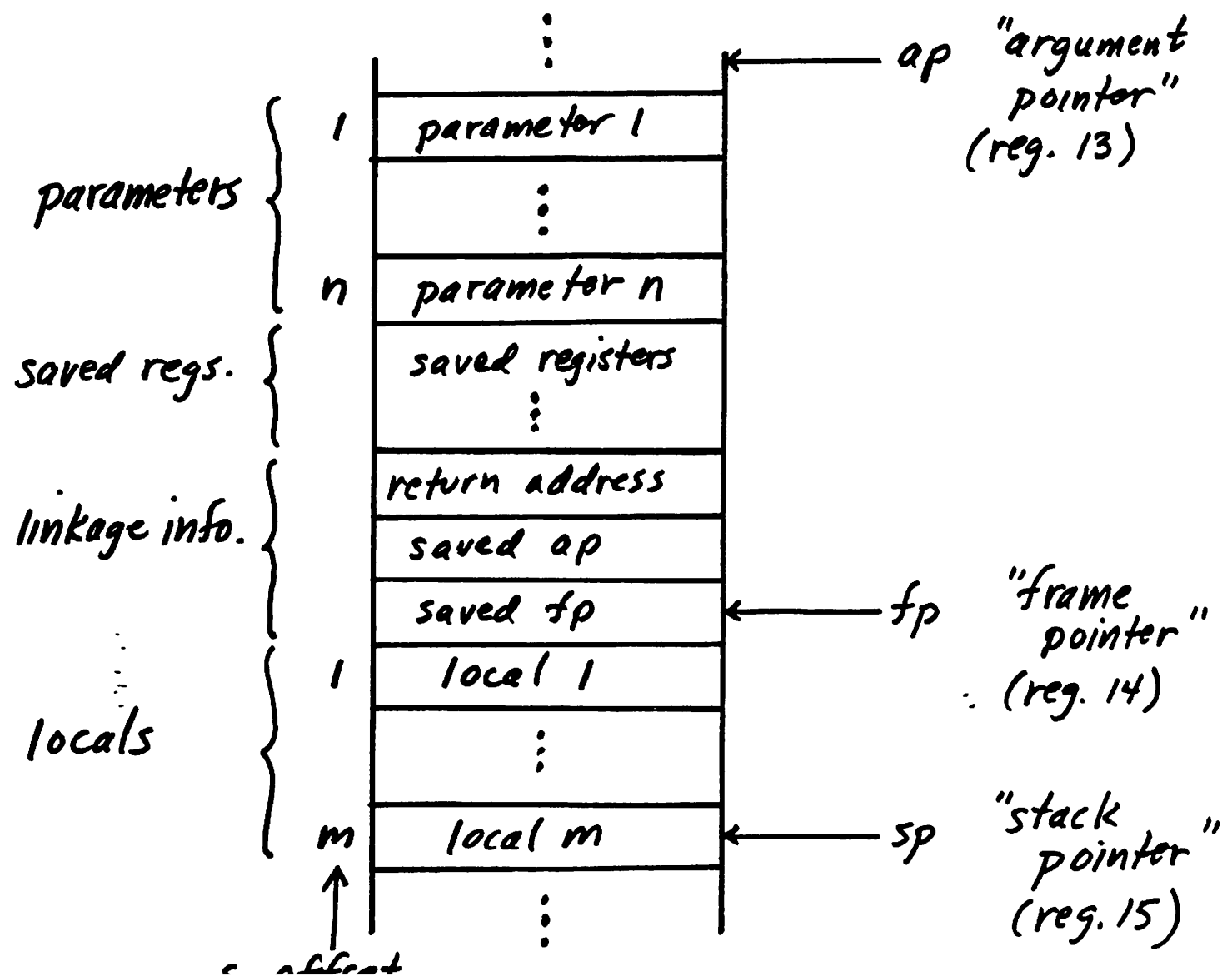
```

move (17) ; load ...

```

# Activation Records

- allocated & initialized by calling sequence & entry code
- deallocated by exit code
- contains parameters, locals, linkage info, saved registers



# Object Code

• commands & "object text"

• commands:

- .def id expr
- .seg id
- .len id n

• expressions: polish suffix

- expr :
- id
  - con
  - expr expr op

op : + - < >

~  
shifts

Example: movei 6, rows      "load address of rows into r6."

op      r      x

0201 27 < 6 23 < + rows +

op      r      x

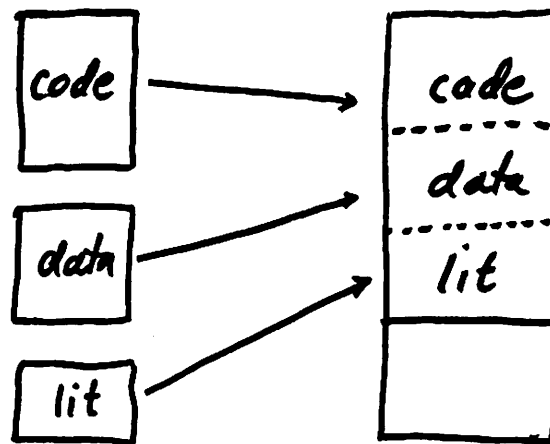
op << 27  
↑  
unsigned  
output two  
16 bit #<sub>2</sub>

020130000000 rows +

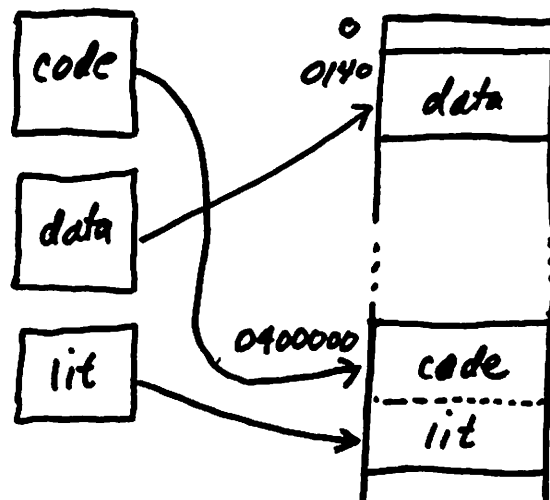


## Segments

- code : executable program code (read-only)
- data : global data & the stack
- lit : constants (read-only)
- map 3 logical segments onto target machine:



CDC Cyber:  
single address  
space, not  
protected



DEC-10 :

"low" segment

"high" segment  
(protected)

(actually, DEC-10  
is paged machine)

## Segments cont'd

- during execution, segments are contiguous
- during compilation, segments are interleaved
- link converts interleaved to contiguous segments
- during compilation, use `.seg` commands to switch segments, e.g.

```

int a;
int g()
{
    a = g("hi there");
    :
}

```

- `.seg data`  
(definition for `a`)
- `.seg code`  
(code for `g`)
- `.seg lit`  
(constant "hi there")
- `.seg code`  
(more code for `g`)



## Example: Local labels

```
for(i = 1; i < 10; i++) {
    sum += x[i];
    if (sum > 100)
        break;
}
```

---

```
i = 1
L1: if (i >= 10) go to L3
    sum += x[i];
    if (sum <= 100) go to L4
    go to L3
L4:L5:L2: i++
          go to L1
L3:
```

---

```
.seg code
i = 1
.def $1 code 5 +
if (i >= 10) go to $3
:
.def $4 code 25 +
.def $5 code 25 +
.def $2 code 25 +
:
025400000000 $1 +      (go to L1)
def $2 code 25 +
```

## Local Labels cont'd

- `def -s $n code loc +`  
     ↑  
     indicates "local" label  
     (deleted from link output)

• must maintain "location counter" for each segment, which is used in definitions.

- location counter incremented for each word of object text emitted or allocated.

Curseg Curloc

```
deflabel(lab)
```

```
int lab;
```

```
{
```

```
output ".def -s $lab name loc +"
```

```
...  
}
```

value of "lab"

name of current  
segment

location  
counter for  
current segment

# Segment Lengths

- at end of object file, issue .len commands giving length (in words) of each segment : e.g.

```
.seg data
{
```

```
.seg code
{
```

```
.seg lit
{
```

```
.seg code
{
```

total # words emitted or allocated

```
.len data 44
```

```
.len code 136
```

```
.len lit 7
```

progend(x)

```
int x;
```

```
{
```

```
for each segment x
```

```
if (x > 0 length)
```

```
output ".len x n"
```

length of x

```
}
```

## Data Segment

- contains definitions only - nothing to emit ("uninitialized data")
- len command must take allocation into account, e.g.

```
int a;
float x[10];
char s[8];
```

```
.seg data
.def a data 0 +
.def x data 1 +
.def s data 11 +
.len data 13
```

```
global(p)
struct symbol *p;
{
    switch to data segment;
    output .def command;
    increment data location counter
    according to p → s_size;
}
if necessary!
```

## Bottom Half Routines

defconst(p)	define a constant
✓ deflabel(lab)	define a label
exprcode(op, r1, r2, lab, p)	expression operator
✓ global(p)	define a global
Jump(lab)	unconditional jump
procbeg(p)	begin procedure
* procend(x)	end procedure
* progbeg()	begin program
✓ progend(x)	end program
retcode(r)	return statement
* rfree(r)	release resource r
✓ - done!	
* - no output; perhaps initialization or other processing	



## Jumps

- unconditional jump:

jrst label

- link object code:

025400000000 label +

jump(lab)

int lab;

{

emit(JRST, 0, NULL, lab, 0);

}

DEC-10 Insts. only

- hint: write emit(op, r, name, off, i) to emit

op r, name ± off (i)

[ permit 0 args for omitted parts ]

in link object text, and increment appropriate location counter.

- use -d option to print expressions in symbolic (assembly language) form

## Procedures

- `procbeg` must
  - initialize code generator
  - generate procedure entry sequence
  - assign frame offsets for parameters and locals.
- procedure entry sequence:

for tracing	{	pointer to "f" in lit segment
entry sequence	{	f: jsp 4, c.savr
		adjsp 15, m
		↑ size (in words)
		of local portion
		of the frame

`emit(jsp, 4, c.savr, 0, 0)`  
`emit(adjsp, 15, 0, 0, m)`

- argument pushed on the stack in calling sequence (return address, too).
- linkage information saved, `r13+r14` set by `c.savr`.

## Procedures cont'd

- assume code location counter is 30,  
lit location counter is 6:

033110000000 G + lit +

.seg code ✓  
0331100000006 lit +

.seg lit  
014600000000 ("f" in 9-bit bytes)

.seg code

ep<sup>out 8</sup> → .def f code 31 +  
0265<sup>0100</sup>200000000 c.savr +  
0105<sup>1111</sup>74000010  
size of locals

- assign offsets 1, 2, ... n to parameters, set s\_offset field. [Done by top half!]
- assign offsets beginning at 1 to :locals, increasing by size of each local, in any order.

## Procedures cont'd

- `retcode` generates return sequence:

```

    move l, r    (if r ≠ 0)
    jsp 4, c.ret
  
```

link code:

```

    0200040000000 r +
    0265200000000 c.retr +
  
```

- omit "move l, r" if  $r == 0$ , e.g.  
`retcode(0)`

```
retcode(r)
```

```
int r;
```

```
{
```

```
    if (r)
```

```
        emit(MOVE, l, NULL, r, 0);
```

```
    emit(JSP, 4, "c.retr", 0, 0)
```

```
}
```

## Expressions (exprcode)

- manage resources, i.e. registers
- leave result of operation in a register
- determine type-specific operation (e.g. int add, float add) given generic operation (O-ADD) and operand types.
- generate correct (but awful!) code

$a = b + c$  (all int)

yields:

	movei	6, a	; load a's lvalue
	movei	7, b	; load b's lvalue
	move	8, (7)	; load b's rvalue
note reuse of reg. 7 →	movei	7, c	; load c's lvalue
	move	9, (7)	; load c's rvalue
	add	8, 9	; a + b + c
	movem	8, (6)	; store in a

could be:

	move	6, b	; b's rvalue
	add	6, c	; add c's rvalue

```
int exprcode(op, r1, r2, lab, p)
int op, r1, r2, lab;
struct node *p;
{
    int r;
    struct symbol *q;

    r = 0;
    switch (op) {
        :
        case O_ADD:
            (set r)
            break;
        :
        default:
            fprintf(stderr, "compiler error\n");
            exit(1);
    }
    return (r);
}
```

## Register Usage

<u>register</u>	<u>use</u>
0	not used
1	function return values
2-5	entry/exit sequence
6-12	saved across calls
13	parameter pointer
14	frame pointer
15	stack pointer

- use 6-12 for evaluating expressions
- if run out of registers, issue error "expression too complicated"
- \* • what if you only had 3 registers? temporaries?
- 13-15 are initialized by "start-up code" prior to calling main; (stack is also allocated).

## Constants

- reference to a constant is to its rvalue:

```

move r, $L
      ↙ store L in
        s-offset field,
        allocate if
        necessary
    
```

output by  
defconst

```

{
  .seg lit
  .def -s $L lit 1 +
  constant
}
    
```

case O.CON:

```
q = (struct symbol *) p → e-left;
```

```
if (q → s-offset == 0)
```

```

    ↙ will change
      labels in
      typval
    q → s-offset = genlabel(1);
    
```

```
emit(MOVE, r = ralloc(1), NULL,
     q → s-offset, 0);
```

```
break;
```

"just better execute"

2, 3, 5, 6

⋮

- `ralloc(1)` - allocates 1 register in 6-12. (useful to allocate lowest # first).



## "Immediate Mode"

"address is value"

- immediate mode: address is the value instead of the address of the value.

```

move 6, $10
:
.seg lit
.def -s $10 lit 1 +
45

```

} load 45  
into reg. 6

also done by

```
movei 6, 45
```

- used to load addresses & small constants ( $\leq 18$  bits).

case 0.CON:

```

q = (struct symbol *) p -> e_left;
if ((q -> s_type == T_CHAR ||
     q -> s_type == T_INT) &&
     (n = atoi(q -> s_name)) <= 0777777)
    emit(MOVEI, r = ralloc(i), NULL,
         n, 0);

```

else {

...

## Lvalues / Addresses

- lvalues are addresses of variables
- lvalue representation must be coordinated with O\_INDIR & O\_ASSIGN.
- simple case: 18-bit address
- hard case: byte pointer (for character strings)
- simple lvalues: depend on scope
  - global: name (defined elsewhere)  
by global
  - param: s\_offset(13)
  - local: s\_offset(14)

case O\_ID:

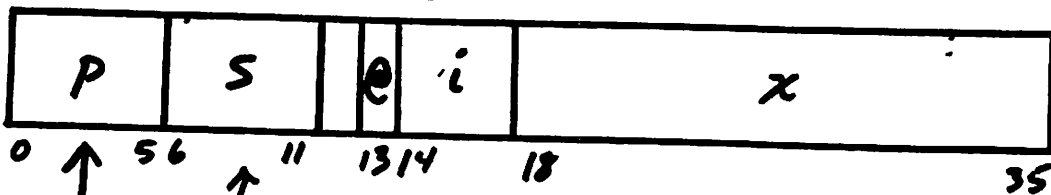
```

r = ralloc(1);
q = (struct symbol *) p->e_left;
if(q->s_flags & GLOBAL)
    emit(MOVEI, r, q->s_name, 0, 0);
else if(q->s_flags & PARAM)
    emit(MOVEI, r, NULL, q->s_offset, 13);
else
    emit(MOVEI, r, NULL, q->s_offset, 14);
    handle byte pointers
break;

```

Byte pointers - interpreted by hardware

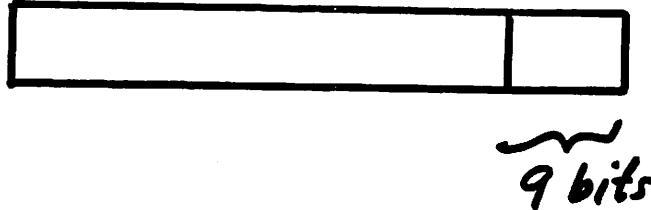
- 36-bit pointer specifying byte size, byte position, address



↑ byte size (9 bits for characters)  
 ↑ number of bits to the right of byte.

## Byte Lvalues

- identifier is a T-CHAR:



must treat  $s[i]$  and  $c$  (both T-CHAR) in same fashion; return byte pointer for both.

return byte pointer to rightmost 9 bits for characters:

in previous box { if ( $g \rightarrow s.type == T\_CHAR$ )  
                   emit(HRLI, r, NULL, 0001100, 0);  
                                   ↓ hrti    6, 0001100    ↑ octal

constructs:



(HRLI = half word, right, to left, immediate!)

- more sophisticated code generator would avoid use of byte pointers here.

## Byte Lvalues

- identifier is an array name: want



i.e. pointer to first byte.

- do not set T-ADDR on O-ID nodes for array names (GLOBAL & LOCAL);  
then  $p \rightarrow e\text{-type} == \text{MKTYPE}(1, \text{T-CHAR})$

in previous  
box

```

      ⋮
    else if (p → e-type == MKTYPE(1, T-CHAR))
      emit(HRLI, v, NULL, 0331100, 0);
  
```

```

    break;
  
```

- note type of result is T-ADDR/T-CHAR because it's an lvalue; use in O-INDIR & O-ASEN.

## Loading Rvalues

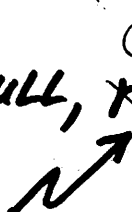
- for characters, use "load byte" instruction; for others, use "load word" instruction.

case O\_INDIR:

```

if ((p->e->left->e->type & ~T_ADDR) == T_CHAR)
    emit(LDB, r=ralloc(1), NULL, r1, 0);
else
    emit(MOVE, r=ralloc(1), NULL, r1, 0);
break;

```


  
 register w/ value  
 (byte pointer or  
 address)


- O\_ASGN is similar:

case O\_ASGN:

```

if ((p->e->left->e->type & ~T_ADDR) == T_CHAR)
    emit(DPB, r=r2, NULL, r1, 0);
else
    emit(MOVM, r=r2, NULL, r1, 0);
break;

```


  
 register w/ righthand side rvalue

defconst

- called at end of compilation with each constant; to be placed in lit segment.
- if referenced, will have non-zero s\_offset field (see p. 21).
- output label and constant value:

T-INT, T-CHAR: print the value, e.g.

```
.def -s $16 lit 1 +
45
```

T-FLOAT: print bit pattern given by atof, e.g.

```
.def -s $16 lit 1 +
0206706314631 ← 56.8
```

MKTYPE(1, T-CHAR): print byte pointer followed by string in 9-bit bytes, e.g. "hello"

byte ptr as value

```
.def -s $16 lit 1 +
0|331|100|000|002 lit +
0150145154154 ← h e l l
0157000000000
```

byte pointer to next word →

next word ↙

↖ null byte

don't forget null byte; e.g. "four" takes 3 words

## Addition

- possible operand types are

ptr T + int

int + int

float + float

- special case ptr char + int;  
treat ptr T + int like int + int.  
(i.e. addresses are added w/ normal  
addition except byte pointers)

CASE 0-ADD:

```
if (p->e->left->e->type == MKTYPE(1, T_CHAR))
```

```
    emit(ADJBP, r=r2, NULL, r1, 0);
```

```
else if (p->e->left->e->type == T_FLOAT)
```

```
    emit(FADR, r=r1, NULL, r2, 0);
```

```
else
```

```
    emit(ADD, r=r1, NULL, r2, 0);
```

```
break;
```

table indexed by  
0-values

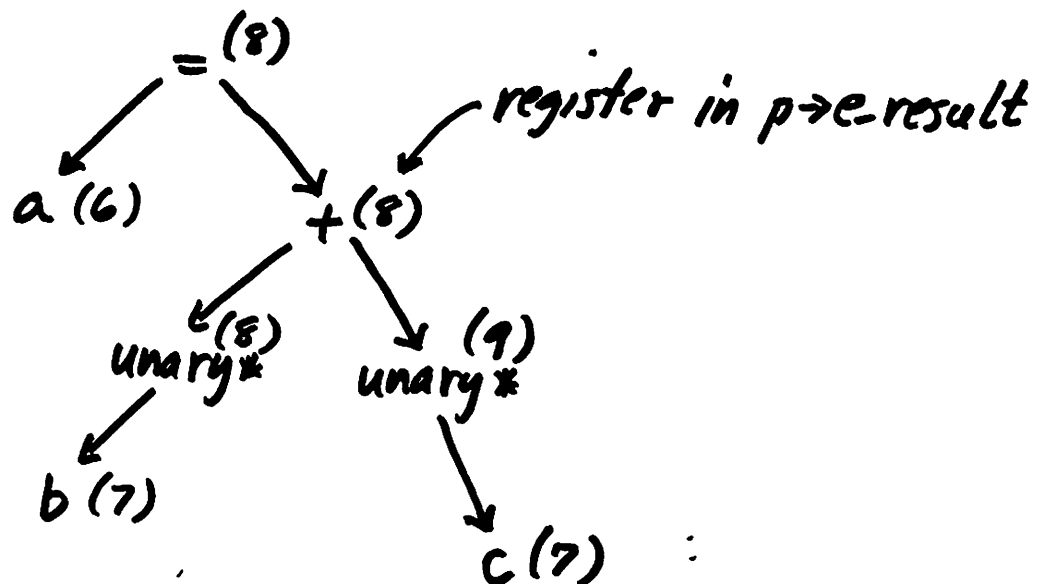
- hint: this code can handle several operators

4  
XOR  
AND  
Cant do  
shifts



## Register Management

- allocate registers (6..12) via `ralloc(i)`
- free registers via `rfree(r)`, called when node is freed during traversal by top half.
- maintain reference count in code generator indicating number of uses of each register. Be careful!
- Example:  $a = b + c$



- note multiple occurrences of register 8, which arise from "copying" register in `a-ADD` & `o-ASGN` code (e.g. `r = r1` in `a-ADD`).

## Register Management, cont'd

- "copying" is equivalent to "allocating"
- reference count must be incremented when register is copied.
- `ralloc(i)`: finds a register `r` with `rcount[r] == 0`, sets `rcount[r]` to 1, and returns `r`. [issue error if none available!]  
 Think about this →
- "copying" `r`: increment `rcount[r]`
- `rfree(r)`: decrement `rcount[r]`

case Q.ADD: (copies `r1` or `r2`)

⋮

else

`emit(ADD, r = r1, NULL, r2, 0);`

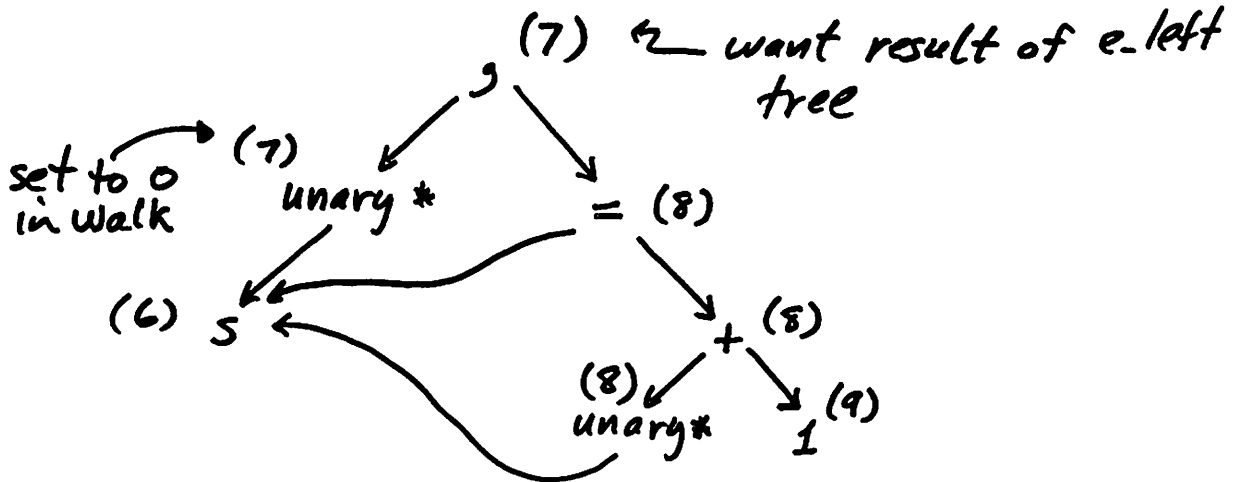
→ `rcount[r]++;`

`break;`

- also need `rcount[r]++` in Q.ASSGN.
- `rfree` called almost exclusively by top half.

## Register Management, cont'd

- problem: root node of  $s++$



- ";" causes a "copy" of register 7, but `exprcode` is not called, so reference count is wrong!
- solution: "move" register 7 to ";" node  
 in walk: could call `exprcode` directly  
 case ';': Better code requires ~~2~~ to copy `regs`

```

r = walk(p->e-left, 0, 0);
walk(p->e-right, 0, 0);
p->e-left->e-result = 0;
... break;
}
p->e-result = r;
return(r);

```

unary '\*'

### freenode

```
freenode(p)
struct node *p;
```

left of p could be  
O-ED or O-CON

```
{
```

```
if (p && --p->e_count == 0) {
```

```
    if (p->e_op != O_ARG && p->e_result)
        rfree(p->e_result);
```

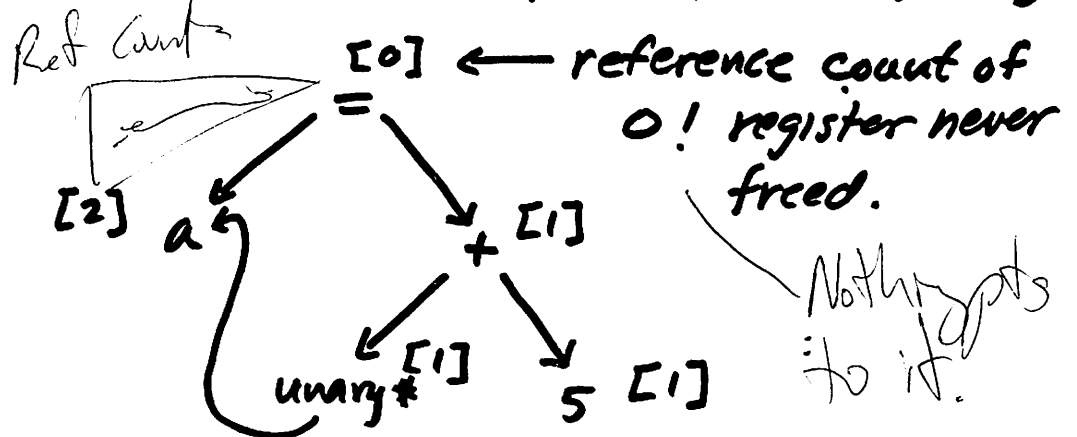
```
    p->e_left = freelist;
```

```
    freelist = p;
```

```
}
```

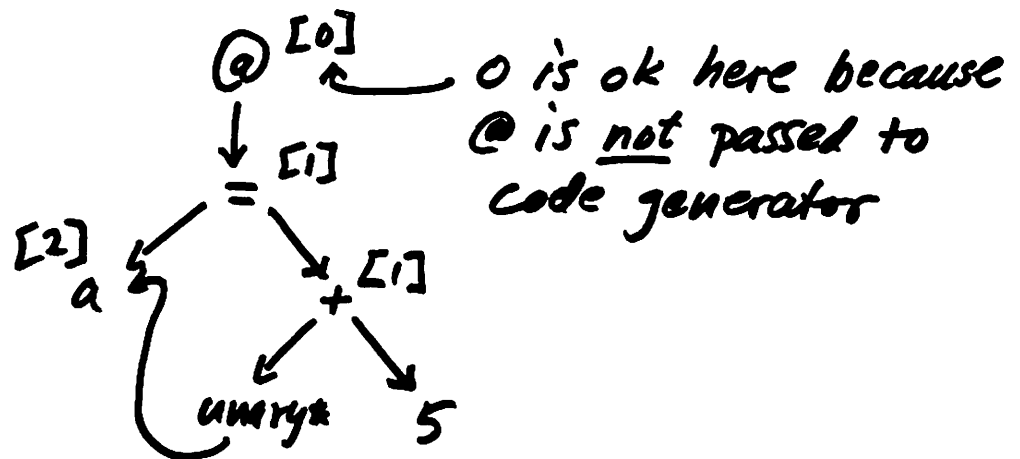
```
}
```

• problem: what about root node: a += 5



## freemove, cont'd

- solution: add another operator at the root of the tree so reference count is 1
- used only for tree traversal to handle register deallocation properly



- in walk:

case '@':

```
r = walk(p->e-left, tlab, flab);
freemove(p->e-left);
```

p->e-result  
is 0.

```
freemove(p);
return(r);
```

why? just netcode?

- add @ node to result of call to expr() (but in right places!)

## Division / Modulus

- integer division (and modulus) requires 2 adjacent registers.

`ldiv 6, 8`

leaves quotient in register 6  
and remainder in register 7.

- if `rcount[r1+1] == 0`, fine otherwise must move `r1`
- `ralloc(2)`: allocate 2 adjacent registers, return first one.
- free extra register, if allocated.
- `0.MOD` is similar.

case O\_DIV:

```

if (p->e->left->e->type == T_FLOAT) {
    emit(FDVR, r=r1, NULL, r2, 0);
    rcount[r]++;
    break;
}

```

*ravail*

```

if (rcount[r1+1] == 0) {
    emit(IDIV, r=r1, NULL, r2, 0);
    rcount[r]++;
}

```

*copy*

```

else {
    emit(MOVE, r=ralloc(2), NULL, r1, 0);
    emit(IDIV, r, NULL, r2, 0);
    rfree(r+1);
}

```

do NOT  
free r1!  
why?



break;

count will be -1  
as it is later freed

left operand is  
in r1

## Comparisons

- use "compare & skip" instructions for all but char \* comparisons.

e.g.  $r1 < r2$ , use

$\begin{array}{l} \text{cmge } r1, r2 \\ \text{jrst } \$lab \end{array} \left. \vphantom{\begin{array}{l} \text{cmge } r1, r2 \\ \text{jrst } \$lab \end{array}} \right\} \begin{array}{l} \text{"skip" if} \\ r1 \geq r2 \end{array}$

- call library routine to compare char \*'s, "jump" on result.

e.g.  $r1 < r2$ , use

```

push 15, r1
push 15, r2
pushj 15, c.scp
adjsp 15, -2      (remove args)
jump 0, $lab

```

$\uparrow$   
 result ( $<0, =0, >0$ ) in reg. 0.



17

Van com.

Derf Trappel

Computer Science 453

Fall 1983

Class Notes, Set 4

October 6, 1983

# Linking

- "standard" translation process:
  - compile source programs,
  - link object programs,
  - load executable programs.

## - linkers:

- combine separately compiled object modules
  - "search" libraries
  - machine dependent
  - idiosyncratic syntax (e.g. DEC-10 .rel files)
  - archaic restrictions (e.g. ordering, naming)  
often influence language design!  
(e.g. 6-character names).
  - cumbersome i/o (binary)
- solution: get a new linker!

## link

- machine independent linker
- language independent
- regular syntax, semi-readable
- no ordering or naming restrictions (e.g. long names permitted).
- features to simplify code generation (e.g. arbitrary # of segments)
- centralizes compiler's "last pass" functions
- does not include "exotic" features
  - overlays
  - report generation
  - "patching"
- Goal: elevate linkers to useful tools, functional standardization of linkers

## Linking & Loading

- object code :
  - output of compilers & assemblers
  - incompletely bound addresses  
e.g. references to external symbols,  
forward references
  - ".rel" files on the DEC-10  
".o" files on UNIX  
"TEXT" files on IBM under CMS
  - often includes "relocation" information  
indicated by flags pointing to locations  
that must be adjusted upon loading.
- executable code:
  - output of linkers
  - (almost) completely bound addresses,  
depends on operating system
  - input to operating system "execute"  
system call.
  - machine dependent
  - ".exe" files on the DEC-10  
"load modules" on the IBM

## Linking & Loading, cont'd

- linking: "linking loader", "relocating loader"  
 object code  $\rightarrow$  executable code  
 e.g. link-10 on the DEC-10
- linking: "link editors", "binder"  
 object code  $\rightarrow$  object code  
 e.g. ld on UNIX, IBM "linkage editor"
- link: a link editor  
 link code  $\rightarrow$  link code
- translation process:
  - compilation: source code  $\rightarrow$  link code
  - linking: link code  $\rightarrow$  link code ...
  - loading: link code  $\rightarrow$  executable code

$\uparrow$  machine dependent (but small)  
 $\uparrow$  machine independent (contents of link code is machine dependent; functions of link are machine independent)

## Link Code

- lines of text (could use binary format).
- interleaved commands and object text.
- syntax:

file : { cmd | expr }

cmd : .def [-s] id expr

| .seg id

| .org con

| .len id con

expr : expr expr ( + | - | < | > )

| id

| con

init char?  
↓

- identifiers: letters, digits, & . ? -
- constants: octal (leading 0), decimal

## Expressions

- Polish postfix - simplifies evaluation
- link makes no architectural assumptions  
e.g. word size, addressability
- expressions represent 1 "interesting cell"  
as determined by compiler & loader;  
i.e. "interesting cell" is machine-dependent.
- typically, "interesting cell" contains  
relocatable address, e.g.

words! →

<u>machine</u>	<u>cell size</u>	<u>address size</u>	<u>word size</u>
DEC-10	18 bits	18 bits	36 bits
PDP-11	16	16	16
Cyber/175	15	18	60
IBM/370	16	24	32
VAX	16	16-32	32

- loaders combine cells into words by  
addition to permit overflow of  
cell boundaries.





## Expressions, cont'd

- cell size determined only by convenience
- link is independent of cell size, but cell size must be fixed for all object files.
- link simply replaces symbols by their values, and evaluates expressions.
- goal: get rid of symbols! (except segment names).

## Segments

- blocks of object text to be loaded into a contiguous region of memory
- segment names & meanings coordinated among compilers & loaders; link itself is independent of the "meaning" of segment names.

• input to link: any number of segments interleaved in any way

• output: segments output as contiguous blocks of object text.

• inter- and intra-segment references made by referring to base location of a segment, which is simply the segment name.

## Segments cont'd

- Suppose compiler emits code for 4 segments:

code - executable program code

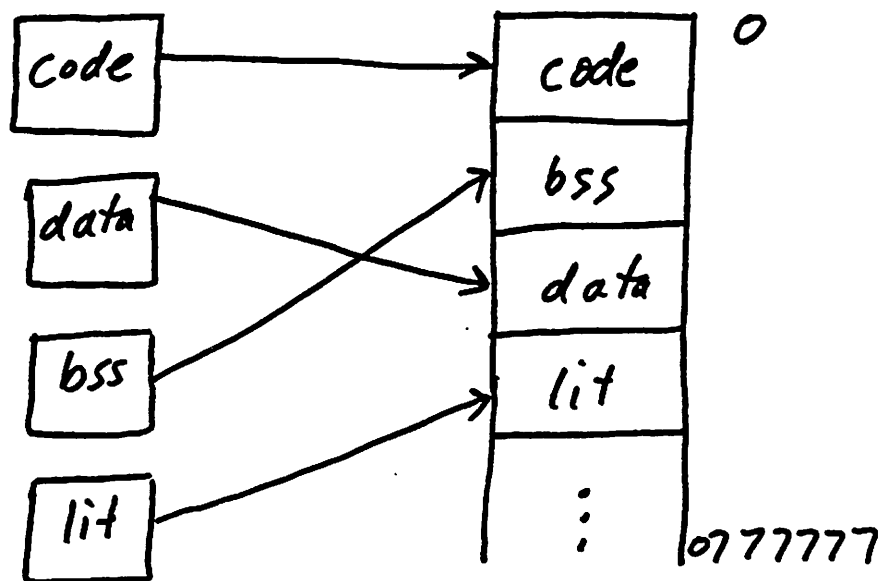
data - initialized data

bss - uninitialized data

lit - literals (constants)

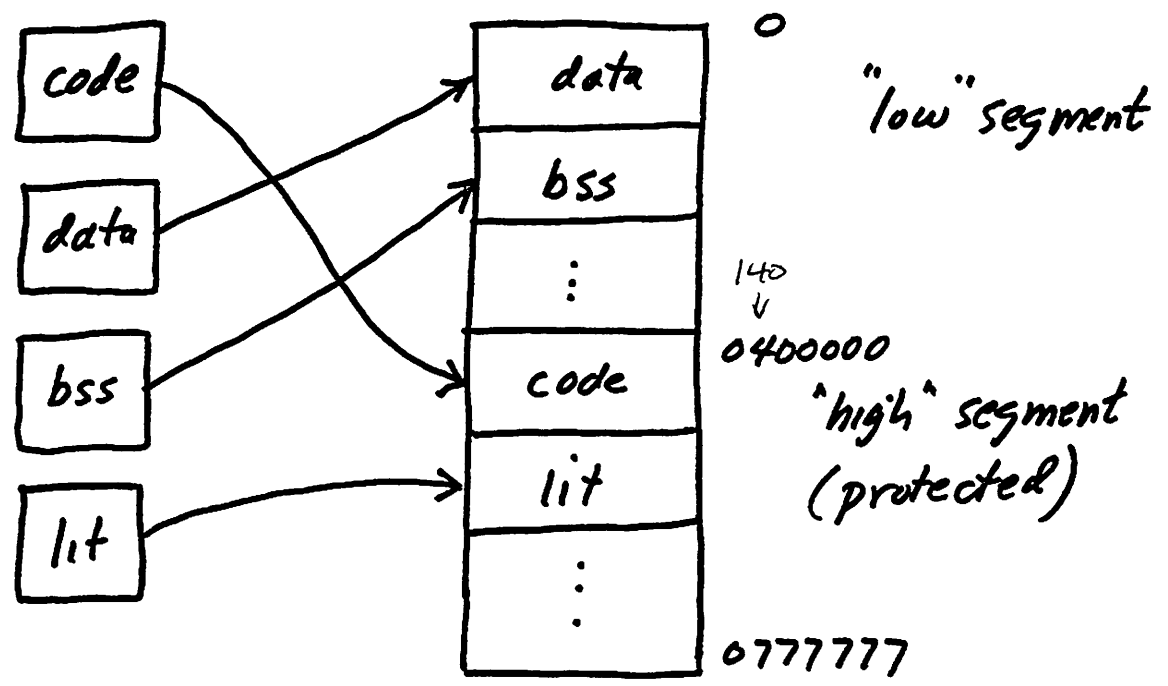
- loader maps these logical segments onto the target machine, using hardware as appropriate.

e.g. Cyber/175: single address space, no protection

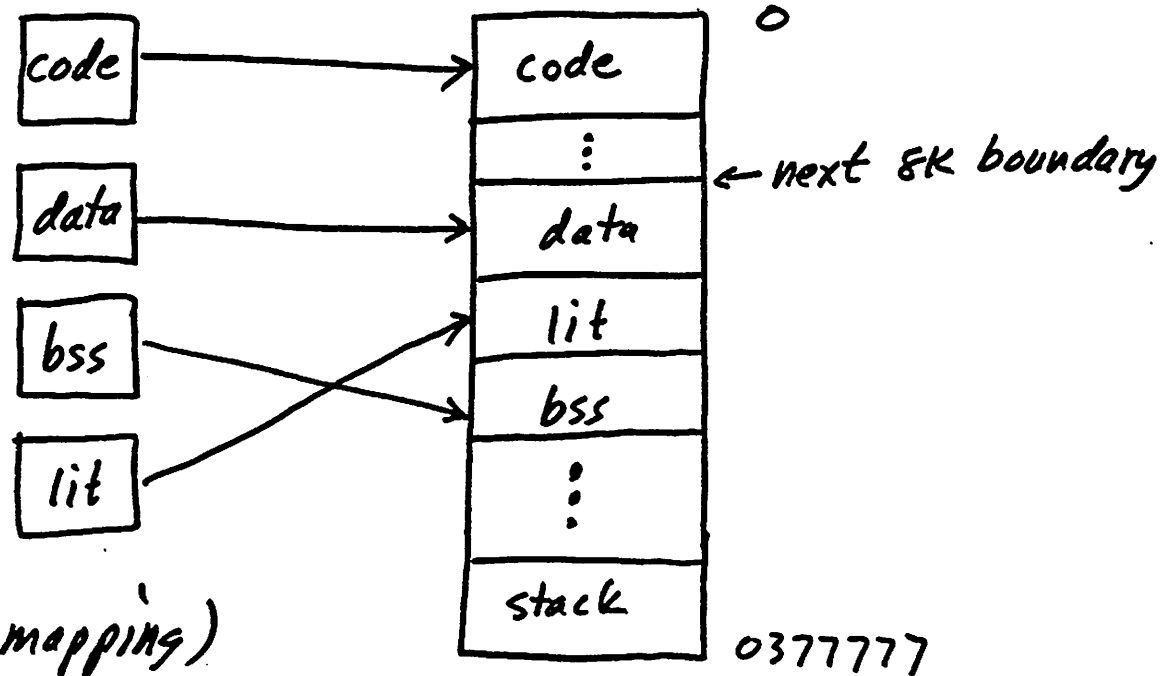


# Segments cont'd

e.g. DEC-10: 2 segments, 1 protected



e.g. PDP-11/70: 8 segments, up to 8 protected



(UNIX mapping)

## Segment Commands

`.seg id`

- start, or resume, placing text in segment id.
- may flip back & forth between segments as necessary to simplify code generation.

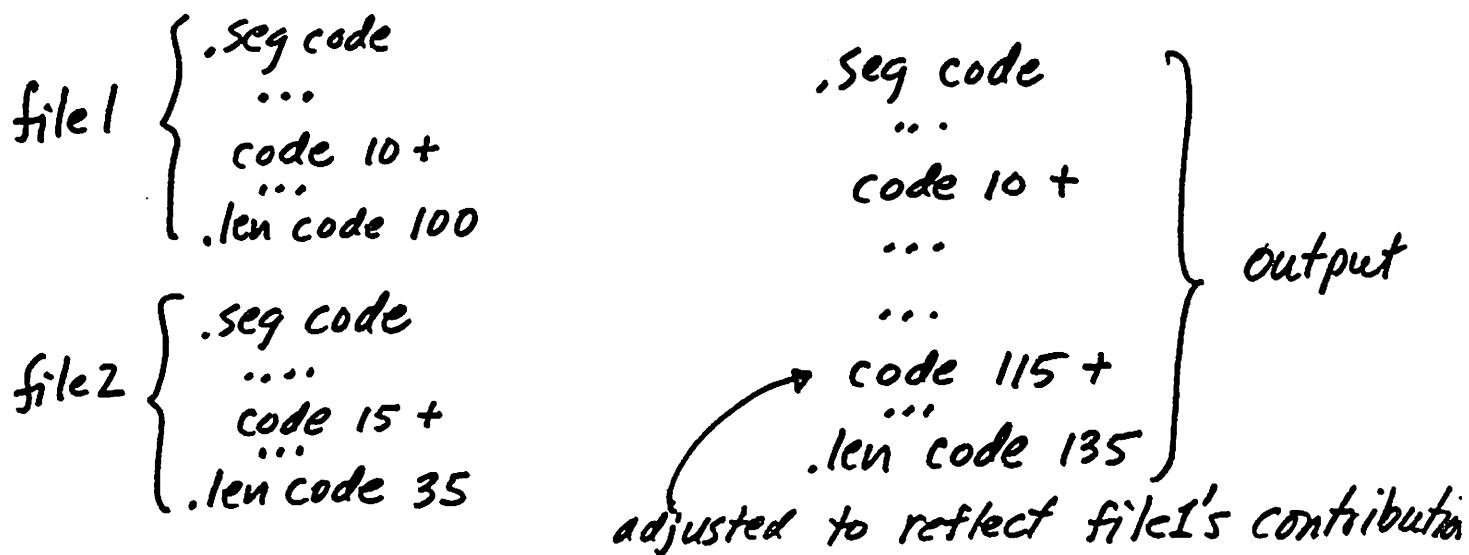
```
input:    int a;
          int b = 6;
          int g()
          {
            ...
            a = g("hi there");
            ...
          }
```

```
output:  .seg bss
          (definition for a)
          .seg data
          (definition for b)
          b
          .seg code
          (code for g)
          .seg lit
          (constant "hi there")
          .seg code
```

## Segment Commands, cont'd

.len id con

- indicate length of segment id in units agreed upon by compilers & loader
- units may not be "interesting cell" e.g. on DEC-10, cell may be 18 bits, lengths in words (36 bits).
- segment lengths used by link to adjust base locations when several object files contribute to the same segment; used to obtain proper relocation of references in the output file



## Segment Commands cont'd

- order of linking affects relocation (but not semantics!)

file1

```
.seg code
...
code 10 +
...
.len code 100
```

file2

```
.seg code
...
code 15 +
...
.len code 35
```

.link file1 file2

```
.len code 135
.seg code
...
code 10 +
...
code 115 +
...
```

.link file2 file1

```
.len code 135
.seg code
...
code 15 +
...
code 45 +
...
```

- len commands appear first in the output so loaders can compute segment addresses, if necessary, and allocate space.



## Segment Commands, cont'd

• org con

- instructs loader to place text at different locations in current segment.
- gives "random access" capability (used for "backpatching").
- link copies org commands to the output, adjusting con if necessary

## Symbols

- `def` commands define symbols as expressions (polish suffix)
- `link` evaluates expressions substituting expressions for symbols as necessary
- symbols are defined in terms of other symbols and segment names (i.e. base locations).

e.g.

```
.def rows data 25 +
      symbol "value"
```

subsequent line of object text

```
rows 4 +
```

becomes (after evaluation)

```
data 29 +
```

Symbols, cont'd

- symbols can be defined & referenced in any order (e.g. forward references)
- in expressions, segment names may be used as symbols.

e.g.

or/

```
.seg code
.def a b
...
a 1 +
...
.def b code 25 +
```

```
.seg code
.def b code 25 +
.def a b
...
a 1 +
...
```

both yield

```
.seg code
...
code 26 +
...
```

- also, local labels for control constructs, e.g. if, while, etc.

## Symbol Sets

• Let

$R$  be set of referenced symbols,

$D$  be set of defined symbols,

$S$  be set of segment names.

- $R$  comes from <sup>symbols in</sup> expressions, <sup>left operands</sup>
- $D$  comes from .def commands
- $S$  comes from .seg & .len commands

- Goal: "get rid of symbols", or a resolved file, i.e. a file for which

$$R \subseteq D \cup S$$

- given values for symbols in  $S$ , a resolved file can be loaded.
- unresolved files have external references, i.e. symbols in  $R$  that are not in  $D \cup S$ .

- linking seeks to combine unresolved files into a single resolved file.
- linking is an iterative process; an unresolved file is valid output.

e.g. given 3 unresolved files that each reference symbols in all 3 files, can link in 1 step:

link file1 file2 file3 -o out  
└───┘  
specifies output file

or in 2 steps:

link file1 file2 -o temp  
 link file3 temp -o out

- if  $R \cup S \subseteq D$ , file is absolute; must be loaded at locations consistent with segment definitions.

## Output

- link produces output in a form that can be loaded in one pass:

.len commands

.def commands

.seg  $id_1$

(text for segment  $id_1$ )

...

.seg  $id_k$

(text for segment  $id_k$ )

(segments are contiguous)

- len commands are first so loaders can know lengths and hence position of all segments.
- def commands come next so symbols are defined before used in object text
- segments appear last, and are contiguous (no segment switching).
- output is valid link input!

## Compiler Construction

- method 1: use assembler to produce object code

source → assembler lang. → object code

e.g. C compiler

- advantages:

- easy to generate assembler lang.
- smaller compiler
- readable output (easy to debug)
- avoids "forward reference" problem
- uses existing tool to do dirty work

- disadvantages:

- slow
- idiosyncratic (e.g. ordering, naming restrictions, inability to handle source constructs)
- fixed output format
- at least 3 passes + linker
- duplication of work (e.g. symbol table)

## Compiler Construction cont'd

- method 2: produce object code directly  
 source  $\rightarrow$  object code  
 e.g. most commercial compilers.

- advantages:

- fast (sometimes)
- fewer passes
- smaller (and fewer) files

- disadvantages:

- larger compiler
- difficult to debug
- idiosyncratic syntax, binary i/o
- silly, archaic restrictions

- use link: retain advantages of both methods <sup>+ disadvantages</sup>

source  $\rightarrow$  link code  $\rightarrow$  executable code

may be fewer passes when linking taken into account; no worse than assembly lang.



## Using link as an Assembler

- consider  $f(a+b, c)$ ; macro-10 assembly language code is (good code!)

```

move 2, a      ; load a
add 2, b       ; compute a+b
push 15, 2     ; push a+b
push 15, c     ; push c
pushj 15, f    ; call f
adjsp 15, -2   ; remove arguments

```

- in link code, "push 15, c" becomes (36-bit cells)

026174000000 c +

- can define symbol "push" to make this more readable:

```
.def push 0261000 18 <
```

...

```
push 15 23 < + c +
```

register 15 in right place.

- definitions such as push need be made only once.

## link as an Assembler, cont'd 23

- can make "15 23 <" part of definition of push, if all push (and pushj, adjsp, etc.) are for same register; e.g.

```
.def push 0261740 18 <  
...  
push c +
```

which is as readable as assembly lang.

- assuming definitions for other opcodes and for registers (e.g. `.def r2 2 23 <`), above sequence is:

```
move r2 a + + (or move r2+ a+)  
add r2 b + +  
push 2 +  
push c +  
pushj f +  
adjsp 0777776 +
```

bogus!

## link as an Assembler, cont'd

- typical usage: simply generate bit patterns in code generators, or have option to produce symbolic output for debugging.
- for assembling "small" assembly language programs, write program using symbolic form (shown above) and link with definitions file, e.g.

```
link defs pgm.asm -o pgm.obj
```

- Experience with Y compiler:

link code gen. is 15% shorter than  
macro-10 code gen.;

link code gen. is 54% shorter than  
binary object code gen.

(lines of code in output module).

## Local Symbol Resolution

- "local" symbols can be hidden from outside by 3 different approaches.

Consider stack module in C:

```
static int sp;
static int stack [STKSIZE];
```

↑ e.g. 100

```
push(x)
```

```
int x;
```

```
{
  ;
}
```

```
pop()
```

```
{
  ;
}
```

```
}
```

static also means that name won't be known outside of the file

- sp & stack are local symbols and are not accessible from outside the stack module.

## Local Symbol Resolution, cont'd

- method 1: "suppress" definitions of `sp` & `stack` and always run output of compiler thru link:

```

.seg data
.def -s sp data 0 +
.def -s stack data 1 +
.seg code
.def push code 0 +
  ;
.def pop code 36 +
  ;
.len data 101
.len data code 80

```

- `-s` option causes link to omit definitions for those symbols from the output; hence, `sp` & `stack` do not appear and cannot be referenced.
- disadvantage: must run link an extra time

## Local Symbol Resolution, cont'd

- method 2: add string to local names (such as the file or module name) to make local names unique:

```
.seg data
.def -s stack.c.sp data 0 +
.def -s stack.c.stack data 1 +
.seg code
.def push code 0 +
  :
.def pop code 36 +
  :
.len data 101
.len code 80
```

- can still use `-s` to reduce # of symbols, but do not have to run link after every compilation.

## Local Symbol Resolution, cont'd

- method 3: "backpatch" forward references when target is known using `.org` command.
- assume expression evaluation starts with "current value" of the cell on the stack; then  

$$\text{code} +$$
means "add code" to current value.

Example: `while (a < 0) S`

L1, L2 `if (a >= 0) goto L3`

`S`

`goto L1`

output:

code location counter	011	<code>movei 6, a</code>	<code>; get a's lvalue</code>
	012	<code>move 7, (6)</code>	<code>; get a's rvalue</code>
	013	<code>movei 8, 0</code>	<code>; get 0</code>
	014	<code>cmpl 7, 8</code>	<code>; a &lt; 0 ?</code>
	015	<code>jrst</code>	<code>; no, goto L3</code>
		...	
		<code>S</code>	
		...	
	021	<code>jrst code 011 ++</code>	<code>; jump to L1</code>
	022		

## Backpatching cont'd

- at location 022, L3 is defined; backpatch all references to L3 via .org commands:

```

.org 015
code 022 + +
.org 022

```

adds "code 022 +" to "jrst" at location 015.

- .org 022 puts location counter "back" for subsequent code
- backpatching specifications (i.e. locations to backpatch) maintained until label is defined, then appropriate sequence of .org commands are emitted.
- complicates compiler.



Computer Science 453  
Fall 1983  
Class Notes, Set 5

October 13, 1983

#1/

1.

do stmt while (expr)

a) code template:

L stmt  
L+1 if (expr) goto L  
L+2

↘ "true" label

b) add to stmt():

added to first(stmt), follow(stmt) ...

Case DO: dostat(loop = genlabel(3));  
break;

↘ "loop handle" L

dostat(lab)

int lab;

{

t = gtoK();

deflabel(lab);

stmt();

deflabel(lab + 1);

mustbe(WHILE);

mustbe('(');

walk(cnode(expr()), lab, 0);

mustbe(')');

deflabel(lab + 2);

#2/

Conditional Expressions

2.

- used wherever a value may be used;  
add "conversion" node in rvalue:

```
struct node *rvalue(p)
```

```
struct node *p;
```

```
{
```

```
...
```

```
if (p->e-type == T_COND)
```

```
    p = node('?', T_INT, p, NULL);
```

```
    return (p);
```

```
}
```

- rvalue adds node with '?' operator that is used in walk to generate code according to following template:

```

    ↖ p->e-left
    if (e) goto L
    ↗ global variable → c.cc = 0
                    goto L+1
L      c.cc = 1
L+1   (return value of c.cc)

```

• in walk:

case '?':

tree for  
lvalue of c.cc { walk(p → e.left, lab = genlabel(z), 0);  
p1 = node(O\_ID, T\_INT | T\_ADDR,  
lookup("c.cc"), NULL);

tree for  
c.cc = 0; { walk(node(O\_ASSIGN, T\_INT,  
p1, constant("0", T\_INT)), 0, 0).

jump(lab+1);  
deflabel(lab);

tree for  
c.cc = 1; { walk(node(O\_ASSIGN, T\_INT,  
p1, constant("1", T\_INT)), 0, 0);

deflabel(lab+1);

tree for  
rvalue of c.cc { r = walk(rvalue(p1), 0, 0);  
break;

#3

Temporaries

- allocate fixed number of temporaries in local portion of frame each large enough to hold largest datatype.
- in ralloc, when out of registers "spill" least-recently-used register into a temporary.
- exprcode returns an index into an array of "resources", which can be registers or temporaries. Level of indirection is necessary to reflect state changes due to spilling.
- emit checks for temporaries, and reloads them into registers (not necessarily the same ones!) if necessary.

## Temporaries, cont'd

Example:  $a = b + c$  w/ 3 registers

```

    movei 6, a      ; a's lvalue
    movei 7, b      ; b's lvalue
    move 8, (7)    ; b's rvalue
    movei 7, c      ; c's lvalue
    spill → movem 6, 1(14) ; spill reg. 6
    move 6, (7)    ; c's rvalue
    add 8, 6       ; b+c
    reload → move 6, 1(14) ; reload a's lvalue
    movem 8, (6)   ; a = b+c
  
```

---

## Implementation

```

int rcount [16];      /* reference count */
int rtime [16];      /* LRU time */
int timer = 0;       /* the clock */

struct resource {
    int count;        /* reference count */
    int reg;         /* register # or */
    int temp;        /* temporary # */
} resources [RESOURCEST];
  
```

6.

*\* ralloc - allocate a register, return resource\**  
int ralloc()

{

int r, r1, min;

min = MAXINT;

for (r = 6; r <= 12; r++)

if (rcount[r] == 0)

break;

else if (rtime[r] < min) {

r1 = r;

min = rtime[r];

}

if (r > 12)

rspill(r = r1);

rcount[r] ++;

rtime[r] = ++timer;

r1 = index of free resource;

resources[r1].count = 1;

resources[r1].reg = r;

resources[r1].temp = 0;

return (r1);

}

```

/*emit - emit an instruction */
emit(op, r, name, off, idx)
int op, r, off, idx;
char *name;

```

resource  
index

```

{
  register numbers → { r = reload(r);
                      idx = reload(idx);

```

```

    if (off is a temporary) {
      idx = 14;
      off = proper frame offset;
    }

```

⋮

```

}

```

- reload ~~may~~ calls ralloc to allocate a register, which can cause a spill.



## Spilling

```

int temps[11]; /* non-zero if busy */
/* rspill - spill register r */
rspill(r)
int r;
{
    int ri, t;

    t = index(1..10) of free temporary;
    temps[t]++;

    emit(MOVEM, r, NULL, toff+t, 14);
                                     ↑
                                     offset-1 to 1st temp.

    for (ri = 0; ri < RESOURCES; ri++)
        if (resources[ri].reg == r) {
            resources[ri].temp = t;
            resources[ri].reg = 0;
            rcount[r]--;
        }
}

```

## Reloading

*/\* reload - reload resource r, if necessary \*/*  
*/\* returns register number \*/*

```
int reload(r)
```

```
int r;
```

```
{
```

```
    int r1, r2, t;
```

```
    if (t = resources[r].temp) {
```

```
        r1 = ralloc();
```

```
        r = resources[r1].reg;
```

```
        emit(MOVE, r1, NULL, toff+t, 14);
```

```
        for (r2 = 0; r2 < RESOURCES; r2++)
```

```
            if (resources[r2].temp == t) {
```

```
                resources[r2].temp = 0;
```

```
                resources[r2].reg = r;
```

```
                rcount[r]++;
```

```
            }
```

```
        rfree(r1);
```

```
        temps[t] = 0;
```

```
    }
```

```
    else r = resources[r].reg;
```

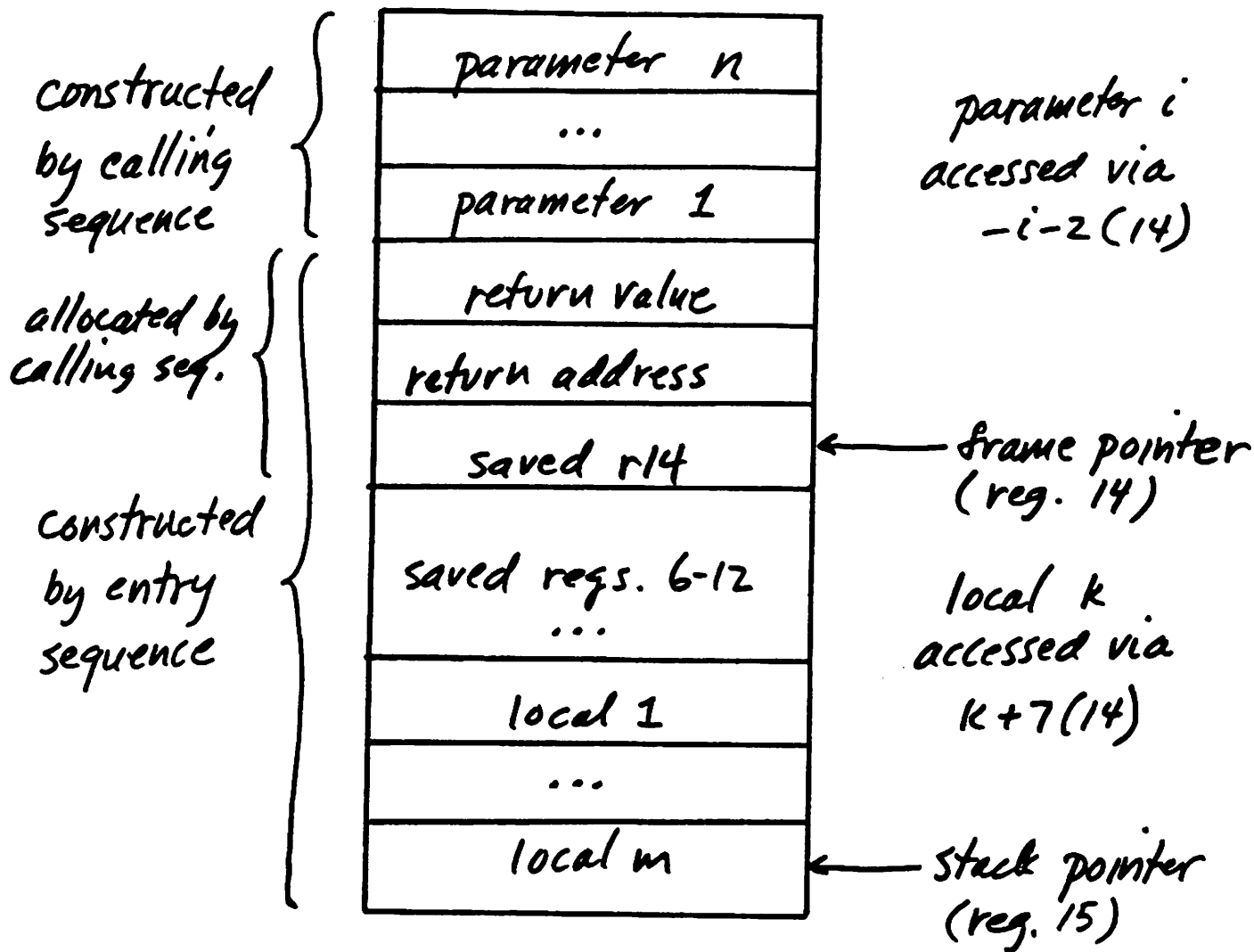
```
    rtime[r] = ++timer;
```

```
    return (r)
```

```
}
```

#41

Frame Construction



- frame construction coordinated among calling, entry, exit sequence to give "best" code.
- must allocate stack space before using it! (nested calls, interrupts!).
- best to allocate/deallocate in same sequence (eg. allocate/deallocate same time as calling seq)

calling sequence: (for n arguments)

```

    adjsp 15, n+3 ; allocate space
    :
    movem r, -i-2(15) ; argument i
    :
    movei 2, $L ; get return address
    jrst f ; jump to procedure
    $L: move r, -2(15) ; get return value
    adjsp 15, -n-3 ; deallocate space

```

could use

```

    adjsp 15, 1
    movem r, 0(15)

```

for each argument, but code is longer.

entry sequence: (m words of locals)

```

f:   movem 2, -1(15) ; save return address
     movem 14, 0(15) ; save frame ptr
     move 14, 15 ; set new frame ptr
     adjsp 15, m+7 ; allocate local space
     movem 6, 1(14) ; save reg. 6
     :
     movem 12, 7(14) ; save reg. 12

```

exit sequence: (return value in r)

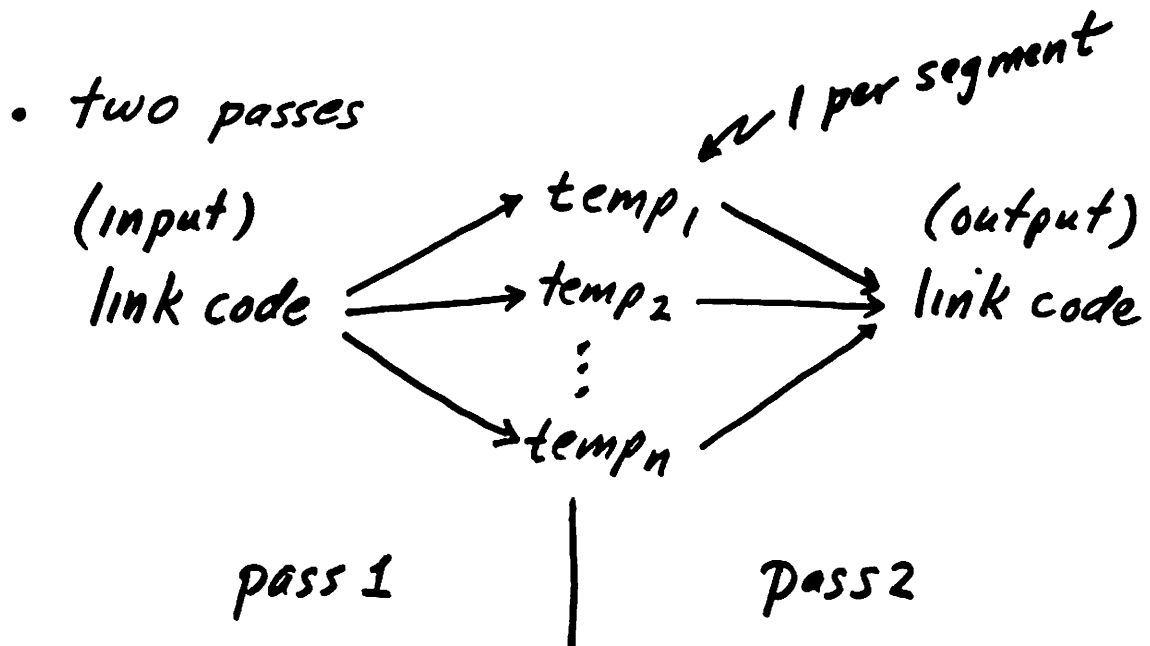
```

movem    r, -2(14)    ; set return value
move     6, 1(14)     ; restore reg. 6
      ;
move     12, 7(14)    ; restore reg. 12
move     15, 14       ; deallocate locals
move     14, 0(15)    ; restore frame ptr.
move     2, -1(15)    ; get return address
jrst     (2)          ; return to caller

```

- note entry/exit sequence is independent of number of actual parameters, and that parameter offsets are independent of number of actual parameters. This makes writing procedures w/varying number of arguments easier.

## Implementation



• pass one:

- distributes text into temporary files, one per segment
- build symbol table representing sets  $R, D, S$
- evaluate definitions
- compute segment lengths

• pass two:

- combines temporary files into output
- evaluates expressions in object text (as far as possible).

## Modules

* link.h	definitions
* lsym.c	symbol table, storage alloc.
lmain.c	main program
lpass1.c	pass one
lpass2.c	pass two
lexpr.c	expression evaluation
llib.c	library searching (later!)

\* - provided

- psuedo-code for pass 1 & 2 given in paper.
- major task: getting symbolic evaluation & relocation correct.

## Symbol Table (lsym.c)

```

struct symbol {
    char *s_name;    /* entry name */
    int s_flags;    /* flags, see below */
    char *s_value    /* value (an expr) */
    struct symbol *s_link;
}

```

/\* flags \*/

```

#define DEF 001    /* symbol is defined */
#define REF 002    /* symbol is referenced */
#define SEGMENT 004 /* symbol is a seg. name */
#define SUPPRESS 010 /* symbol has -s on */

```

- idea: treat segment name like other symbols during evaluation
- routines: struct symbol \*lookup(s):  
looks up s, installing it if necessary.  
also stlookup, dump, foreach, alloc, initsym  
(see lsym.c)



Main Program (Suggested 32  
organization)

```
main(argc, argv)
int argc;
char *argv[];
{
    initsym();
    init1();
    for (each argument)
        if (argument is an option)
            process option;
        else
            pass1(argument);
    if (no file arguments)
        pass1(standard input);
    end1();
    init2(output file);
    pass2(output file);
    end2(output file);
}
```

-U - flags++  
-o ofile is next  
arg

- init1, pass1, end1 in lpass1.c
- init2, pass2, end2 in lpass2.c

## Segment Information

- \* side effect of pass one is a segment table:

```

struct segment {
    FILE *x-fp;      /* temp file pointer */
    char *x-fname;  /* temp file name */
    int x-length;   /* length in current file */
    int x-base;    /* accumulated length */
    struct symbol *x-p /* symbol table
                       value of segn entry */
}
  
```

x-fname : construct using segment name,  
 e.g. "name.seg"; unlink (delete)  
 after pass 2 (unless -d).

basics!

name

x-length, x-base : set by .len commands;  
 used as "value" of segment name  
 symbol.

test for segment or ~~test~~  
 maintain value

## Pass One

- distribute text into segment temporary files.
- evaluate definitions
- build R, S, D (i.e. build symbol table).
- compute segment lengths.
- implement relocation (i.e. adjust segment base locations).

init1 - init segment table

pass1 read object opened on f  
exec pass1 operators

end1 - close temp files

```

pass1(f)
FILE *f;
{

```

effect of ~~letter~~ len ends  
must be index of location

```

    for(each segment p)
        p->x_length = 0;
    curseg = NULL; /* "current" segment */
    while (fgets(line, MAXLINE, f) != NULL)
        if (line is a command)

```

symbols start  
w/ a letter  
, digit, etc.

```

        else {
            ...
            just for -v?
            builds set R

```

eval enters symbols  
on symbol table

```

            add symbols in line to
            symbol table;
            fputs(line, curseg->x_fp);

```

thus evaluate

p => segments;  
p < segments + users; p++

adjust  
segment  
base location

```

    for(each segment p) {
        p->x_base += p->x_length;
        for(each segment p1)

```

write .len  
commands  
to temp  
file.

all temp  
files have  
len ends  
for all segments

```

        fprintf(p->x_fp,
            ".len %s %d\n",
            p1->x-p->s.name,
            p1->x_length);
    }
}

```

Pass one, cont'd

```

if (line is ".seg id") {
    p = lookup "id" in segment table;
    if ("id" is a new segment) {
        p->x_fname = "id.seg";
        p->x_fp = fopen(p->x_fname, "w");
        if (p->x_fp == NULL)
            issue fatal error & exit
        p->x_base = p->x_length = 0;
        p->x_p = lookup("id");
        enter p into segment table;
        fprintf(p->x_fp, ".seg %s\n",
                p->x_p->s_name);
    }
    curseg = p;
}

```

Pass one, cont'd

```
else if (line is ".len id n") {
```

```
    p = lookup "id" in segment table;
```

```
    if ("id" is a new segment) {
```

```
        ... (same as in .seg command)
```

```
    }
```

```
    p → x-length += n;
```

```
}
```

.len commands  
can precede  
.seg commands

- len commands can occur anywhere (and repeatedly), but their effect is deferred until the end of pass 1 when x-base is incremented by the sum of the n's in len commands.

## Pass One, cont'd

38

```
else if (line is ".def [-s] id expr") {
```

```
    sp = lookup("id");
```

```
    if (sp->s_flags & DEF)
```

```
        fprintf(stderr, "multiply defined
```

```
symbol: %s\n", sp->s_name);
```

use strtok

to stick  
away

the string

```
    sp->s_value = evaluate("expr");
```

```
    sp->s_flags |= DEF;
```

```
    if (-s specified)
```

```
        sp->s_flags |= SUPPRESS;
```

```
}
```

```
else
```

```
    fputs(line, curseg->x_fp);
```

→ .org commands & "comments".

any other dot commands

(can also turn on  
SUPPRESS to avoid  
output in pass 2)

move to 36,  
don't want  
output a lfp  
for a seg

# Pass One: Example (from paper)

- two files, a.o & b.o:

a.o

b.o

```

.seg code
.def start code
0263 f +
.seg data
.def -s a-l data
data 01 +
.org 6
code
.len code 1
.len data 7

```

```

.seg code
.def f code
0201 b-l +
.seg data
.def -s b-l data
start ←
.len code 1
.len data 1

```

inter-segment reference

link b.o a.o { Note: b.o is linked first }

- after pass 1, symbols are

<u>name</u>	<u>flags</u>	<u>value</u>
f	DEF   REF	"code"
a-l	DEF   SUPPRESS	"data 1 +"
start	DEF   REF	"code 1 +"
b-l	DEF   REF   SUPPRESS	"data"
code	DEF   REF   SEGMENT	-
data	DEF   REF   SEGMENT	-

"code 0 +"



## Example, contd

- temporary files:

<u>code.seg</u>		<u>data.seg</u>
.seg code		.seg data
0201 b-2 +	} b.0	start
.len code 1		.len code 01
.len data 1		.len data 1
0263 f +	} a.0	data 01 +
.len code 1		.org 6
.len data 7		code
		.len code 1
		.len data 7

- Pass two combines code.seg & data.seg into single output file, evaluating expressions in object text.

## Expression Evaluation

- problem: must avoid "double relocation" while still permitting "active substitution" of symbol values.

Example: in pass two, will evaluate

0263 f +

substituting value of "f" yields

0263 code +

relocating substituted  
value is active  
sub, other is of 123456

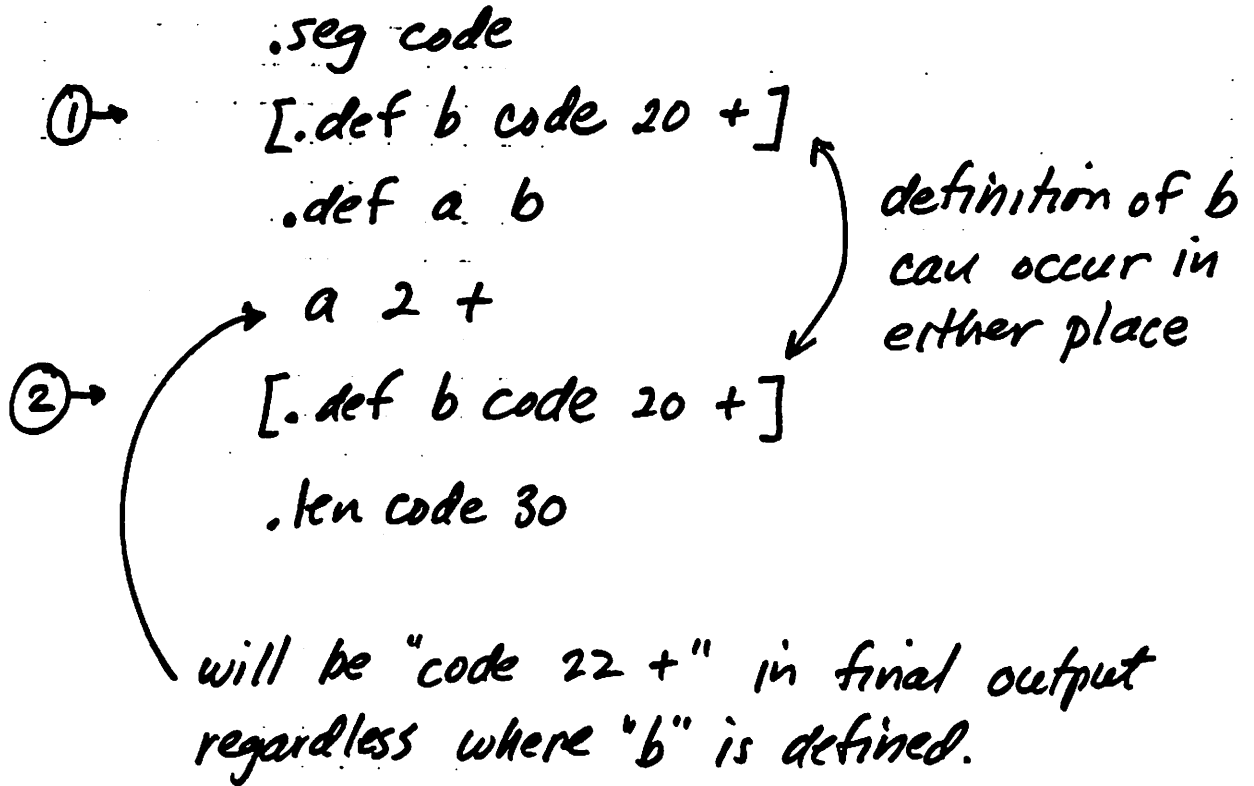
but value of "code" must not be substituted because appropriate relocation of segment names in definitions is done in pass one.

On the other hand, in expression

code

pass two must evaluate code to obtain proper relocation in object text:

code 01 +

another example:

- in position ①, a is defined as "code 20 +" and "a 2 +" is evaluated as

$a\ 2\ + \rightarrow \text{code } 20\ +\ 2\ +$

- in position ②, a is defined as "b" and "a 2 +" is evaluated as

$a\ 2\ + \rightarrow b\ 2\ + \rightarrow \text{code } 20\ +\ 2\ +$   
~~code 20 + 2 +~~      ok, don't get

## Possible Solutions

- special case segment names so their "evaluation" depends on context (i.e. definition in pass 1, object text evaluation in pass 2)
- add syntax to expressions (for use only within link) that suppresses substitution of value, e.g.

~~code~~ .def f code <sup>seg</sup>

evaluation of "code" yields "code n +"  
 where "n" is x-base field for segment "code".

'name suppresses further evaluation & is removed when definitions are printed in pass 2.

- expression evaluation must recognize name and passes it through (but should turn on REF bit?).

code

## Pass Two

44

- output final segment lengths
- output definitions (unless suppressed)
- combine segment temporary files into output
- evaluate expressions in object text

### Pass 2 Routines:

- `init2(f)` - output segment lengths and symbol definitions to `f`
- `pass2(f)` - combine temporary files into output file `f`
- `end2()` - remove temporary files (unless `-d`).

## Pass Two, cont'd

45

```
init2(f)
FILE *f;
{
    for (each segment p)
        fprintf(f, ".len %s %d\n",
                p->x-p->s-name,
                p->x-base);
    foreach (putdef, DEF, f);
}
```

- x-base field has final length for each segment
- putdef outputs definition if  $s\_flags \& (\text{SEGMENT} | \text{SUPPRESS}) == 0$ .

```
pass2(f)
```

```
FILE *f
```

```
{
```

```
    FILE *fp;
```

```
    for (each segment p) {
```

```
        fp = fopen(p->x_fname, "r");
```

```
        if (fp == NULL)
```

```
            issue fatal error & exit;
```

```
        for (each segment p1)
```

```
            p1->x_base = 0;
```

```
        while (fgets(line, MAXLINE, fp))
```

```
            if (line is a command)
```

```
                ...
```

```
            else {
```

```
                s = evaluate(line);
```

```
                → output s to f;
```

```
            } strip '"'s from line
```

```
        fclose(fp)
```

```
    }
```

```
}
```

.len commands  
in temp files  
cause x\_base  
to be incremented

## Pass Two cont'd

47

```
if (line is ".len id n") {
```

```
    p1 = lookup "id" in segment table;
```

```
    p1->x_base += n;
```

```
}
```

```
else if (line is ".org n") {
```

```
    n += p->x_base;
```

```
    fprintf(f, ".org %d\n", n);
```

```
}
```

```
else
```

```
    fputs(line, f);
```



## Example

48

- temporary files for b.o & a.o input:

code.seg

```
.seg code
0201 b-l +
.len code 1
.len data 1
0263 f +
.len code 1
.len data 7
```

data.seg

```
.seg data
start
.len code 1
.len data 1
data 01 +
.org 6
code
.len code 1
.len data 7
```

- output:

```
.len code 2
.len data 8
.def f code
.def start code 01 +
.seg code
0201 data +
0263 code +
.seg data
code 01 +
data 02 +
.org 7
```

## Library Searching

- link object modules based on current value of  $R - (DUS)$ , i.e. symbols referenced but not defined; sort of a "data-driven include mechanism".

- libraries are searched during pass one:

link a.o b.o -l subs.lib

performs pass 1 on relevant object files in "subs.lib".

- more than one library can be searched, and same library may be searched more than once:

link a.o -l lib1 b.o -l lib1 -l lib2

- function of library searching is simple: "if library member defines a symbol in  $R - (DUS)$ , link that member";
- form of library affects search

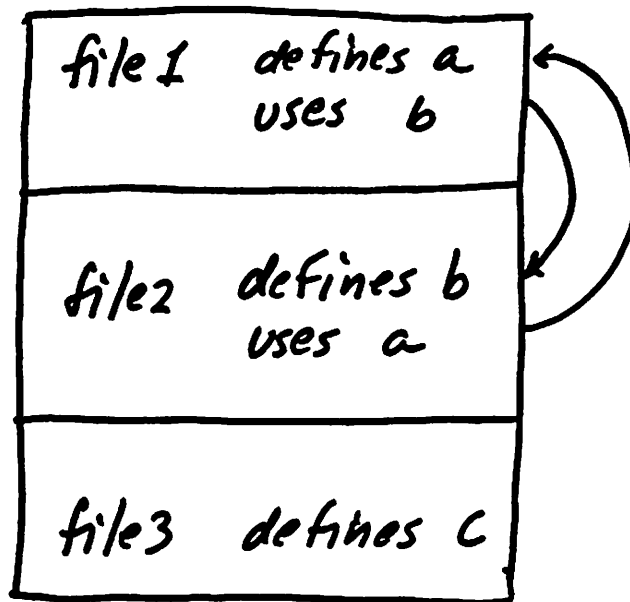
## Library Representations

50.

- simple concatenation of object files (e.g. DEC-10):

for each object file in library  
if this file defines an undefined symbol  
link this file

- object files must be concatenated in topological order:



note: library searching affects R, D, & S!

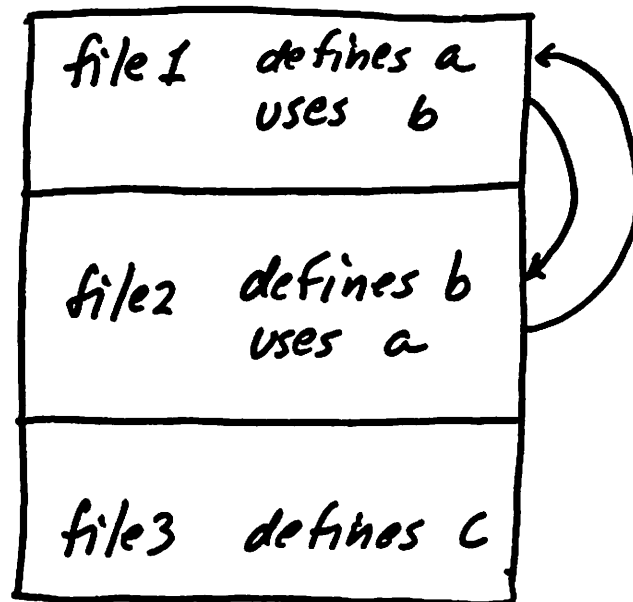
- if a is undefined, both file 1 & file 2 are linked; if only b is undefined, file 2 is linked and a remains undefined.

## Library Representations

- simple concatenation of object files (e.g. DEC-10):

for each object file in library  
if this file defines an undefined symbol  
link this file

- object files must be concatenated in topological order:



note: library  
searching  
affects R,  
D, & S!

- if a is undefined, both file 1 & file 2 are linked; if only b is undefined, file 2 is linked and a remains undefined.

## Library Representations, cont'd

- simple concatenation may require several passes to get everything: at worst use
  - for each undefined symbol
    - for each object file in library
      - if this file defines symbol
        - link it
- topologically sorted concatenation of object files can be searched in one pass.
- better organization: concatenation of object files preceded by index or directory of (symbol, file) pairs for all symbols defined in the library:
  - for each undefined symbol
    - if symbol defined in library
      - link corresponding file
- can use either random access or multiple passes over the library in search

**Computer Science 453**  
**Fall 1983**

**Class Notes, Set 6**

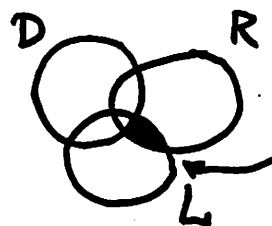
**November 1, 1983**

## Implementation

- library contains entries, each of which is an object file that contains symbol definitions.

- let  $L$  = set of symbols defined by all entries in the library.

- want to load those entries that define symbols in  $(L \cap R) - (D \cup S)$  [assume "US" below].



- algorithm:

open library

load index, i.e. construct  $L$

while  $((L \cap R) - D$  is non-empty)

for (each <sup>library</sup> entry  $q$ )

if ( $q$  defines a symbol in  $(L \cap R) - D$ ,

pass 1 ( $q$ )

expunge symbols in  $L - (R \cup D)$

reclaim space

↑ symbols only in  $L$

may involve  
random-access  
or multiple  
passes.

$(L \cap R) - D$   
↑  
undef. syms in  
library

## Implementation, cont'd

- additional symbol table information:

```

struct symbol {
    char *s_name;
    int s_flags;
    char *s_value;
    struct entry *s_entry;
    /* library entry structure */
    struct symbol *s_link;
};
  
```

(see lsym.c)

```

#define LIB 020 /* symbol defined
                in a library */
                ↑
                member of L
  
```

- library entry information (see link.h)

```

struct entry {
    char *e_name; /* entry name */
    long e_pos; /* location in library */
    int e_size; /* length of entry */
    int e_count; /* # symbols in (L ∩ R) - D
                 ? →
    struct entry *e_link; /* next entry */
};
  
```



## Library Format

- archive files:

header

entry

header

entry

...

- first entry is "index": ↙ length of entry in chars.

#-h- index size ↙ date information

# symbol name

... ↗

entry name where "symbol" is defined

- remaining entries are object files.
- no position information; must be computed while reading library.

```

lib(name)
char *name;
{
    FILE *f;

    f = fopen(name, "r");
    if (f == NULL) {
        fprintf(stderr, "can't search %s\n",
                name);
        return;
    }

    getindex(f); /* construct index structure */
    search(f);
    fclose(f);
    expunge();
    free index structure;
}

```

- | 500  
call lib(500) ?

## Processing the Index

```
getindex(f)
```

```
FILE *f;
```

```
{
```

```
    struct symbol *p;
```

```
    struct entry *q;
```

```
    for (each line in index) {
```

```
        q = pointer to entry structure for "name"
```

```
        if (q is a new entry) {
```

```
            fill in q's fields;
```

```
            add q to list of entries.
```

```
        }
```

```
        p = lookup("symbol");
```

```
        p->s_flags |= LIB;
```

```
        p->s_entry = q;
```

```
    }
```

```
}
```

- end up with all symbols in symbol table, <sup>and</sup> with <sup>them</sup> LIB set and s\_entry pointing to entry structure, and a list of entry structures.

2nd field



note - not p or

## Searching the Library

57.

- two phase search:
  - 1) read library sequentially, linking necessary modules, and computing e-pos & e-size for second phase.
  - 2) read library randomly, as necessary to link remaining required modules.
- "linking" an entry involves using pass 1 to read the entry.
- linking an entry may change  $(L \cap R) - D$ ; must recompute that set after calling pass 1. (hint: use  $g \rightarrow e\text{-count}$  field).

```

search(f)
FILE *f;
{

```

```

    struct entry *q;

```

```

    /* phase one */

```

```

    for (each entry) {

```

```

        q = pointer to corresponding
            entry structure;

```

```

        q → e-size = size from header;

```

```

        q → e-pos = ftell(f);

```

```

        if (q defines a symbol in (L ∩ R)-D)
            pass1 (this entry);

```

```

        else skip this entry;

```

```

    }

```

While there are  
symbols in the  
library that  
are referred  
but not defined

```

    /* phase two */

```

Referred, in lib

```

    while ((L ∩ R)-D is non-empty)

```

```

        for (each entry q in list of entries)

```

```

            if (q defines a symbol in (L ∩ R)-D) {

```

```

                fseek(f, q → e-pos, 0);

```

```

                pass1 (this entry);

```

```

            }

```

```

    }

```

in BH I get ulbrun, I'm  
going to let you know  
I'm going to ruin your  
life this term -

①  
8/23T

Intro DEC-10 - 27, 30, 31

STUM

Don't miss a class

1200 lines

OH TH - 2:00  
Mad on 10 (1055, 21)

HW -  $\frac{1}{2}$  /  $\frac{1}{2}$     3 Exams - Equal weight

A Bare machine

~~applets~~

"bare" machine

---

operating system

system calls, { i/o services 452  
                          { processes

---

basic system software 453

e.g. compiler, linker, editor

perf. monitoring, etc., i.e. essential tools

---

"application" software

eg. software tools, editors,  
DBMS, etc.



8/23T

### Cozater

- For ~~system~~ prog lang.
- must produce reasonable be relative size; using ~~reporting~~ must produce correct code.

### For CS 453

- write a rd cozater for subset of C in (B) C on the OEL-10.
- 1-pass, emits link object code
- machine independent parser ("top half")  
machine dep. code gen.
- classical design

### Linkers

- machine-independent linker editor
- links textual object lang.
  - ~~open~~ ~~for~~
  - md, control
  - ~~ed~~ mi form
- permits "iterative link" i.e. output can be used as input
- ~~includes library search~~



8/23 ✓

- requires syntax, and loader
- includes library searching  
(~~data~~ libraries or add. files)
- replace traditional use of an assembler

### Debuggers

- low level debug - at the instruction level
- use to build high-level debugger pipe at the C statement level
- features
  - editor style interface
  - breakpoint
  - single step
  - procedure trace
- requires computer words
- editor interface - like editing a running program
- programmer defines debugging functions (print lib, etc.)

1250

5

8/23

## Monday

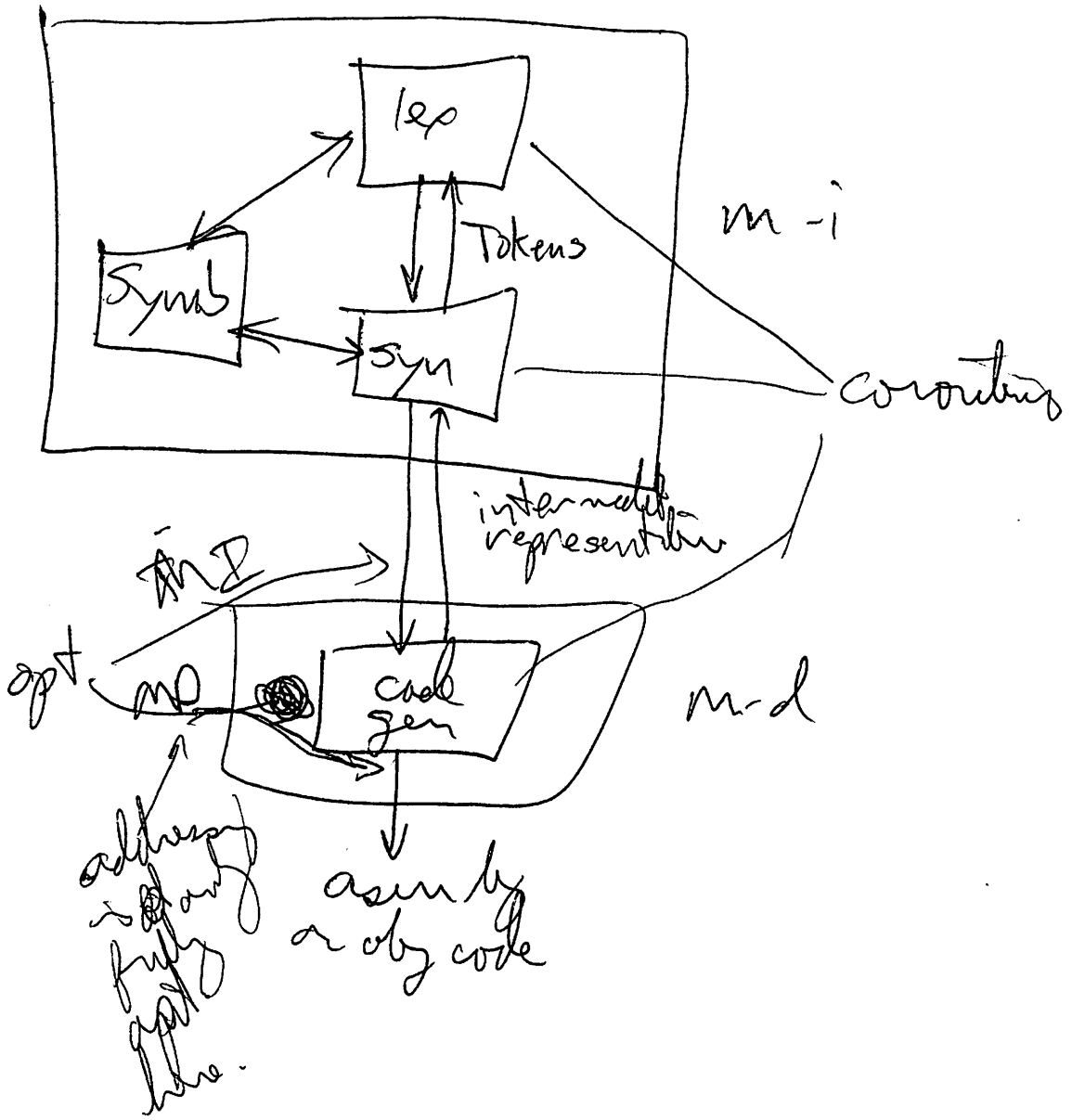
- shell level perf. analysis
- execution frequency & shell & procedure
- time history, sample taking
- data presentation, tables, program sheets, call graphs.

## com logs

- Unix Shell; shell program
- display shells, port & execs
- interactive & batch output
- language design / binary issues /  
data processing
- unifs prog & command languages;  
BZ & related systems

187230

# Compiler Architecture



~~we might be able to~~ tell how  
to go on, but all you'll be  
able to do is ride."

(1983)

8/25H

Notes in quick copy center  
switch.ini [1055, 21]

Lex Analysis

pairs of values (type, value)

constant = ; con, "value"

↑  
pointer into symbol table

"Why does the shell always walk behind the screen?"

Returns tokens, sets global var pointer

must handle errors

Trees will be built before it is emitted

Code Gen

- read IR (abs. exp. trees), misc dirs
- manage mach. resources
- assign mach. dep. parms e.g. parms  
frame layout
- emit output

IR:  
goals  
types  
paths  
parms  
trees

10/10/18

the more things I do  
the more I like them

the more I do  
the more I like them

the more I do

the more I do the more I like them

the more I do the more I like them

the more I do

the more I do the more I like them

the more I do the more I like them

the more I do

the more I do the more I like them

the more I do

Most Comp. Classes spend a lot of time on parsing - parsing is slow

8/25/14

## Symbol Table

- maintain semantics referred an identifier, const, e.g. type, address, size, etc.
- maintain multiple instances of some identifier (local, global)
- maintain single-copy of all strings
  - keep all unique strings in index
  - permits fast string comparison
  - permits ~~copy~~ constant pool

grammar: precise specification of syntax

language: set (possibly  $\infty$ ) of sentences of terminal symbols

## Top Down Parsing

using the next input symbol and current derivation to properly "guess" the next derivation step

implement parsing a rec. proc for each NT

## Pitfalls

1. Left recursion
2. Backtrack
3. ambiguity  $E \rightarrow E + E \mid E * E$
4. Error msg.

Handwritten notes at the top of the page, possibly including a date or title.

4/28/8

no longer stress, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

part of the process, it's all  
part of the process, it's all

8/25/18

$$A \rightarrow \alpha \beta$$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

$\beta$  followed by  $\alpha$  or more  $\alpha$ 's

\*  $\text{First}(\alpha)$  = terminals that begin sentences derived from  $\alpha$

$$\text{Thm: } \text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$$

Comp.:

- ① if  $x$  is a terminal,  $\text{first}(x) = \{x\}$
- ② if  $x \rightarrow a\alpha$  where  $a$  is a terminal, add  $a$  to  $\text{first}(x)$
- ③ if  $x \rightarrow \epsilon$ , add  $\epsilon$  to  $\text{first}(x)$
- ④ if  $x \rightarrow \gamma_1 \gamma_2 \dots \gamma_{i-1} \gamma_i \dots \gamma$



1/2/8

1/2/8

1/2/8

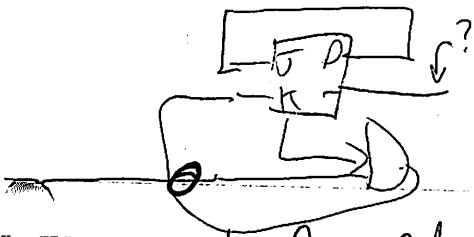
1/2/8

1/2/8

1/2/8

1/2/8

1/2/8



~~STH SA~~

8/30

\* looks like same in old as in use.

struct sym \*p → type of \*p is ~~sym~~ sym

(\*p).s\_name ≡ p → s\_name

p → x (\*p).x

Follow (A) = terms that can follow

"Don't panic, the data in this class hasn't been stored"

"Well, let's just worry about it until we figure it out"

struct symbol \* (\*f)();

ptr to struct returns a ptr to a symbol

\* (\*f)();

pointer to (\*f)();

pointer to function ret (\*f)

ptr to func returns a ptr to a symbol.

18

18

18

18

18

18

18

18

18

18

18

18

9/11

New cc.h csym.h

O-ADDR - removed

(struct sym \*)

char \*alloc;

refs are  
~~in~~  
~~to~~  
globals  
static  
travels

can read  
rules file  
to get global  
& static vars

P = (struct symbol \*) malloc (sizeof (struct sym))

↑  
symbol pointer

↑  
type cast

↑  
char point #bytes

convert to ~~the~~ type.

"A type conversion fcn."

"Have you tried it? It wontes for me."

~~Don't~~

Source lay eyes my ~~now~~

global defs + funcs first,  
typed func later

HIP

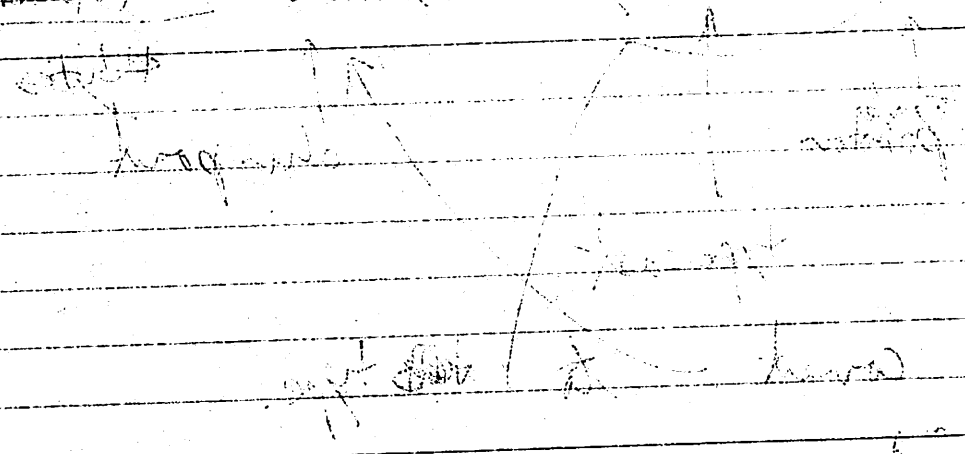
W. J. ...

Journal - 1911

...  
...  
...

...  
...

(...)



...

...  
...

9/24

"We need to order some  
chalk that doesn't have lead"

size:

int a → 36

char s[63] → 54

not used for functions

x & p is int

x p is pointer to int

p is pointer to p to int

for each - csync. loops through  
symbol table. also see day 1

11/10/19

group of 10 people  
last week had 10 people

10 people  
10 people  
part of the group

10 people  
10 people  
10 people

10 people  
10 people  
10 people

9/6T

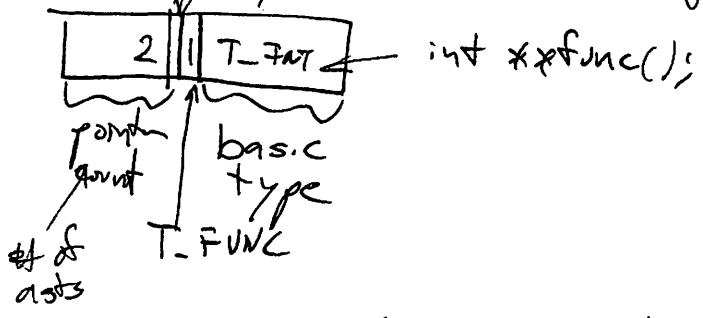
install

```

int f(a, b)
int a;
long x;
} lookup
if any params modified

```

← S\_ADDR - T\_LVAL from an lval



MKTYPE(2, T\_FUNC) | T\_FUNC

Array is a pointer type  
 Size field not used for params  
 Size field for array param is ptr size

int \*X[10]; - 10 \* PTR\_SIZE

char \*S[8] - 7 \* PTR\_SIZE

REF flag - symbol is referenced



10/10

10/10  
10/10  
10/10

10/10  
10/10

10/10  
10/10  
10/10

10/10  
10/10

10/10  
10/10

10/10  
10/10

10/10  
10/10

10/10  
10/10

10/10  
10/10

10/10  
10/10

9/6T

Must stack loop heads, but  
recursion should handle it.

Added body any trouble because  
addy loop change stack  
frame size

Td/P

that's about good state to all  
to all and that's name

and that's about good state to all  
to all and that's name

~~x[i]~~  
\*(x+i)

9/8H

\*f() += a;

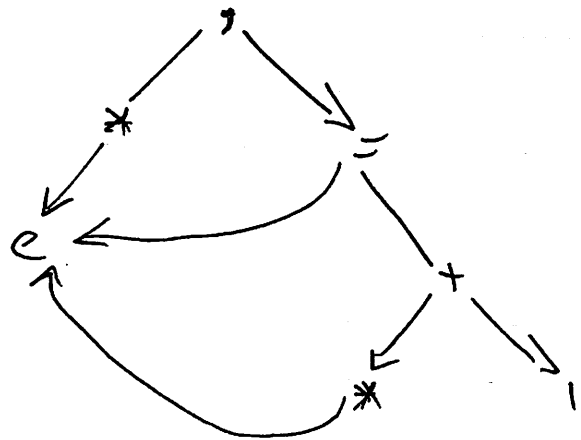
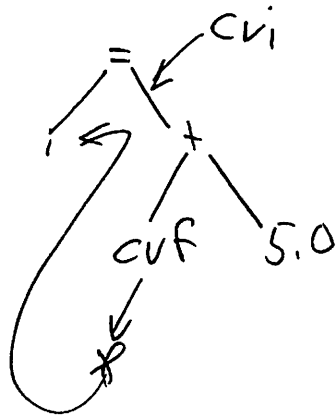
e1 op = e2

return e\_result field rather than walking the tree

enode does type checking + conversion of lval to rvals

+ + e  $\leftrightarrow$  e + = 1  
↑  
lvalue

i += 5.0



48' P

$(i+1)^*$

$p = + (1) 7 *$

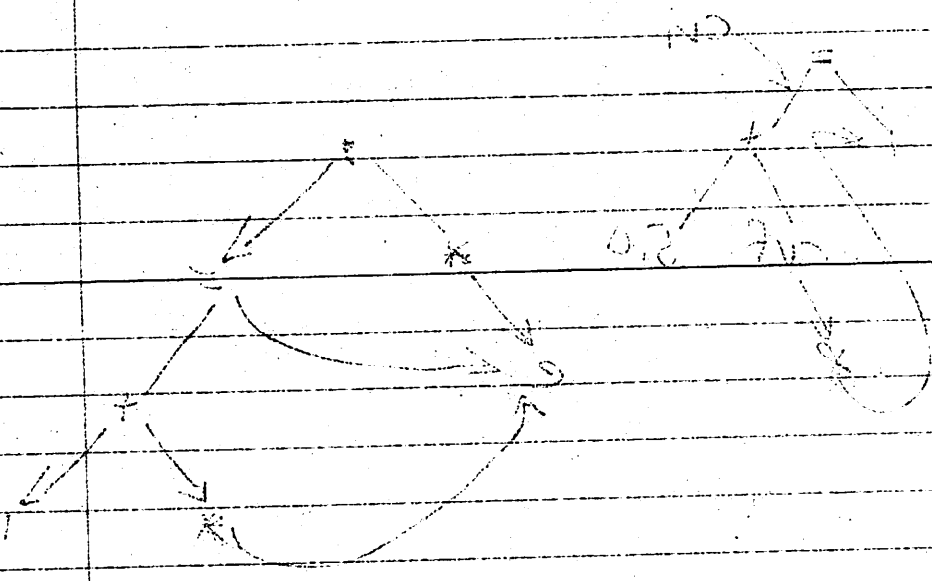
$\int 9 = 90 \quad 19$

method used then ...  
with all ...

...  
...

$12 + 9 \rightarrow 21$

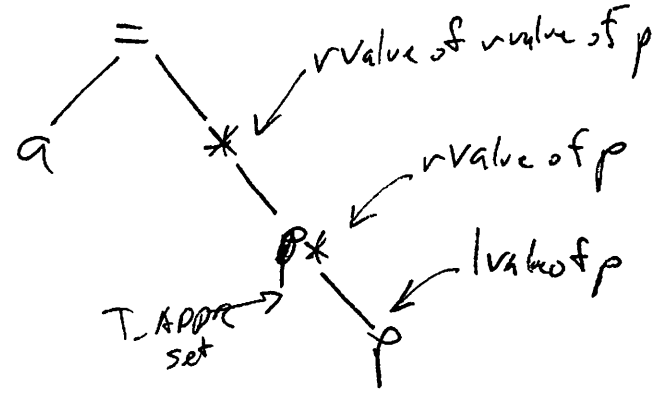
$6 \times 2 = 12$



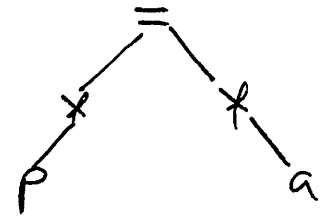
9/13T

left to right or right to left  
eval order depends on  
dir. of stack growth

$$a = *p$$



$$*p = a$$



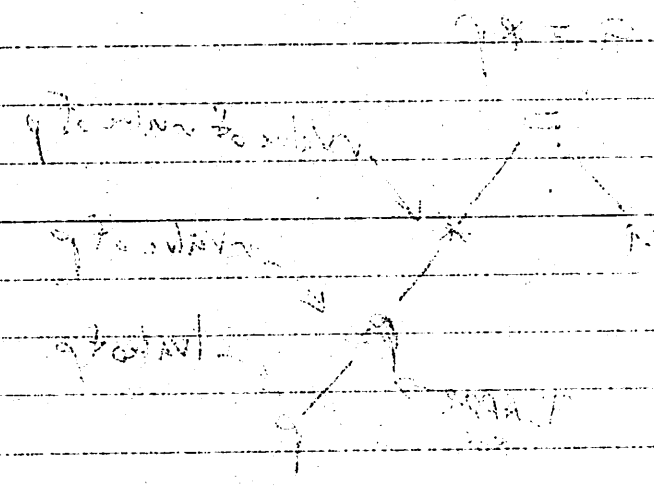
~~wants rvalue which is~~

See 63

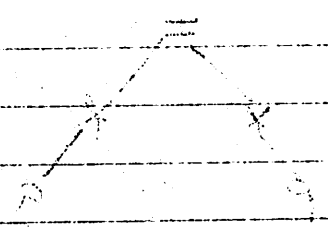
lvalue(-) - so value is to check 0-ADDR

21/11/21

Full of paper in hand of the  
 a paper should be with the  
 through the data of the



$p = 9 \times x$



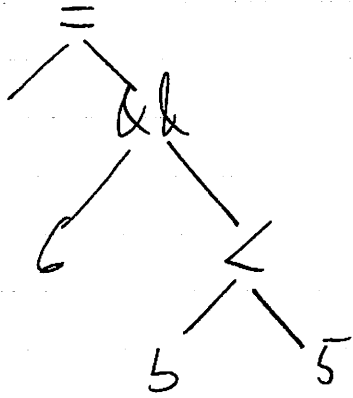
the data of the

$C = 2 \times 5$

Since the data of the

19/13T

$a = 6 \ \&\& \ b = 5$

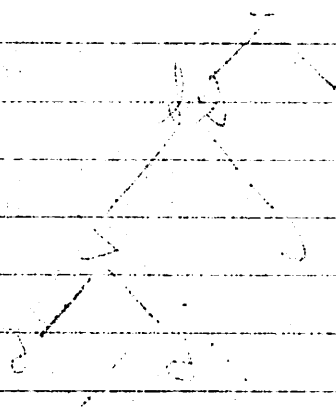


• assignment of const. values



1/1/11

1/1/11



1/1/11

$x > y$   $2 > 1$

$x > y$   
9/15/14

50% day for 3 days

Tues <sup>midnight</sup> at midnight

Three files named

#include "ccl.h [1033, 047113]"

"This is probably the hardest assignment"

Start coding on Tuesday

Augmented assignment

$\rho = \text{cost}$

~~and~~ no runtime representation

"If handles cost - only error handling."

Counts. Prizes on interest stuff, not error handling.

IS handles all progress w/o free work, don't need freesty of nodes

Feb 2  
1917

Spoke of ...  
...  
...

"The ..."  
"The ..."  
...  
...

...  
...  
...  
...

...  
...  
...  
...

"All of the biggies will  
be here - drinky - that's  
the major occ. at these t'gs"

"I don't go to talks, I just  
go to the bar."

FCS

9/20T

REF - a bit to tell if the id was referred to.  
REF on + DEF off - undefined

Ref a bad address

Ref through 0  $\rightarrow$  1  $\frac{SP}{8}$   $\rightarrow$  Boom

$a < 1$ ;  $\rightarrow$  Value expected

$f(a < 1)$ ;

Not in test1-6.c - not see it.

our code - 2, 3 times larger than optimal

"Nothing worse than fast code that doesn't work"

op arg1, arg2 - a two address machine

add + store - three address

segment - contiguous region of code or data

a.v. - local data + return information

Task 10

Let  $\mu$  be the mean of a normal distribution with variance  $\sigma^2$ . We want to test  $H_0: \mu = \mu_0$  against  $H_1: \mu > \mu_0$ . The test statistic is  $T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$ . The power function is  $\beta(t) = P(T > t | \mu)$ .

Let  $\mu_1 > \mu_0$  and  $\alpha = P(T > t_0 | \mu_0)$ . Find  $\beta(t_0 | \mu_1)$ .

$$\beta(t_0 | \mu_1) = P\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} > t_0 \mid \mu = \mu_1\right)$$

$$= P\left(\bar{X} > \mu_0 + t_0 \frac{\sigma}{\sqrt{n}} \mid \mu = \mu_1\right)$$

$$= P\left(\frac{\bar{X} - \mu_1}{\sigma/\sqrt{n}} > \frac{\mu_0 - \mu_1 + t_0 \frac{\sigma}{\sqrt{n}}}{\sigma/\sqrt{n}}\right)$$

Let  $Z = \frac{\bar{X} - \mu_1}{\sigma/\sqrt{n}}$ . Then  $Z \sim N(0, 1)$  under  $H_1$ .

$$\beta(t_0 | \mu_1) = P\left(Z > \frac{\mu_0 - \mu_1 + t_0 \frac{\sigma}{\sqrt{n}}}{\sigma/\sqrt{n}}\right)$$

Let  $\delta = \frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}$ . Then  $\beta(t_0 | \mu_1) = P(Z > t_0 - \delta)$ .

$$= 1 - \Phi(t_0 - \delta)$$

$$= \Phi(\delta - t_0)$$

Since  $\alpha = P(T > t_0 | \mu_0) = 1 - \Phi(t_0) = \Phi(-t_0)$ .

$$\beta(t_0 | \mu_1) = \Phi\left(\frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}} - t_0\right)$$

9/20T

Ex: [see opcode. h]

move 2, x ; load x into 2  
add 6, 7 ; add reg 7 to reg 6  
movem 13, 16(14) ; store reg 13  
addi 6, 40 ; add 40 reg 6  
move 6(7) ; load 6 in word ptr at by 7

addi 6, 040(7) - 6 = 740

DEC has 502

S\_offset will hold frame offsets

(10/17)

W. George ...

$f(x) = x^2 + 2x + 1$   
 $f'(x) = 2x + 2$   
 $f''(x) = 2$   
 $f(0) = 1$   
 $f'(0) = 2$   
 $f''(0) = 2$

002 ...

... ..

9/22/4

calling  $\{2,1\}$  - call sequences



AS/P

Handwritten notes on lined paper, including a circled "AS/P" and some illegible text.

9/27T

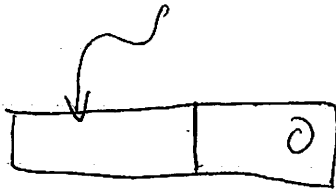
"If that's a C.C. cookie, you damn well better give me one!"

#undef AND

#undef LSH

#include "opcode.h [055, 2]"

MOVSI 6, left-18



1/2/19

1. The first part of the book is...

Chapter 1

The first part of the book is...

Chapter 2

Chapter 3

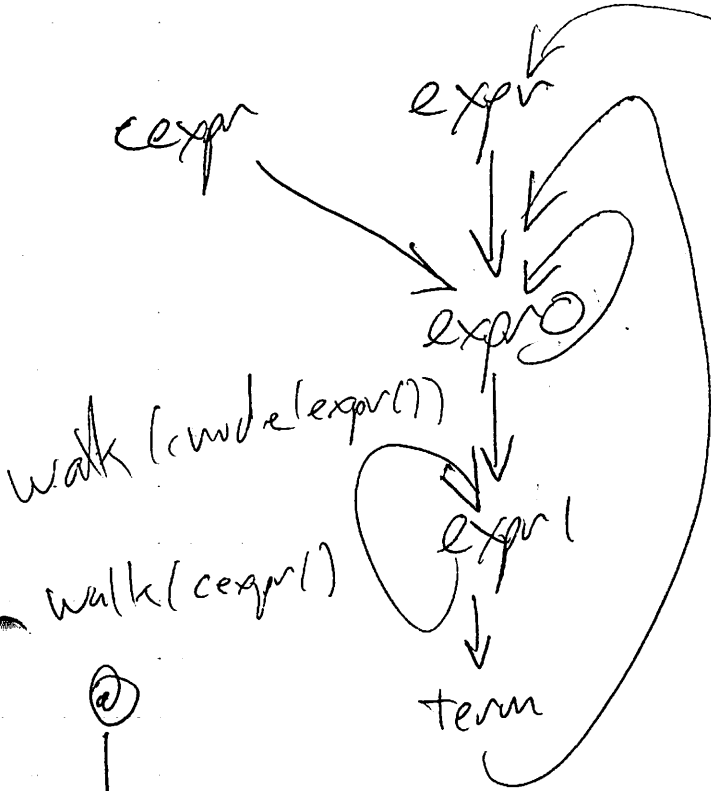


"Will look at ~~that~~ the next set of dis → "I'll pass out air sockets today"

"I'll like to find the SOB at DEC that put this on the machine" (ADJBP)

9/29/84

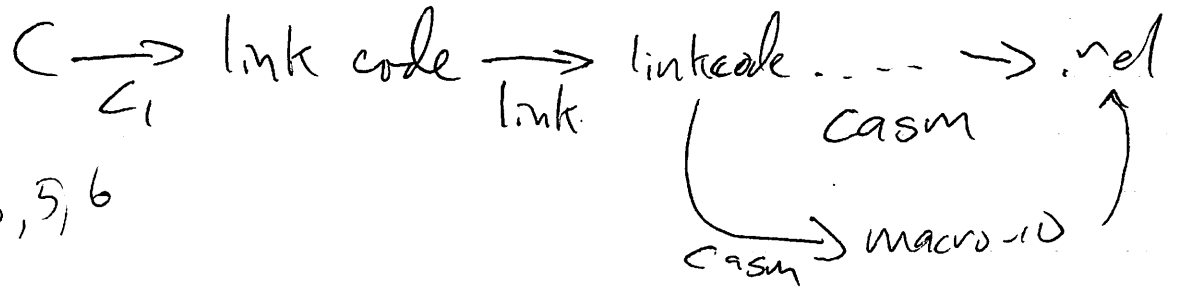
(e.g.p.c. →



return node(0, 1, int, expr(), 0))

Exam on Tuesday  
No q's about reg alloc.

when node(0, 0, code(cexpr()), int) 0 2, 3, 5, 6



cl < test2.c | linker | casm > Check casm output  
ccomp test2.mac vtblib, rel[21]  
i.exe -T 1020

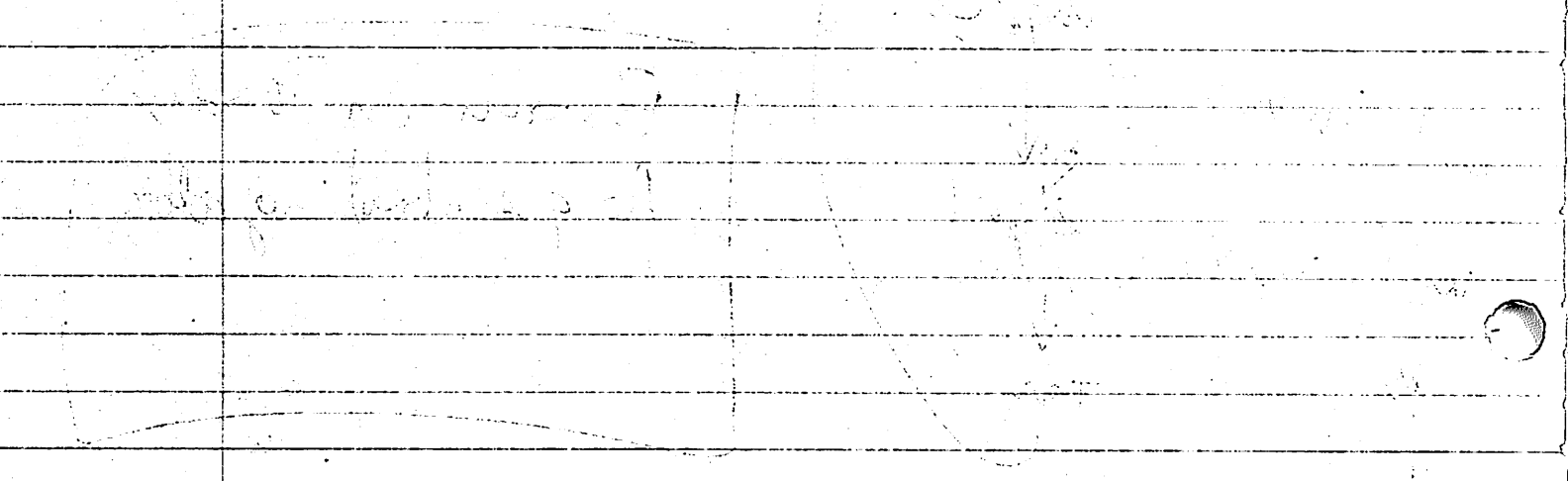
```
main()
{
  puts("hello world\n");
  puts(s);
  done #3;
}
```

```
while (*s)
  putc(*s++);
```

Handwritten notes at the top of the page, including a circled number '25'.

Handwritten notes in the second section of the page.

Handwritten notes in the third section of the page.



Handwritten notes in the fifth section of the page, including arrows and a circular diagram.

Handwritten notes in the sixth section of the page.

Handwritten notes at the bottom of the page.

Link

10/6/4

\* How could overlays be added?

"So the overall goal of this other than to parody 455 may be to..."

"link editor" implies more than "loader"

UNIX: exec code  $\rightarrow$  object code  
w/ no unbracketed addresses

All expressions represent  
one ~~cell~~ cell

move 0100 + 15 +  
%x 02 -

0  
0206

lint assumes each input expression

loader does  
9 } 18 bit  
5 }

via (a << 15) + b

lint replaces symbol by these values

1/11/11

1. The first part of the book is a history of the world from the beginning of time to the present day. It is written in a simple and easy-to-understand style.

2. The second part of the book is a collection of stories and legends from different cultures. These stories are often very interesting and provide a glimpse into the lives of people in the past.

3. The third part of the book is a collection of poems and songs. These are often very beautiful and provide a glimpse into the lives of people in the past.

4. The fourth part of the book is a collection of facts and figures. These are often very interesting and provide a glimpse into the lives of people in the past.

5. The fifth part of the book is a collection of recipes and cooking instructions. These are often very simple and easy to follow.

6. The sixth part of the book is a collection of puzzles and games. These are often very fun and provide a good way to pass the time.

3c      Sc

10/11T

Ex prob graded on 0-75 basis

~~Simple Opt~~  
~~Calcutting Paper~~  
~~AR Head~~  
~~Dist taught~~  
~~Wbt Word VA~~ STATICS.  
~~Reps 0-5 on VA~~  
~~2000~~  
 Yarn TR      -1  
~~2000~~      83-1  
~~Lab exercises~~  
~~As 5000~~

~~proc~~  
~~proc~~      ~~proc~~ (bp)  
~~proc~~      ~~proc~~ (bp)  
~~proc~~      ~~proc~~ (bp)

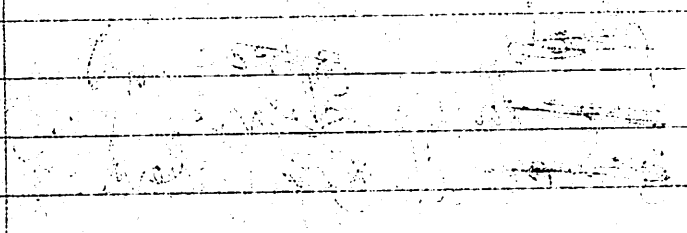
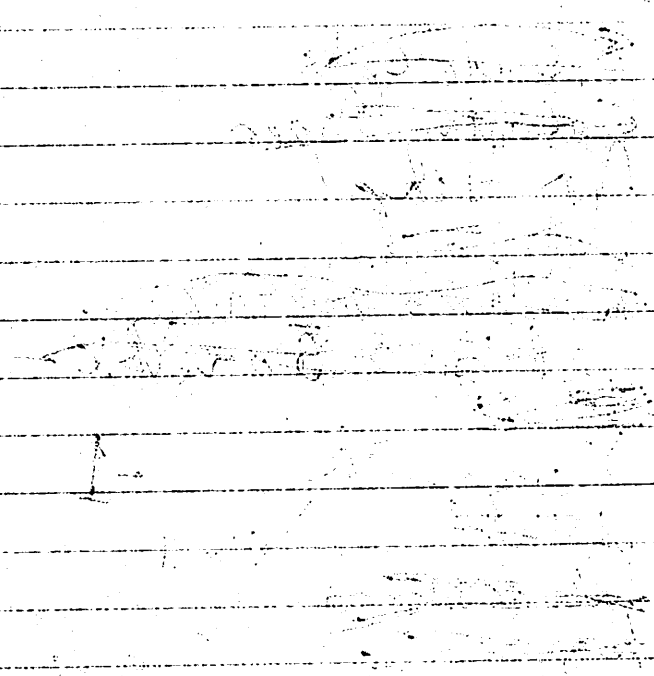
C2      ~~200~~      -d  
             -m



7/1/10

2008

Handwritten notes at the top of the page, possibly including a date or title.



Handwritten text or notes in the lower right quadrant of the page.

~~At 5:00~~  
No Class 10/20 #

10/13 #

Elegant link for computer writers, blurry link from users point of view

link [file1 -v| -o file]

Issue Warning ~~error~~ if ~~error~~ can't find an input file

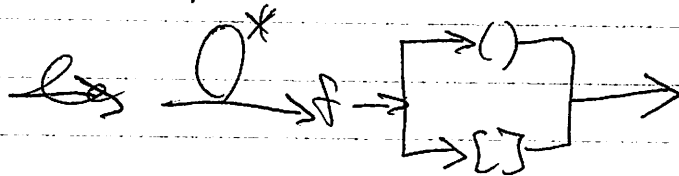
link a -v b -v c -v -o foo

~~at 5:00~~  
take last output file

Value of segment name  $\rightarrow$  base location of segment

bss - "block started by symbol"

UCP AR  
ARCW  
550a  
550b  
473  
520  
453  
700



4/5/01

... ..

... ..

... ..

... ..

... ..

... ..

... ..

... ..

... ..

bdd

"You alright boy?" - look to bdd

10/18

"For those of you who haven't heard about it, you've probably

It was a L home grown ~~here~~ used to be used in some of the classes.   
 ~~here~~   
 ~~in some of the classes.~~   
 ~~probably~~   
 ~~locky~~

Went to univ until poss 2.

~~Fixed~~  
Always give error for two defines  
w/ same name

70/100

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

Handwritten notes, possibly describing a process or experiment.

"~~Comments~~ What are comments?"

"Comments are for humans"

10/25T

when -u is seen, output undefined symbols

⓪

LOS, 213 -d b.lc a.lc

-d don't delete temp files for segment

don't worry about names > 6 chars

A maze of segment ptrs

Don't worry about evaluation until end

10/23/21

10/23/21

10/23/21

10/23/21

10/23/21

10/23/21

10/23/21

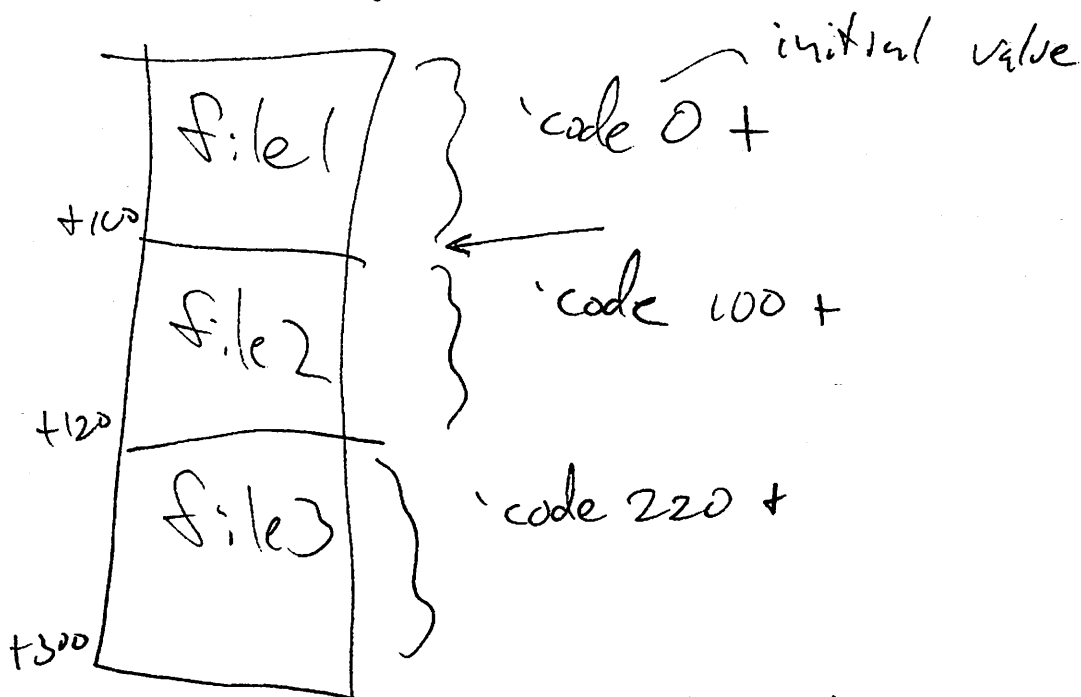
10/27H

~~a code +~~

```
.def a b 2 +  
.def b code 4 + → b = 'code 4 +  
a code + 2 +  
b code + 2  
'code 4 +
```

code can't be defined until no files are left.

Can't run dump after value has been substituted for code



= S-dump symbol table



10/11

+ abcd

+ 12 d  
+ f d o d h  
+ 1 d o d  
+ 1 d o d  
+ 1 d o d

with the help of these

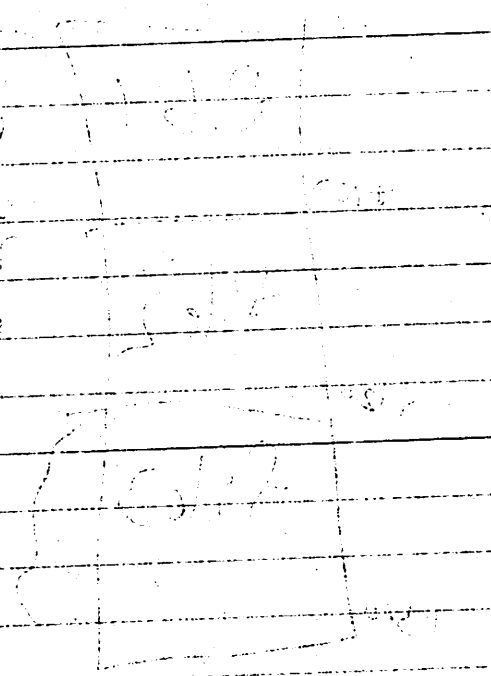
we can find the

also

+ 1 d o d

+ 1 d o d

+ 1 d o d



with the help of these

(only)

10/27H

DEF bit on segs for -u  
 no point to evaluate stuff in pass 1  
 since no eval in  
 values of code

'code 0 +

evaluate definitions  
 pass 1 evals definition

a 2 + 4 +

Do only test subs, but no  
 arith. evals.

Use octal for output

```
def a code
ses code
```

1

2

```
.len code 2
ses code
```

---

```
9
code
```

.len code 2

4/5/00

in the area of the 1991

major effect at the time of

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

the area of the

II/IT

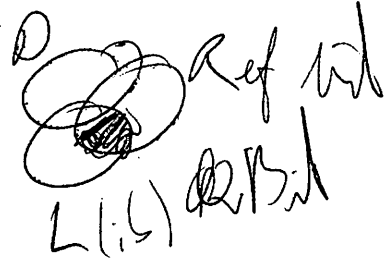
BUS:

058 2007--6

Library contains entries, each of which is an object file w/ definitions

$L =$  set of symbols defined by all entries in library

Ref lib



would load  $(L \cap R) - (D \cup S)$

needs  
hard access  
or mult  
passes

while  $((L \cap R) - D)$  is not empty  
for each entry  $q$   
if  $q$  defines a symbol in  $(L \cap R) - D$   
pass  $q$

expunge symbols in  $L - (R \cup D)$

↑ symbols only in L

7/11

What is the difference between  
a simple and a complex sentence?

A simple sentence has one  
independent clause. A complex  
sentence has at least one  
independent clause and one  
dependent clause.

Example: The cat sat on the mat.  
The cat sat on the mat when the  
owner came home.

Identify the independent and  
dependent clauses in the  
following sentences.

1. Although it was raining, the  
game went on.

2. She finished her homework  
before she went to bed.

3. The car started when the  
engine was turned on.

RL/IT

add:

struct symbol &S-link  
to symbol structure

#define LIB 020

```
struct entry {
    char *e_name
    long e_pos
    int e_size
    int e_cont
    struct entry *e_link
}
```

Sum of all cont fields is # of cont symbols

loc in lib  
beg of entry  
# of sym in (L/R) - 0

# -h- index size  
# symbol name

\* { Test Q: @

Q: Can't things other than object text be in library

VBS on OS is byte stream

And when we're where we want to know where are ...

Test Q: what can be done to improve the index, size + pos info?

11/11

July 23rd 1954

Dear Mother  
I received your letter of the 21st and was  
glad to hear from you. I am well and  
hope these few lines will find you the same.  
I am still in the hospital and will be  
discharged in a few days.

I am sure you will be glad to hear  
that I am getting better. I will be home  
in a few days and will be able to see  
you again.

I love you very much and hope to see  
you soon. Write soon.  
Your affectionate son,  
John Doe

U/LT

might want to def endo struct  
after it's built

two phase search, 2nd phase  
is only used for more  
topological sorted fits

will tell  $\phi$  how far  
to read

Must recompute (LR) - D after  
each whole in pass

Can't do algebra on #'s returned  
by step



TW

Handwritten notes at the top of the page, possibly including a title or date.

Handwritten notes in the middle section of the page.

Handwritten notes in the lower middle section of the page.

Handwritten notes in the bottom middle section of the page.

Handwritten notes at the bottom of the page.

11/3H

fseek will be there Tuesday

.copy = dskd: [1000, 102, cbout] st drah  
cli5, rel  
root, rel

.as dsk csc

| 1-4.1c

.lib - various libraries

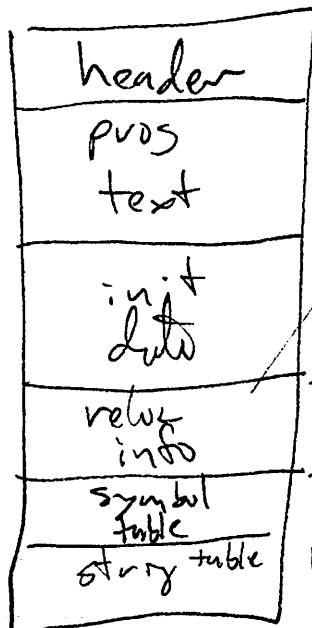
Exam - everything + anything about ~~lib~~  
linking + library searching

No Code Gen. Just link + load on test.

The ans. to how it's done in  
the real world is that it's  
done wrong.

Print some  
h files

UNIX OBJECTS:



used segment names  
to ~~code~~ cause rels.  
i.e. our seg. names ~~at~~  
were our rels info.

} got  
} got  
} got

115/11

put out with on the street

what's the so-called job?

low die

low die

see you go.

what's the job?

what's the job of the...  
the so-called job?

test me... what's the job?

what's the job of the...

the job of the...

what's the job of the...

what's the job of the...

what's the job of the...

what's the job of the...

what's the job of the...

what's the job of the...

low

low

low

low

low

low

low

low

low

symbol table indexes into string table

string table is only for identifiers

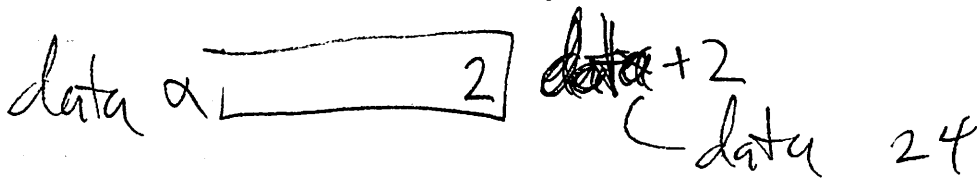
.o + bin files have the same format

If no relocation info, file can be executed

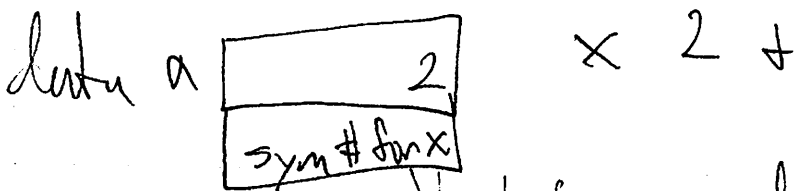
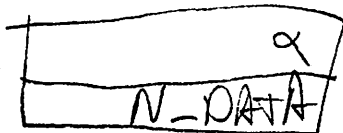
strip removes reloc, symbol, str stuff

Magic #'s are 4xx because they can be identified by first byte

\* Might think about how to do common in Link.



reloc



get value for  $\alpha \times$  + add to  $\alpha$

every occurrence of  $\alpha$  has a relocation word,  $\alpha \times$  to it.

1/11/20

In each of the following  
cases, find the value of x

1. In the figure below, find the value of x

2. In the figure below, find the value of x

3. In the figure below, find the value of x

4. In the figure below, find the value of x

5. In the figure below, find the value of x

6. In the figure below, find the value of x

7. In the figure below, find the value of x

8. In the figure below, find the value of x

9. In the figure below, find the value of x

10. In the figure below, find the value of x

11. In the figure below, find the value of x

12. In the figure below, find the value of x

13. In the figure below, find the value of x

11/34

Want to avoid having to relocate upon  
expectation  
— Somewhat about 1k, 8k boundaries

DEC-10

VMX format is somewhat similar  
to bus

file is seq of blocks:

0	type code, leg not cont overl
1	18 2 bit byte hunks
2-19	16 <del>leg</del> words of type-dep data

REQUEST LIBRARY block directed  
link to search library

Rad 50 - 50 is outal

"If I ever catch up of you doing  
this, I'll take your degree away."

"It is well known that LC are faster  
to type because they aren't as big."

18/11

Wages & salaries in manufacturing industry  
about 18% of GDP

Wages & salaries in services industry  
about 25%

Wages & salaries in construction industry  
about 10%

Wages & salaries in agriculture industry  
about 5%

Wages & salaries in public sector  
about 15%

11/3

~~DL~~ Linking

do all by at exec time  
usually (best?) works w/  
hardware support

DL can be done on per-object  
file basis



2111

2005.11.14

out on the road

in about 15 minutes

high up in the air

11/8T

Didn't touch his pass except for  
making sure it didn't read too far  
when symbols are <sup>defined</sup> referenced

R fread after my pass on  
each entry.

- l - list of everything defined in a module

\* Common question for exam"  
4 questions - closed book

"Heisenberg -- If you touch something, you're probably  
going to break it"

\* How about using debugger on  
the debugger?

Symbol name conflict

\* Making debuggers robust is very hard



Cyber

1/11

Handwritten notes, possibly starting with "The first..."

Handwritten notes, possibly starting with "The second..."

Handwritten notes, possibly starting with "The third..."

Handwritten notes, possibly starting with "The fourth..."

Handwritten notes, possibly starting with "The fifth..."

Handwritten notes, possibly starting with "The sixth..."

Handwritten notes, possibly starting with "The seventh..."

Handwritten notes, possibly starting with "The eighth..."

11/8T

400100p  
p/i

200217 - print  
400100 - sp

" - and queens are hard to place."

18

18

18

11/15T

Won't use ~~test~~ lib4.lib

Several undefined errors

db is a throw-away program

? - bad address

?? - bad command

task is to write db1

db1 should give "tty:" in "vt"

1/11

1/11/11

1/11/11

1/11/11

1/11/11

1/11/11

1/11/11

11/17H

"If we needed seabeth last time  
we need answers logs today."

-switch - j command + at very start  
C version debugger → debug

Now - call Emacs from Iron  
Fix? to call ExecuteBand



ATLANTA

Mr. J. B. ...

...

...

11/22T

set/pit - set flag ~~add~~  
dbi - init addr <sup>address</sup> 0

test2.val rtf16.val db.val

setup struct area  
compute two ind

~~use atoi~~

change operand is out/led  
all #'s?

atoi(xcp)

~~setadistcp)~~  
~~char xcp;~~

atoi(xcp)  
char xcp;

~~int base~~  
int n = 0; base  
char \*p = xcp; hc  
~~char hc~~  
while (!isspace(\*p)) p++

~~\*(xcp)~~  
base = \*xcp == '0' ? 8 : 10;  
hc = '0' + (base - 1);

11/5/11

the job is: 1/10/11  
of all the time

work had to be done

the time was spent  
in the field

the time was spent

in the field

the time

was spent

in the field

the time

was spent

in the field

the time

was spent

accdS

11/225

Focus is used to guide interpretation  
of symbols.

\* EA → break part spec?

could map draw #'s to strings, @  
(at not my opp.)

TSSW

the independent thing at least for my part  
is to be done

long long time ←

all about it a few years ago  
I applied for it

11/29

4001003  
b

How could vars be added at  
run time?

Should uninterpreted code in lieu  
of w/o it.

RT Symbol table could stay in disk;  
good for large sym tables. 17

How to strip named segments  
from lib you.

How can this be used to make  
a linked list.

Some integration question.

Mode to system

RS/III

10-10-44

to follow at ...

and ...

17

... ..

... ..

... ..

... ..

12/15

code used by edb is prealgebra  
(might have bypassed gallery but  
that seems unlikely to me.)

"I'll take off points if you  
forget to turn, unless you  
find something major with white  
base, my eye will be so  
banned"



7/21

... ..  
... ..  
... ..

... ..  
... ..  
... ..  
... ..

12/64

We have some vacant TA slots -  
if you're interested, you can  
pick up ~~the~~ ball and chain at the office

ED Read problem 4 of first exam  
+ answer

Exam: only db, cdb

10:30 am Dec 16

Closed book, closed notes

"Obviously this class presents a  
challenge to both student  
+ teacher."

How could *cozylike* be less painful  
"This is basically a two pass class."

10/10  
The first thing I noticed  
was the smell of the  
ocean.

The second thing I noticed  
was the sound of the  
waves crashing against the rocks.

The third thing I noticed  
was the taste of the  
salt water.

The fourth thing I noticed  
was the feel of the  
sand under my feet.

The fifth thing I noticed  
was the sight of the  
sun setting over the horizon.

The sixth thing I noticed  
was the sound of the  
seagulls flying overhead.

Registers

sp

dp

ap

pk - trick - last

Handwritten scribble, possibly "Handwritten"

Handwritten scribble, possibly "Handwritten"

~~New structure~~

~~#13~~ 12/13

3. ~~O\_CALL~~ in

{, 2, 3} RTL IN C

Unit  
switch

O\_CALL  
a PARM  
ARG

~~find~~

~~find call~~ ARG

~~for~~  
parse  
eval

f(1, 2)

(x f)()

(x f)(1, 2, 3)

O\_CALL

func

args

25%

build stack frame  
fix return ~~the~~ address  
on return get no

to switch / fix up.

switch (func, substate)

switch (substate, junk)

---



12/13

⊙ Breakpoint - integrated



511

logarithmic distribution