

COMPILERS AND INTERPRETERS



- Pointers, the addresses we can see
- Programs that write other programs
- Managing the details

A compiler is a program that, when fed itself as input, produces ITSELF!



Then how was the first compiler written?



Problem set #2 is out.
(Warning it is significantly more time consuming than problem set #1)



A WARM UP QUESTION

```
x:      .word    42
```

```
main:   lw       t0, x
```

```
        li       t1, x
```

```
*halt:  j         halt
```

miniRISCv

What is in registers t0 and t1
after these two instructions
are executed?





AN ASIDE: LET'S C

C is the ancestor to most languages commonly used today.

{Algol, Fortran, Pascal} → C → C++ → Java → Rust

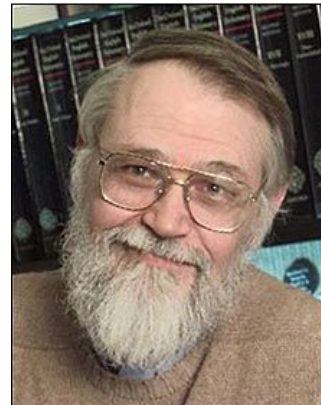
C was developed to write the operating system UNIX.

C is still widely used for "systems" programming

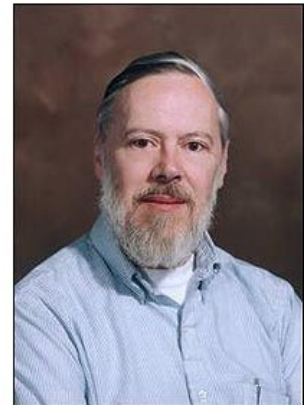
C's developers were frustrated that the high-level languages available at the time, lacked many of the capabilities of assembly.

An advantage of high-level languages is that they are portable (i.e. not ISA specific).

C, thus, was an attempt to create a portable blend of a "high-level language" and "assembler"



Brian Kernighan



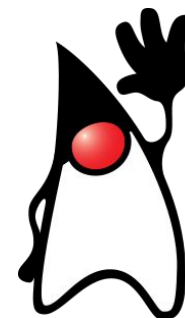
Dennis Ritchie

C BEGAT JAVA



C++ was envisioned to add Object-Oriented (OO) concepts from *Simula* and *CLU* on top of C

Java was envisioned to be more purely OO, and to hide the details of memory management as well as Class/Method/Member implementation



For our purposes C is almost identical to JAVA except:

- C has "functions", whereas JAVA has "methods".
- C has explicit variables that contain the **addresses** of other variables or data structures in memory.
- JAVA hides addresses under the covers.

RECALL YOUR FIRST C POINTER!



Let's start with a feature that Java does not have-- called "pointers"

```
int i = 4;    // simple integer variable
int a[10];   // array of integers (a is a pointer)
int *p;      // pointer to integer (s)
```

`*(expression)` means the "contents of address computed by expression".

`a[i] ≡ *(a+i)`

Array variables are our first hint that "pointers" exist. The name of an array tells us where a collection of indexable variables could be found.

`a` is a constant of type "int *"

We now know that *all* variables are shorthands for addresses in memory.

`a[i] = a[i+1] ≡ *(a+i) = *(a+i+1)`

Normal variables are just the 0th element of a length "1" array.





OTHER POINTER RELATED SYNTAX

```
int i;           // simple integer variable
int a[10];      // array of integers
int *p;         // pointer to integer(s)

p = &i;         // & means address of (not AND)
p = a;         // no need for & on a
p = &a[5];      // address of 6th element of a
*p = 1;        // change value of that location
*(p+1) = 1;    // change value of next location
p[1] = 1;      // exactly the same as above
p++;          // step pointer to the next element
(*p)++;       // increments contents of location
*p++;        // get contents, and then modifies p
```



The ampersand operator, "&", means "give me the address of this variable reference". Whereas the star operator, "*", means "give me the contents of the memory location implied by the expression". These are VERY different things. Not to mention, "&" and "*" can sometimes be confusing because of their other uses as "anding" and "multiplying" operators.



LEGAL USES OF POINTERS

```
int i;           // simple integer variable
int a[10];      // array of integers
int *p;         // pointer to integer(s)
```

So what happens when: `p = &i;`

What is value of `p[0]`?

What is value of `p[1]`?

`p[0]` is **always** an alias for the variable `i` in this context. `p[1]` **could** reference `a[0]`, but don't count on it.





C POINTERS VS. OBJECT SIZE

```
int i;           // simple integer variable
int a[10];      // array of integers
int *p;         // pointer to integer(s)
```

```
i = *p++;
```

Does "p++" really add 1 to the pointer?
NO! It adds 4. Why 4?

```
char *q;
```

```
...
```




The "char" type is slightly different than the type of the same name in Java. C chars are 8-bit signed bytes. Java chars are 16-bits and hold only Unicode variables (they have no sign). Java has a type called "byte" that is most similar to a C "char".

```
q++; // really does add 1
```





CLEAR1,2,3, ALL ARE VALID C!


```
void clear1(int array[], int size) {  
    for (int i = 0; i < size; i++)  
        array[i] = 0;  
}
```

{  Written using "Array" semantics

```
void clear2(int array[], int size) {  
    for (int *p = array; p < array + size; p++)  
        *p = 0;  
}
```

{  Written using C "Pointer" semantics.

```
void clear3(int *array, int size) {  
    int *end = array + size;  
    while (array < end)  
        *array++ = 0;  
}
```

{  Array is just a pointer.



POINTER SUMMARY

- In the "C" world and in the "machine" world:
 - a pointer is just the **address** of an object in memory
 - size of pointer is fixed and architecture dependent, regardless of size of object that it points to
 - to get to objects of the same type, we offset by increments of the object's size in bytes
 - Ex: to get the the i^{th} object add $i * \text{sizeof}(\text{object})$
- More details:
 - $\text{int } R[5] \equiv R$ (i.e. R is an int^* to 20 bytes of storage)
 - $R[i] \equiv *(R+i)$ (array offsets are just pointer arithmetic)
 - $\text{int } *p = \&R[3] \equiv p = (R+3)$ (p points to 3rd element of R)



INDIRECT ADDRESSING

- What we want:
 - The contents of a memory location held in a register
- Examples:

“C”

```
int x = 10;

void main() {
    int *y = &x;
    *y = 2;
}
```

“RISC-V Assembly”

```
x:      .word 10
main:   addi sp,sp,-8
        sw   ra,4(sp)
        li   t0,x          # int *y = &x;
        sw   t0,(sp)
        li   t1,2         # *y = 2;
        lw   t0,(sp)
        sw   t1,(t0)
        lw   ra,4(sp)
        addi sp,sp,8
        jalr x0,ra
```

Loads the
“address”
of x into y,
not its
contents



• Caveats

- You must make sure that the register contains a *valid address* (double, word, or short *aligned* as required)

COMPILERS AS TEMPLATE MATCHERS



The basic task of a compiler is to scan a file looking for particular sequences of operators and keywords called **templates**.

The first major sort of template is an **expression**. We've already played around converting C expressions to assembly language. A compiler does basically the same thing.

```
int x, y;  
y = (x-3)*(y+123456);
```



Here the compiler noticed that the desired constant was too big to fit as an immediate constant, so it creates a new variable, *c*, to keep track of this constant. (Usually, compiler generated constant names are cryptic, so you can't generate them by chance).

```
x: .word 0  
y: .word 0  
c: .word 123456  
  
...  
lw      t0, x  
addi    t0, t0, 3  
lw      t1, y  
lw      t2, c  
add     t1, t1, t2  
mul     t0, t0, t1  
sw      t0, y
```

Once a template is matched, a compiler emits a specific code sequence.



C ARRAYS

The C source code

```
int hist[100];  
int score = 92;  
...  
hist[score] += 1;
```

might translate to:

```
hist:    .space 100  
score:  .word 92  
  
li      t3, hist  
lw      t2, score  
slli    t1, t2, 2  
add     t1, t3, t1  
lw      t0, (t1)  
addi    t0, t0, 1  
sw      t0, (t1)
```

score:

| | | | |
|--|--|--|----|
| | | | 92 |
| | | | |
| | | | |
| | | | |

⋮

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

hist:

Address:

CONSTANT base address + scaled VARIABLE offset



C "STRUCTS"

- C "structs" are lightweight "container objects" - objects with members, but no methods.
- There is special "Java-like" syntax for accessing particular members: *variable.member* (actually, Java's dot operator "." is borrowed from C)
- You can also have pointers to structs.

C provides an new operator to access them:

pointerVariable->member

In place of the alternative syntax:

*(*pointerVariable).member*

An implied dereference with an implied offset. Similar to $*(p+i)$



Here the dereference is more explicit by requires more typing.



```
struct Point {  
    int x, y;  
} P1, P2, *p;  
...  
P1.x = 157;  
...  
p = &P1;  
p->y = 123;
```



STRUCTS IN ACTION

```

struct Point {
    int x, y;
} P1, P2, *p;
...
P1.x = 157;
...
p = &P1;
p->y = 123;

```

Address:
VARIABLE base address + **CONSTANT** offset

might translate to:

```

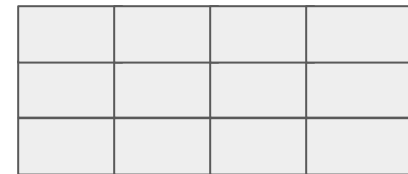
P1: .space 8
P2: .space 8
p: .space 4
...

```

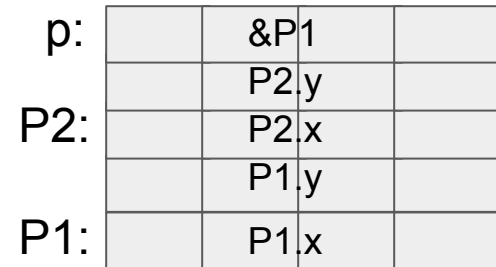
```

addi    t1,x0,P1      # P1.x = 157
addi    t0,x0,157
sw      t0,(t1)
sw      t1,p          # p = &P1
lw      t2,p          # p->y = 123
addi    t0,x0,123
sw      t0,4(t2)

```



⋮





C "IF" TO ASSEMBLY TRANSLATION

C code:

```
if (expr) {  
    STUFF  
}
```

C code:

```
if (expr) {  
    STUFF1  
} else {  
    STUFF2  
}
```

RISC-V assembly:

```
(compute expr in t0)  
beq t0,x0,Lendif  
(compile STUFF)
```

Lendif:

RISC-V assembly:

```
(compute expr in t0)  
beq t0,x0,Lelse  
(compile STUFF1)
```

j Lendif

Lelse:

```
(compile STUFF2)
```

Lendif:

Note: the branches used in assembly "SKIP" code blocks rather than cause them to be executed. This often results in a complement test!





C "WHILE" LOOPS

C code:

```
while (expr) {  
    STUFF  
}
```

Assembly:

```
Lwhile:  
    (compute expr in t0)  
    beq    t0,x0,Lendw  
    (compile STUFF)  
    j      Lwhile  
Lendw:
```

Alternate Assembly:

```
    j      Ltest  
Ltest:  
    (compile STUFF)  
Ltest:  
    (compute expr in t0)  
    bne    t0,x0,Lwhile  
Lendw:
```

Compilers spend a lot of time optimizing in and around loops.

- moving all possible computations outside of loops
- unrolling loops to reduce branching overhead
- simplifying expressions that depend on "loop variables"



C "FOR" LOOPS

- Most high-level languages provide loop constructs that establish and update an iterator, which controls the loop's behavior

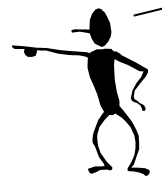
```
for (initialization; conditional; afterthought) {  
    STUFF;  
}
```

Assembly:
(compile **initialization**)

```
Lfor:  
    (compute conditional in t0)  
    beq t0, z0, Lendfor  
    (compile STUFF)  
    (compile afterthought)  
    j    Lfor  
Lendfor:
```

For loops are the most commonly used form of iteration found programming languages.

Their advantage is readability. They bring together the three essential components of iteration, setting an initial value, establishing a termination condition, and giving an update rule.



Ahhh, but one other iteration forms there are!

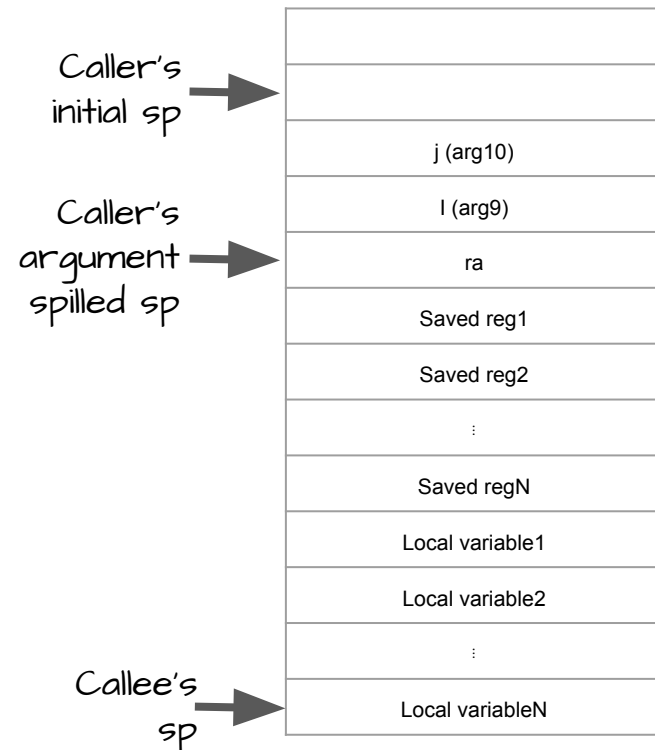




SOME AFTERTHOUGHTS

What if you have more than 8 arguments? The caller allocates space on the stack before calling the function

```
# x = func(a,b,c,d,e,f,g,h,i,j)
lw    a0,a(gp)
lw    a1,b(gp)
lw    a2,c(gp)
lw    a3,d(gp)
lw    a4,e(gp)
lw    a5,f(gp)
lw    a6,g(gp)
lw    a7,h(gp)
addi  sp,sp,-8 # spill onto stack
lw    t0,i(gp)
sw    t0,(sp)
lw    t0,j(gp)
sw    t0,4(sp)
jal   ra,func
addi  sp,sp,8
```





NEXT TIME

- The details behind assemblers
- 2-pass and 1-pass assembly
- Linkers and dynamic libraries

