

CS:APP2e Web Aside ASM:EASM: Combining Assembly Code with C Programs*

Randal E. Bryant
David R. O'Hallaron

June 5, 2012

Notice

The material in this document is supplementary material to the book Computer Systems, A Programmer's Perspective, Second Edition, by Randal E. Bryant and David R. O'Hallaron, published by Prentice-Hall and copyrighted 2011. In this document, all references beginning with "CS:APP2e " are to this book. More information about the book is available at csapp.cs.cmu.edu.

This document is being made available to the public, subject to copyright provisions. You are free to copy and distribute it, but you should not use any of this material without attribution.

1 Combining Assembly Code with C Programs

In the early days of computing, most programs were written in assembly code. Even large-scale operating systems were written without the help of high-level languages. This becomes unmanageable for programs of significant complexity. Since assembly code does not provide any form of type checking, it is very easy to make basic mistakes, such as using a pointer as an integer rather than dereferencing the pointer. Even worse, writing in assembly code locks the entire program into a particular class of machine. Rewriting an assembly language program to run on a different machine can be as difficult as writing the entire program from scratch.

Aside: Writing large programs in assembly code.

Frederick Brooks, Jr., a pioneer in computer systems wrote a fascinating account of the development of OS/360, an early operating system for IBM machines [2] that still provides important object lessons today. He became a devoted believer in high-level languages for systems programming as a result of this effort. **End Aside.**

Early compilers for higher-level programming languages did not generate very efficient code and did not provide access to the low-level data representations, as is often required by systems programmers. Programs

*Copyright © 2010, R. E. Bryant, D. R. O'Hallaron. All rights reserved.

requiring maximum performance or requiring low-level access to data structures were still often written in assembly code. Nowadays, however, optimizing compilers have largely removed performance optimization as a reason for writing in assembly code. Code generated by a high quality compiler is generally as good or even better than what can be achieved manually. The C language has largely eliminated low-level data structure access as a reason for writing in assembly code. The ability to access low-level data representations through unions and pointer arithmetic, along with the ability to operate on bit-level data representations, provide sufficient access to the machine for most programmers. For example, almost all of the code for modern operating systems, including Linux, Windows, and MacOS, are written in C.

Nonetheless, there are times when writing in assembly code is the only option. This is especially true when implementing an operating system. For example, there are a number of special registers storing process state information that the operating system must access. There are either special instructions or special memory locations for performing input and output operations. Even for application programmers, there are some machine features, such as the values of the condition codes, that cannot be accessed directly in C.

The challenge then is to integrate code consisting mainly of C with a small amount written in assembly language. In this document, we will describe two such mechanisms. The first is to write a few key functions in assembly code, using the same conventions for argument passing and register usage as are followed by the C compiler. The linker then serves to combine the two forms of code into a single program. This approach is often feasible for simple functions, and it does not require any GCC-specific constructs. An alternative to writing an entire function in C is to embed assembly code within a C program. GCC supports *inline assembly* via the `asm` directive. Inline assembly allows the user to insert assembly code directly into the code sequence generated by the compiler. Features are provided to indicate to the compiler how to interface to the inserted code. The resulting code, of course, only runs on a specific class of machines, but we will see, for example, that it is often possible to have inline assembly that compiles properly on both IA32 and x86-64 machines. The `asm` directive is also specific to GCC, creating an incompatibility with many other compilers. Nonetheless, this can be a useful way to keep the amount of machine-dependent code to an absolute minimum.

Our presentation is drawn both from the GCC documentation [3], as well as the book by Blum [1]. The former is, of course, the authoritative reference, but it does not provide any examples. Blum's book, on the other hand, provides many practical tips and examples.

2 Program Example

For our presentation, we will develop several implementations of functions with the following prototypes. These examples provide real-life cases where gaining access to the condition codes will enable us to monitor the status of a computation.

```
/* Multiply x and y. Store result at dest.
   Return 1 if multiplication did not overflow
*/
int tmult_ok(int x, int y, int *dest);

/* Multiply x and y. Store result at dest.
```

```

    Return 1 if multiplication did not overflow
*/
int umult_ok(unsigned x, unsigned y, unsigned *dest);

```

Each function is to compute the product of arguments x and y and store the result in the memory location specified by argument $dest$. As return values, they should return 0 when the multiplication overflows, requiring more than 32 bits to represent the true product, and 1 when it does not. We have separate functions for signed and unsigned multiplication, since their overflow conditions differ. We examined ways to determine whether a multiplication has overflowed using C (see CS:APP2e Problem 2.35 and 2.36), but all of these methods require performing additional operations to check the result of a multiplication.

Examining the documentation for the x86 multiply instructions `mull` and `imull`, we see that both set the carry flag `CF` when they overflow. Thus, by inserting code that checks this flag after performing a multiplication of x and y , we should be able to easily test for multiplicative overflow.

3 Handwritten Assembly-Code Functions

Although writing complex programs entirely in assembly code is a daunting task, we can often narrow the amount of functionality that needs to be expressed in assembly code to a small amount and then write this code as functions in a separate file. The compiled C code is combined with the assembled assembly code by the linker. For example, if file `p1.c` contains C code and file `p2.s` contains assembly code, then the compilation command

```
unix> gcc -o p p1.c p2.s
```

will cause file `p1.c` to be compiled, file `p2.s` to be assembled, and the resulting object code to be linked to form an executable program `p`.

In order for the assembler to generate information required by the linker about a function, we must declare the function to be *global*. Whereas in C, any function is global unless it is declared to be static, the assembler assumes that a file is only locally available to functions within the same file, unless it is declared global. If we have assembly code for a function `fun` in a file, then we should precede it with the declaration

```
.globl fun
```

We have found that even when writing functions in assembly code, it is best to let GCC do as much of the work as possible. Toward that end, it often helps to write a function in C similar to the desired functionality and then generate an initial version of the assembly code by running GCC with the command-line option `-S`. This code provides a good starting point for fetching arguments, allocating and deallocating the stack frame, and so forth. It is much easier to edit existing assembly code than to start from scratch.

As an example, consider the following approximation to our function `tmult_ok`:

```

/* Starter function for tmult_ok */
int tmult_ok_asm(int x, int y, int *dest) {
    int p = x*y;
    *dest = p;
}

```

```

    return p > 0;
}

```

This function does much of what we want for `tmult_ok`—it multiplies arguments `x` and `y`, stores the product at `dest`, and returns a 0 or 1 based on the result. Its only shortcoming is that it checks the wrong property, but this is just one part of the overall computation.

GCC generates the following assembly code for the initial function:

```

1 .globl tmult_ok_asm
2 tmult_ok_asm:
3   pushl   %ebp
4   movl    %esp, %ebp
5   movl    12(%ebp), %eax
6   imull   8(%ebp), %eax
7   movl    16(%ebp), %edx
8   movl    %eax, (%edx)
9   testl   %eax, %eax
10  setg    %al
11  movzbl  %al, %eax
12  popl    %ebp
13  ret

```

Note the presence of the `.globl` declaration in this code. Only two lines of this code need to be changed. Line 9 sets condition codes based on the 32-bit product of `x` and `y`. We want to eliminate this instruction and instead rely on the condition-code values set by the `imull` instruction. Line 10 sets the low-order byte of register `%eax` based on the zero and sign flags. We want to set the byte based on the carry flag.

Examining CS:APP2e Figure 3.11, we see that the instruction `setae` can be used to set the low-order byte of a register to 0 when the carry flag is set and to 1 otherwise. We can therefore make small edits to the assembly code to get the desired function:

```

1 # Hand-generated code for tmult_ok
2 .globl tmult_ok_asm
3 tmult_ok_asm:
4     pushl   %ebp
5     movl    %esp, %ebp
6     movl    12(%ebp), %eax # Get y
7     imull   8(%ebp), %eax # Multiply by y
8     movl    16(%ebp), %edx # Get dest
9     movl    %eax, (%edx)  # Store product at dest
10 # Deleted code
11 #     testl   %eax, %eax
12 #     setg    %al
13 # Inserted code
14     setae   %al           # Set low-order byte
15 # End of inserted code
16     movzbl  %al, %eax     # Zero remaining bytes
17     popl    %ebp
18     ret

```

We show this code in the exact form it appears in the file, rather than the more stylized way we have presented assembly code. As this example shows, it helps to add annotations to the assembly code as documentation. Anything to the right of the symbol '#' is treated as a comment by the assembler.

Practice Problem 1:

Create an x86-64 implementation of `tmult_ok` suitable for assembling and linking with C code. You might find it useful to start with assembly code generated for a similar function, as we did for the IA32 code.

4 Basic Inline Assembly

The basic form of inline assembly is to write code that looks like a procedure call:

```
asm( code-strings );
```

The term *code-strings* denotes an assembly code sequence given as one or more quoted strings (with no delimiters between them.) The compiler will insert these strings verbatim into the assembly code being generated, and hence the compiler-supplied and the user-supplied assembly will be combined. The compiler does not check the string for errors, and so the first indication of a problem might be an error report from the assembler.

In an attempt to use the least amount of both assembly code and detailed analysis, we attempt to implement `tmult_ok` with the following code:

```
/* First attempt. Does not work */
int tmult_ok1(int x, int y, int *dest)
{
    int result = 0;
    *dest = x*y;
    asm("setae %al");
    return result;
}
```

The strategy here is to exploit the fact that register `%eax` is used to store the return value. Assuming the compiler uses this register for variable `result`, our intention is that the first line of the C code will set the register to 0. The inline assembly will insert code that sets the low-order byte of this register appropriately, and the register will be used as the return value.

Unfortunately, the generated code does not work as desired. In running tests, it returns 0 every time it is called. On examining the generated assembly code for this function, we find the following:

```
IA32 code for tmult_ok1 (Does not work)
x at %ebp+8, y at %ebp+12, dest at %ebp+16
1 tmult_ok1:
2   pushl   %ebp
```

```

3  movl    %esp, %ebp
4  movl    12(%ebp), %eax    Get y
5  imull   8(%ebp), %eax    Multiply by x
6  movl    16(%ebp), %edx   get dest
7  movl    %eax, (%edx)    store product at dest
   Code generated by asm
8  setae   %al              Set low-order byte
   End of asm-generated code
9  movl    $0, %eax        Set result to 0
10 popl    %ebp
11 ret

```

GCC has its own ideas of code generation. Instead of setting register `%eax` to 0 at the beginning of the function, the generated code does so at the very end (line 9), and so the function always returns 0. The fundamental problem is that the compiler has no way to know what the programmer’s intentions are, and how the assembly code should interface with the rest of the generated code. Clearly, more sophisticated mechanisms are required to embed assembly code within C code.

5 Extended Form of `asm`

GCC provides an extended version of `asm` that allows the programmer to specify which program values are to be used as operands to an assembly code sequence, and which registers are overwritten by the assembly code. With this information the compiler can generate code that will correctly set up the required source values, execute the assembly instructions, and make use of the computed results. It will also have information it requires about register usage so that important program values are not overwritten by the assembly code instructions.

The general syntax of an extended `asm` directive is

```
asm( code-strings [ : output-list [ : input-list [ : overwrite-list ] ] ] ) ;
```

where the square brackets denote optional arguments. The directive contains one or more strings giving stylized versions of the lines of assembly code. These are followed by optional lists of *outputs* (i.e., results generated by the assembly code), *inputs* (i.e., source values for the assembly code), and registers that are overwritten by the assembly code. These lists are separated by the colon (‘:’) character. As the square brackets show, we only include lists up to the last nonempty list.

The syntax for a code string is reminiscent of that for the format string in a `printf` statement. It consists an assembly code instruction, but the operands are written in a symbolic form, with references to the operand expressions in the output and input lists. Within the assembly code, we give names to the operands. Earlier versions of GCC required these names to be of the form `%0`, `%1`, and so on, up to `%9`, based on the order of the operands in the two lists. Since version 3.1, a more descriptive naming convention is supported, where names are written with the notation `[%name]`. Register names such as “`%eax`” must be written with an extra ‘%’ symbol, such as “`%%eax`.”

If multiple instructions are given, it is critical that return characters be inserted between them. The conventional method of doing this is to finish all but the final string with the sequence “`\n\t`,” so that the

generated assembly code follows the normal formatting conventions for assembly code [1].

The output input list is a comma-separated list, with each list element of the form `[name] "=r" (expr)`, giving the name of the operand, the fact that it is an output, and the C expression indicating the destination for the instruction result. It can be any assignable value (known in C as an *lvalue*). The compiler will generate the necessary code sequence to perform the assignment. The tag `"=r"` indicates that the instruction will write its result in an integer register. The input list entries have the same format, except that each entry has tag `"r"`, indicating that it will be read from an integer register. Each input operand can be any C expression. The compiler will generate the necessary code to evaluate the expression. The overwrite list consists of a comma-separated list of register names, with each in double quotes.

As an illustration, the following is a better implementation of `tmult_ok` using the extended assembly directive to indicate to the compiler that the assembly code generates the value for the variable `result`:

```
int tmult_ok2(int x, int y, int *dest)
{
    int result;

    *dest = x*y;
    asm("setae %%b1          # Set low-order byte\n\t"
        "movzbl %%b1,%[val] # Zero extend to be result"
        : [val] "=r" (result) /* Output */
        : /* No inputs */
        : "%b1" /* Overwrites */
        );
    return result;
}
```

We see that the `asm` directive has two code strings: one for the `setae` instruction, and one to zero-extend the low-order byte to form the result. We can see that we have chosen register `%b1` as the destination of the `setae` instruction and as the source of the `movzbl` instruction. The register operand must be written as `"%b1."` Observe also that we can add comments to our code, and that all but the last lines are terminated with `"\n\t."` We have given the name `val` to indicate the final value generated by the code. This is shown in the output list as corresponding to program variable `result`. We also indicate that our code will overwrite register `%b1`.

When we compile this code for IA32, GCC generates the following assembly code:

```
IA32 code for tmult_ok2
x at %ebp+8, y at %ebp+12, dest at %ebp+16
1 tmult_ok2:
2  pushl   %ebp
3  movl    %esp, %ebp
4  pushl   %ebx
5  movl    12(%ebp), %eax          Get y
6  imull   8(%ebp), %eax          Multiply by x
7  movl    16(%ebp), %edx          Get dest
8  movl    %eax, (%edx)           Store product at dest
Code generated by asm
```

```

 9   setae %bl           Set low-order byte
10   movzbl %bl,%eax    Zero extend to be result
    End of asm-generated code
11   popl   %ebx
12   popl   %ebp
13   ret

```

We see here that GCC has designated register `%eax` to hold program value `result`. We also see that the program saves the value of `%ebx` on the stack (line 4) and restores it at the end (line 11.) Since this register is a callee-saved register, and we have indicated that our code will overwrite its low-order byte (register `%bl`), GCC takes the necessary steps to preserve its value. The code compiles correctly for x86-64, as well, something that would not be possible if we wrote the entire function in assembly code.

As a further refinement, we can simplify the code even more and make use of GCC's ability to work with different data types. GCC uses the type information for an operand in determining which register to substitute for an operand name in the code string. In the version given as function `tmult_ok2`, it used a 32-bit register `%eax`, based on the fact that variable `result` has data type `int`. Instead, we can use a variable `bresult` of type `unsigned char` in the output list, and have this operand be the destination of the `setae` instruction:

```

/* Uses extended asm to get reliable code */
int tmult_ok3(int x, int y, int *dest)
{
    unsigned char bresult;
    *dest = x*y;

    asm("setae %[b]          # Set result"
        : [b] "=r" (bresult) /* Output */
        );

    return (int) bresult;
}

```

The compiler will use a single-byte register identifier as the destination for the `setae` instruction, and then use this register as the source operand of a `movzbl` instruction to implement the casting of `bresult` to data type `int`. This simplified form also avoids the need for us to make use of a specific register, and hence we need not specify any overwritten registers.

One would expect the same code sequence could be used for `umult_ok`, but GCC uses the `imull` (signed multiply) instruction for both signed and unsigned multiplication. This generates the correct value for either product, since we are only concerned with the low-order 32 bits, but it sets the carry flag according to the rules for signed multiplication. We therefore need to include an assembly-code sequence that explicitly performs unsigned multiplication using the `mull` instruction as documented in CS:APP2e Figure 3.9. This instruction is only available in single-operand form, implicitly using register `%eax` as a source, and registers `%eax` and `%edx` as destinations. All of this requires a more elaborate `asm` directive:

```

int umult_ok(unsigned x, unsigned y, unsigned *dest)

```

```

{
    unsigned char bresult;

    asm("movl %[x],%%eax      # Get x\n\t"
        "mull %[y]           # Unsigned multiply by y\n\t"
        "movl %%eax,%[p]     # Store low-order 4 bytes at dest\n\t"
        "setae %[b]         # Set result"
        : [p] "=r" (*dest), [b] "=r" (bresult) /* Outputs */
        : [x] "r" (x), [y] "r" (y)           /* Inputs */
        : "%eax", "%edx"                  /* Overwrites */
        );

    return (int) bresult;
}

```

This code makes use of many of the features of the extended `asm` directive. The two output operands are given symbolic names `p` (the product) and `b` (the status byte), while the two input operands have symbolic names `x` and `y`. We can see that output operand `p` is associated with the expression `*dest`, while `b` is associated with local variable `bresult`. We need to list both registers `%eax` and `%edx` on the overwrite list.

To see how the compiler generates code in connection with an `asm` directive, here is the code generated for `umult_ok`:

```

    IA32 code for umult_ok
    x at %ebp+8, y at %ebp+12, dest at %ebp+16
1 umult_ok:
2  pushl   %ebp
3  movl   %esp, %ebp
4  pushl   %ebx
5  movl   12(%ebp), %ebx      Get y
6  movl   8(%ebp), %ecx      Get x
    Code generated by asm
7  movl   %ecx,%eax          Copy y
8  mull   %ebx               Unsigned multiply by y
9  movl   %eax,%ecx          Copy low-order 4 bytes
10 setae  %bl                Set result
    End of code generated by asm
11 movl   16(%ebp), %eax     Get dest
12 movl   %ecx, (%eax)       Store product at dest
13 movzbl %bl, %eax          Zero-extend result
14 popl   %ebx
15 popl   %ebp
16 ret

```

We can see that GCC has chosen the following register allocations for the operands: `%ecx` for `x`, `%ebx` for `y`, `%ecx` for `p`, and `%bl` for `b`.

6 Concluding Remarks

We have explored two ways to combine assembly code with C code to generate a program. Writing a complete function in assembly code as a separate file has the advantage that it uses existing and familiar technology: the assembler and the linker. Using the facility for GCC to insert assembly code directly in a C function has the advantage that we can greatly limit the amount of machine-specific code.

Although the syntax of the `asm` directive is somewhat arcane, and its use makes the code less portable, it can be very useful for writing programs that access machine-level features using a minimal amount of assembly code. We have found that a certain amount of trial and error is required to get code that works. The best strategy is to compile the code with the `-S` command-line option and then examine the generated assembly code to see if it will have the desired effect. It is important to note that in processing an `asm` directive, GCC has no real understanding of the assembly code it is generating. It merely follows a set of syntactic rules for replacing the symbolic names of operands with different register identifiers. The code should be tested with different settings of switches such as with different levels of optimization.

Practice Problem 2:

Write the complete function for `umult_ok` in IA32 assembly code. How does the effort required for this compare to using the `asm` directive?

Practice Problem 3:

Use the `asm` directive to implement an IA32 function with the prototype

```
/* Multiply two n-bit numbers to get 2n-bit result,
   where n = 8*sizeof(unsigned long)
*/
typedef unsigned long ulong_t;
void ulmult_full(ulong_t x, ulong_t y, ulong_t *dest);
```

This function should compute the full 64-bit product of its arguments and store the result in the destination array, with `dest[0]` having the low-order four bytes and `dest[1]` having the high-order four bytes.

Practice Problem 4:

Implement the function `ulmult_full` in x86-64 code, computing the 128-bit product of 64-bit values `x` and `y`.

Practice Problem 5:

X86 machines have a parity flag PF as one of the condition codes. Every arithmetic or logical operation sets this flag when the low-order eight bits of a result have even parity, meaning that they contain an even number of ones. Whether the operation computes an 8-bit result or a larger one, the parity flag depends only on the low-order 8 bits.

Consider the following function prototype:

```
int odd_parity(unsigned long x);
```

This function should determine whether its argument has odd parity, meaning that it contains an odd number of ones. We want this function to operate correctly when compiled for either IA32, in which case the argument is 32 bits, or for x86-64, in which case the argument is 64 bits.

Write a C function including an `asm` directive to implement this function, using the parity flag to compute the parity 8 bits at a time.

Solutions to Practice Problems

Problem 1 Solution: [Pg. 5]

The x86-64 code is much simpler than the IA32 code, due to the passing of arguments through registers:

```
1 .globl tmult_ok_asm
2 tmult_ok_asm:
3     imull    %edi, %esi
4     movl    %esi, (%rdx)
5     setae   %al    # Set low-order byte
6     movzbl  %al, %eax
7     ret
```

Problem 2 Solution: [Pg. 10]

We can once again start with code generated for a similar function by GCC, but more extensive editing is required:

```
1 # Hand-generated code for umult_ok
2 .globl umult_ok_asm
3 umult_ok_asm:
4     pushl   %ebp
5     movl   %esp, %ebp
6     movl   12(%ebp), %eax # Get y
7     movl   8(%ebp), %ecx # Get x
8     mull   %ecx          # Unsigned multiply
9     movl   16(%ebp), %edx # Get dest
10    movl   %eax, (%edx)  # Store product at dest
11    setae  %al          # Set low-order byte
12    movzbl %al, %eax    # Zero remaining bytes
13    popl   %ebp
14    ret
```

Problem 3 Solution: [Pg. 10]

This function bears many similarities to the function `umult_ok`, except that we want to store both the high and the low-order words of the product:

```

/* 32-bit version */
void ulmult_full(ulong_t x, ulong_t y, ulong_t *dest)
{
    asm("movl %[x],%%eax      # Get x\n\t"
        "mull %[y]           # Unsigned multiply by y\n\t"
        "movl %%eax,%[lo]    # Store low-order 4 bytes\n\t"
        "movl %%edx,%[hi]    # Store high-order 4 bytes"
        : [lo] "=r" (dest[0]), [hi] "=r" (dest[1]) /* Outputs */
        : [x] "r" (x), [y] "r" (y) /* Inputs */
        : "%eax", "%edx" /* Overwrites */
        );
}

```

Problem 4 Solution: [Pg. 10]

The code for this function is very similar to the 32-bit version, except that we want to use the “q” form of the instructions, and the “r” form of the registers:

```

/* 64-bit version */
void ulmult_full(ulong_t x, ulong_t y, ulong_t *dest)
{
    asm("movq %[x],%%rax      # Get x\n\t"
        "mulq %[y]           # Unsigned multiply by y\n\t"
        "movq %%rax,%[lo]    # Store low-order 8 bytes\n\t"
        "movq %%rdx,%[hi]    # Store high-order 8 bytes"
        : [lo] "=r" (dest[0]), [hi] "=r" (dest[1]) /* Outputs */
        : [x] "r" (x), [y] "r" (y) /* Inputs */
        : "%rax", "%rdx" /* Overwrites */
        );
}

```

Problem 5 Solution: [Pg. 10]

There are many solutions to this problem. This one is perhaps the simplest. It simply shifts the successive bytes of argument *x* into the low-order byte and thereby tests the parity of each non-zero byte. We use the `testb` instruction to test each byte, since the parity flag depends only on the low-order byte, and this makes the code portable between IA32 and x86-64.

```

/* Using ASM to access parity flag */
int odd_parity(unsigned long x) {
    int result = 0;
    while (x != 0) {
        char bresult;
        unsigned char bx = x & 0xff;
        asm("testb %[bx],%[bx] # Test value of low-order byte\n\t"
            "setnp %[v]      # Set if odd parity"
            : [v] "=r" (bresult) /* Output */
            : [bx] "r" (bx) /* Input */
            );
        result ^= (int) bresult;
    }
}

```

```
        x = x >> 8;  
    }  
    return result;  
}
```

References

- [1] R. Blum. *Professional Assembly Language*. Wiley, 2005.
- [2] F. P. Brooks, Jr. *The Mythical Man-Month, Second Edition*. Addison-Wesley, 1995.
- [3] *GCC Online Documentation*. Available at <http://gcc.gnu.org/>.