

Generating Code and Running Programs

COMS W4115

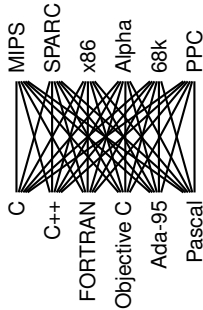


Prof. Stephen A. Edwards
Fall 2006
Columbia University
Department of Computer Science

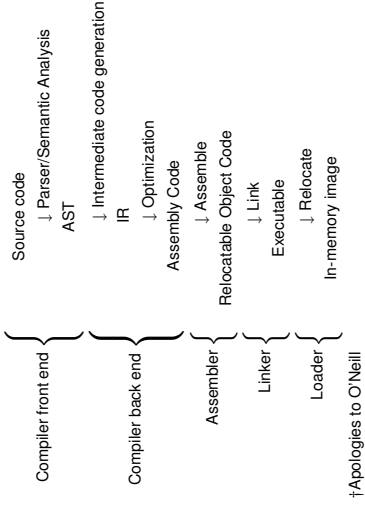
Portable Compilers

Building a compiler a large undertaking; most try to leverage it by making it portable.

Instead of



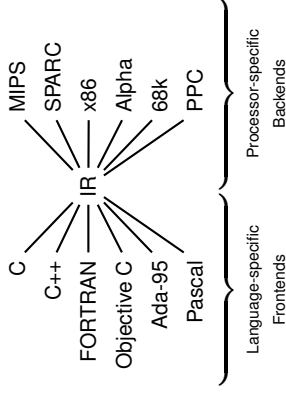
A Long K's Journey into Byte†



†Apologies to O'Neill

Portable Compilers

Use a common intermediate representation.



Compiler Frontends and Backends

The front end focuses on analysis:

- lexical analysis
- parsing
- static semantic checking
- AST generation

The back end focuses on synthesis:

- Translation of the AST into intermediate code
- optimization
- assembly code generation



Intermediate Representations/Formats

Stack-Based IRs

Advantages:

- Trivial translation of expressions
- Trivial interpreters
- No problems with exhausting registers
- Often compact

Disadvantages:

- Semantic gap between stack operations and modern register machines
- Hard to see what communicates with what
- Difficult representation for optimization

Stack-Based IR: Java Bytecode

```
int gcd(int a, int b) {
    while (a != b) {
        if (a > b)
            a -= b;
        else
            b -= a;
    }
    return a;
}

# javap -c Gcd
Method int gcd(int, int)
  0 goto 19
  3 iload_1 //Push a
  4 iload_2 //Push b
  5 if_icmple 15 //if a <= b goto 15
  8 iload_1 //Push a
  9 iload_2 //Push b
  10 isub //a - b
  11 istore_1 //Store new a
  12 goto 19
  15 iload_2 //Push b
  16 iload_1 //Push a
  17 isub //b - a
  18 istore_2 //Store new b
  19 iload_1 //Push a
  20 iload_2 //Push b
  21 if_icmpe 3 //if a != b goto 3
  24 iload_1 //Push a
  25 ireturn //Return a
```



Register-Based IR: Mach SUIF

```
int gcd(int a, int b) {
    while (a != b) {
        if (a > b)
            a -= b;
        else
            b -= a;
    }
    return a;
}

gcd_gcdtmp0:
    sne $vr1, $s2 <- gcd, a, gcd, b
    seq $vr0, $s2 <- $vr1, $s2, 0 //!(a != b) goto tmp1
    brne $vr0, $s2, gcd_gcdtmp1 //!(a != b) goto tmp1
    sll $vr3, $s2 <- gcd, b, gcd, a
    seq $vr2, $s2 <- $vr3, $s2, 0 //!(a < b) goto tmp4
    brne $vr2, $s2, gcd_gcdtmp4 //!(a < b) goto tmp4
    mtk 2, 4 //Like number 4
    sub $vr5, $s2 <- gcd, b, gcd, a
    mov gcd_gcdtmp2 <- $vr5, $s2
    mov gcd, a <- gcd_gcdtmp2 // a = a - b
    jmp gcd_gcdtmp5
gcd_gcdtmp4:
    mtk 2, 4
    sub $vr5, $s2 <- gcd, b, gcd, a
    mov gcd_gcdtmp3 <- $vr5, $s2
    mov gcd, b <- gcd_gcdtmp3 // b = b - a
    jmp gcd_gcdtmp5
gcd_gcdtmp1:
    mtk 2, 8 //Return a
    ret gcd, a
```



Register-Based IRs

Most common type of IR

Advantages:

- Better representation for register machines
- Dataflow is usually clear
- Disadvantages:
 - Slightly harder to synthesize from code
 - Less compact
 - More complicated to interpret

Typical Optimizations

Folding constant expressions

1+3 → 4

Removing dead code

if (0) { ... } → nothing

Moving variables from memory to registers

```
ld [%fp+68], %i1
sub %i0, %i1, %i0 → sub %o1, %o0, %o1
st %i0, [%fp+72]
```

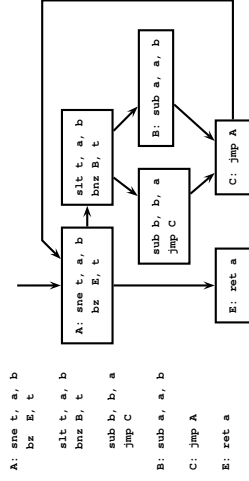
Removing unnecessary data movement

Filling branch delay slots (Pipelined RISC processors)

Common subexpression elimination;

Control-Flow Graphs

A CFG illustrates the flow of control among basic blocks.



Optimization

```
int gcd(int a, int b) {
  while (a != b) {
    if (a < b) b -= a;
    else a -= b;
  }
  return a;
}
```

First version: GCC on SPARC
Second version: GCC-O7



Basic Blocks

```
int gcd(int a, int b) {
  while (a != b) {
    if (a < b) b -= a;
    else a -= b;
  }
  return a;
}
```



The statements in a basic block all run if the first one does.
Starts with a statement following a conditional branch or is a branch target.
Usually ends with a control-transfer statement.

Machine-Dependent vs. -Independent Optimization

No matter what the machine is, folding constants and eliminating dead code is always a good idea.

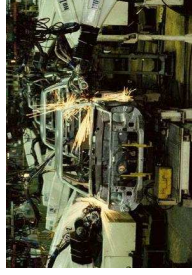
```
a = c + 5 + 3;
if (0 + 3) {
  b = c + 8;
} → b = a = c + 8;
```

However, many optimizations are processor-specific.
Register allocation depends on how many registers the machine has

Not all processors have branch delay slots to fill

Each processor's pipeline is a little different

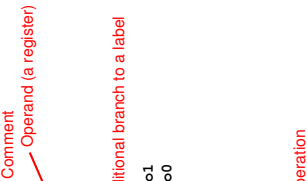
Assembly Code and Assemblers



Assembly Code

Most compilers produce assembly code: easier to debug than binary files.

```
! gcd on the SPARC
gcd:
  cmp %o0, %o1
  be .LL8
  nop
  ble,a .LL2
  sub %o1, %o0, %o1
  sub %o0, %o1, %o0
.LL2:
  cmp %o0, %o1
  bne .LL9
  nop
.LL8:
  retl
nop
```



Role of an Assembler

Translate opcodes + operand into byte codes

```

Instruction code
Address ↓
0000 80A20009    cmp    %o0, %o1
0004 02800008    be     .LL8
0008 01000000    nop
.LL9:
000c 24800003    ble,a .LL2
0010 92224008    sub   %o1, %o0, %o1
0014 90220009    sub   %o0, %o1, %o0
.LL2:
0018 80A20009    cmp   %o0, %o1
001c 12BFFFFC    bne  .LL9
0020 01000000    nop
.LL8:
0024 81C3E008    retl
0028 01000000    nop

```

Encoding Example

sub %o1, %o0, %o1

Encoding of "SUB" on the SPARC:

10	rd	000100	rs1	0	reserved	rs2	4
31	29	24	18	13	12		

rd = %o1 = 01001
rs1 = %o1 = 01001
rs2 = %o0 = 00100

10 01001 000100 01001 0 00000000 01000
1001 0010 0010 0010 0100 0000 0000 1000
= 0x922228004

Role of an Assembler

Transforming symbolic addresses to concrete ones.

Example: Calculating PC-relative branch offsets.

```

000c 24800003    ble,a .LL2
0010 92224008    sub   %o1, %o0, %o1
0014 90220009    sub   %o0, %o1, %o0
.LL2:
0018 80A20009    cmp   %o0, %o1

```

LL2 is 3 words away

Role of an Assembler

Most assemblers are "two-pass" because they can't calculate everything in a single pass through the code.

```

.LL9:
000c 24800003    ble,a .LL2
0010 92224008    sub   %o1, %o0, %o1
0014 90220009    sub   %o0, %o1, %o0
.LL2:
0018 80A20009    cmp   %o0, %o1
001c 12BFFFFC    bne  .LL9
0020 01000000    nop
.LL8:
0024 81C3E008    retl
0028 01000000    nop

```

Don't know offset of LL2

Role of an Assembler

Constant data needs to be aligned.

```

char a[] = "Hello";
int b[3] = { 5, 6, 7 };
.section ".data"
global a
.type a, #object
.size a,6
a: .asciz "Hello"
.global b
.align 4
.type b, #object
.size b,12
0008 00000005    .word 5
000c 00000006    .word 6
0010 00000007    .word 7

```

Assembler directives

Bytes added to ensure alignment

Role of an Assembler

The MIPS has pseudoinstructions:

```

li $t14, 0x12345abc
ori $t14, 0x5abc
lui $t14, 0x1234

```

expands to

Load the immediate value 0x12345abc into register 14."

Load the upper 16 bits, then OR in the lower 16"

MIPS instructions have 16-bit immediate values at most

RISC philosophy: small instructions for common case

Role of an Assembler

Eight "general-purpose" 32-bit registers:

```

eax ebx ecx edx ebp esi edi esp
esp is the stack pointer
ebp is the base (frame) pointer
addl %eax, %edx
movl 20(%ebp), %eax

```

Load word at ebp+20 into eax

Role of an Assembler

Where to put temporary results? Our compiler will just put them on the stack; a typical default.

```

int bar(int g, int h, int i, int j, int k, int l)
{
    int a, b, c, d, e, f;
    a = foo(g);
    b = foo(h);
    c = foo(l);
    d = foo(j);
    e = foo(k);
    f = foo(l);
    return a + (b + (c + (d + (e + f))));
}

```

Role of an Assembler

Optimization: Register Allocation

```

int bar(int g, int h, int i, int j, int k, int l)
{
    int a, b, c, d, e, f;
    a = foo(g);
    b = foo(h);
    c = foo(l);
    d = foo(j);
    e = foo(k);
    f = foo(l);
    return a + (b + (c + (d + (e + f))));
}

```

Role of an Assembler

Optimization: Register Allocation

```

.LL9:
000c 24800003    ble,a .LL2
0010 92224008    sub   %o1, %o0, %o1
0014 90220009    sub   %o0, %o1, %o0
.LL2:
0018 80A20009    cmp   %o0, %o1
001c 12BFFFFC    bne  .LL9
0020 01000000    nop
.LL8:
0024 81C3E008    retl
0028 01000000    nop

```

Know offset of LL9

Role of an Assembler

Optimization: Register Allocation

```

.LL9:
000c 24800003    ble,a .LL2
0010 92224008    sub   %o1, %o0, %o1
0014 90220009    sub   %o0, %o1, %o0
.LL2:
0018 80A20009    cmp   %o0, %o1
001c 12BFFFFC    bne  .LL9
0020 01000000    nop
.LL8:
0024 81C3E008    retl
0028 01000000    nop

```

Role of an Assembler

Quick Review of the x86 Architecture

```

Eight "general-purpose" 32-bit registers:
eax ebx ecx edx ebp esi edi esp
esp is the stack pointer
ebp is the base (frame) pointer
addl %eax, %edx
movl 20(%ebp), %eax

```

Role of an Assembler

Quick Review of the x86 Architecture

```

Eight "general-purpose" 32-bit registers:
eax ebx ecx edx ebp esi edi esp
esp is the stack pointer
ebp is the base (frame) pointer
addl %eax, %edx
movl 20(%ebp), %eax

```

Unoptimized GCC on the x86

```

movl 24(%ebp), %eax      % Get k
pushl %eax              % Push argument
call foo                % e = foo(k);
addl $4, %esp           % Make room for e
movl %eax, %eax         % Does nothing
movl %eax, -20(%ebp)    % Save return value on stack

movl 28(%ebp), %eax     % Get l
pushl %eax              % Push argument
call foo                % f = foo(l);
addl $4, %esp           % Make room for f
movl %eax, %eax         % Does nothing
movl %eax, -24(%ebp)    % Save return value on stack

movl -20(%ebp), %eax    % Get f
movl -24(%ebp), %edx    % Get e
addl %edx, %eax         % e + f
movl %eax, %edx         % Accumulate in edx
addl -16(%ebp), %edx    % d + (e+f)
movl %edx, %eax        % Accumulate in edx
    
```

Optimized GCC on the x86

```

movl 20(%ebp), %edx     % Get j
pushl %edx              % Push argument
call foo                % d = foo(j);
movl %eax, %esi         % save d in esi

movl 24(%ebp), %edx     % Get k
pushl %edx              % Push argument
call foo                % e = foo(k);
movl %eax, %ebx         % save e in ebx

movl 28(%ebp), %edx     % Get l
pushl %edx              % Push argument
call foo                % f = foo(l);

addl %ebx, %eax         % e + f
addl %esi, %eax         % d + (e+f)
    
```

Unoptimized vs. Optimized

```

movl 24(%ebp), %eax
pushl %eax
call foo
movl %eax, %ebx

movl 28(%ebp), %eax
pushl %eax
call foo
addl $4, %esp
movl %eax, %eax
movl %eax, -20(%ebp)

movl 28(%ebp), %eax
pushl %eax
call foo
addl $4, %esp
movl %eax, %eax
movl %eax, -24(%ebp)

movl -20(%ebp), %eax
movl -24(%ebp), %edx
addl %edx, %eax
movl %eax, %edx
addl -16(%ebp), %edx
movl %edx, %eax

addl %ebx, %eax
addl %esi, %eax
    
```



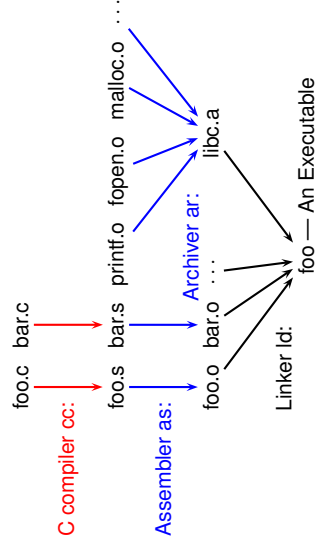
Linking

Goal of the linker is to combine the disparate pieces of the program into a coherent whole.

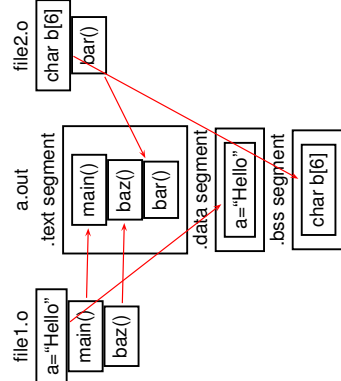
```

file1.c:                               libc.a:
#include <stdio.h> #include <stdio.h> int
char a[] = "Hello"; extern char a[];   printf(char *s, ...)
extern void bar();                      {
int main() {                             static char b[6];
    bar();                                } /* ... */
}                                          void *bar() {
                                          strcpy(b, a);
                                          char *
                                          baz(b);
void baz(char *s) {                       strcpy(char *d, char *s)
    printf("%s", s);                      {
}                                          /* ... */
}                                          }
    
```

Separate Compilation and Linking



Linking



Object Files

Relocatable: Many need to be pasted together. Final in-memory address of code not known when program is compiled

Object files contain

- imported symbols (unresolved "external" symbols)
- relocation information (what needs to change)
- exported symbols (what other files may refer to)

Object Files

```

file1.c:                               libc.a:
#include <stdio.h> #include <stdio.h> int
char a[] = "Hello"; extern char a[];   printf(char *s, ...)
extern void bar();                      {
int main() {                             static char b[6];
    bar();                                } /* ... */
}                                          void *bar() {
                                          strcpy(b, a);
                                          char *
                                          baz(b);
void baz(char *s) {                       strcpy(char *d, char *s)
    printf("%s", s);                      {
}                                          /* ... */
}                                          }
    
```

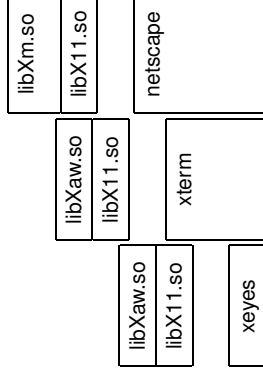

Shared Libraries: First Attempt

Obvious disadvantage: must ensure each new shared library located at a new address.

Works fine if there are only a few libraries; tended to discourage their use.

Shared Libraries

Problem fundamentally is that each program may need to see different libraries **each at a different address**.



Position-Independent Code

Solution: Require the code for libraries to be position-independent. **Make it so they can run anywhere in memory.**

As always, add another level of indirection:

All branching is PC-relative

All data must be addressed relative to a base register.

All branching to and from this code must go through a jump table.

Position-Independent Code for bar()

Normal unlinked code

```

save %sp, -112, %sp
sethi %hi(0), %o0
R_SPARC_HI22 .bss
mov %o0, %o0
R_SPARC_LO10 .bss
sethi %hi(0), %o1
R_SPARC_HI22 a
mov %o0, %o1
call 14
R_SPARC_WDISP30 strcpy
nothi %hi(0), %o0
R_SPARC_HI22 .bss
mov %o0, %o0
R_SPARC_LO10 .bss
call 24
R_SPARC_WDISP30 baz
nop
ret
restore
  
```

gcc -fpic -shared

```

save %sp, -112, %sp
sethi %hi(0x10000), %i7
call %o0 | add PC to %i7
add %i7, 0x198, %i7
ld [ %i7 + 0x20 ], %o0
ld [ %i7 + 0x24 ], %o1
call 10a24 | strcpy
nop
ld [ %i7 + 0x20 ], %o0
call 10a3c | baz
nop
ret
restore
  
```

Annotations:
 - Red arrow points to `call 10a24 | strcpy` with text "Actually just a stub"
 - Red arrow points to `call 10a3c | baz` with text "call is PC-relative"