

Chapter 1

x86 Assembly

...

This chapter was derived from a document written by Adam Ferrari and later updated by Alan Batson, Mike Lack and Anita Jones

1.1 Introduction

This small guide, in combination with the material covered in the class lectures on assembly language programming, should provide enough information to do the assembly language labs for this class. In this guide, we describe the basics of 32-bit x86 assembly language programming, covering a small but useful subset of the available instructions and assembler directives. However, real x86 programming is a large and extremely complex universe, much of which is beyond the useful scope of this class. For example, the vast majority of real (albeit older) x86 code running in the world was written using the 16-bit subset of the x86 instruction set. Using the 16-bit programming model can be quite complex – it has a segmented memory model, more restrictions on register usage, and so on. In this guide we restrict our attention to the more modern aspects of x86 programming, and delve into the instruction set only in enough detail to get a basic feel for programming x86 compatible chips at the hardware level.

1.2 Registers

Modern (i.e 386 and beyond) x86 processors have 8 32-bit general purpose registers, as depicted in Figure 1.1. The register names are mostly historical in nature. For example, EAX used to be called the “accumulator” since it was used by a number of arithmetic operations, and ECX was known as the “counter” since it was used to hold a loop index. Whereas most of the registers have lost their special purposes in the modern instruction set, by convention, two are reserved for special purposes – the stack pointer (ESP) and the base pointer (EBP).

In some cases, namely EAX, EBX, ECX, and EDX, subsections of the registers may be used. For example, the least significant 2 bytes of EAX can be treated as a 16-bit register called AX. The least significant byte of AX can be used as a single 8-bit register called AL, while the most significant byte of AX can be used as a single 8-bit register called AH. It is important to realize that these names refer to the same

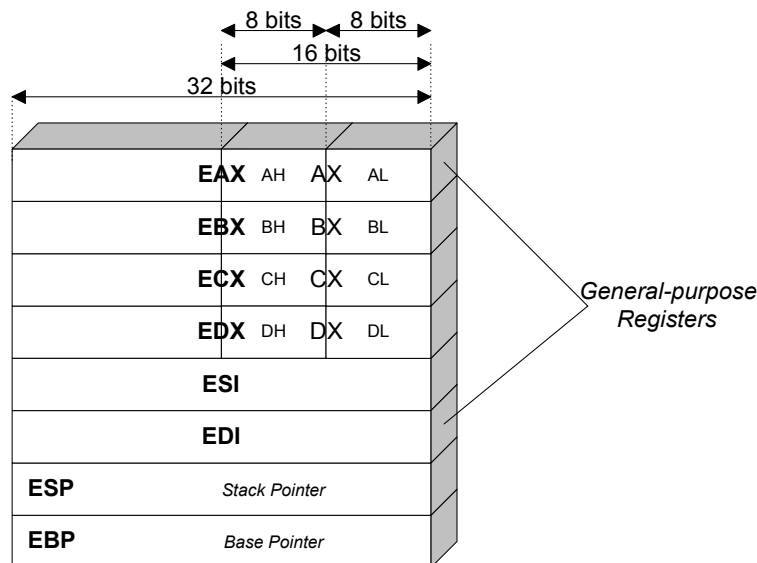


Figure 1.1: The x86 register set

physical register. When a two-byte quantity is placed into DX, the update affects the value of EDX (in particular, the least significant 16 bits of EDX). These “sub-registers” are mainly hold-overs from older, 16-bit versions of the instruction set. However, they are sometimes convenient when dealing with data that are smaller than 32-bits (e.g. 1-byte ASCII characters).

When referring to registers in assembly language, the names are not case-sensitive. For example, the names EAX and eax refer to the same register.

1.3 Memory and Addressing Modes

1.3.1 Declaring Static Data Regions

You can declare static data regions (analogous to global variables) in x86 assembly using special assembler directives for this purpose. Data declarations should be preceded by the `.DATA` directive. Following this directive, the directives `DB`, `DW`, and `DD` can be used to declare one, two, and four byte data locations, respectively. Declared locations can be labeled with names for later reference - this is similar to declaring variables by name, but abides by some lower level rules. For example, locations declared in sequence will be located in memory next to one another. Some example declarations are depicted in Listing 1.1.

The last example in Listing 1.1 illustrates the declaration of an array. Unlike in high level languages where arrays can have many dimensions and are accessed by indices, arrays in assembly language are simply a number of cells located contiguously in memory. Two other common methods used for declaring arrays of data are the `DUP` directive and the use of string literals. The `DUP` directive tells the assembler to duplicate an expression a given number of times. For example, the statement “4 DUP(2)” is equivalent to “2, 2, 2, 2”. Some examples of declaring arrays are depicted in Listing 1.2.

Listing 1.1: Declaring x86 memory regions

```

section .data
var      DB      64      ; Declare a byte containing the value 64. Label the
                        ; Memory location "var".
var2     DB      ?       ; Declare an uninitialized byte labeled "var2".
        DB      10      ; Declare an unlabeled byte initialized to 10. This
                        ; byte will reside at the memory address var2+1.
X        DW      ?       ; Declare an uninitialized two-byte word labeled "X".
Y        DD      3000    ; Declare 32 bits of memory starting at address "Y"
                        ; initialized to contain 3000.
Z        DD      1,2,3   ; Declare three 4-byte words of memory starting at
                        ; address "Z", and initialized to 1, 2, and 3,
                        ; respectively. E.g. 3 will be stored at address Z+8.

```

Listing 1.2: Declaring x86 arrays in memory

```

section .data
bytes    DB      10      DUP(?) ; Declare 10 uninitialized bytes starting at
                        ; the address "bytes".
arr      DD      100     DUP(0) ; Declare 100 4 bytes words, all initialized
                        ; to 0, starting at memory location "arr".
str      DB      'hello', 0    ; Declare 5 bytes starting at the address
                        ; "str" initialized to the ASCII character
                        ; values for the characters 'h', 'e', 'l',
                        ; 'l', 'o', and '\0' (NULL), respectively.

```

1.3.2 Addressing Memory

Modern x86-compatible processors are capable of addressing up to 2^{32} bytes of memory; that is, memory addresses are 32-bits wide. For example, in Listings 1.1 and x86-declaring-arrays.lst, where we used labels to refer to memory regions, these labels are actually replaced by the assembler with 32-bit quantities that specify addresses in memory. In addition to supporting referring to memory regions by labels (i.e. constant values), the x86 provides a flexible scheme for computing and referring to memory addresses:

X86 Addressing Mode Rule – Up to two of the 32-bit registers and a 32-bit signed constant can be added together to compute a memory address. One of the registers can be optionally pre-multiplied by 2, 4, or 8.

To see this memory addressing rule in action, we'll look at some example `mov` instructions. As we'll see later in Section 1.4.1, the `mov` instruction moves data between registers and memory. This instruction has two operands – the first is the destination (where we're moving data *to*) and the second specifies the source (where we're getting the data *from*). Some examples of `mov` instructions using address computations that obey the above rule are shown in Listing 1.3.

Some examples of incorrect address calculations are shown in Listing 1.4.

Listing 1.3: Valid x86 addressing modes

mov eax, [ebx]	<i>; Move the 4 bytes in memory at the address contained ; in EBX into EAX</i>
mov [var], ebx	<i>; Move the contents of EBX into the 4 bytes at memory ; address "var" (Note, "var" is a 32-bit constant).</i>
mov eax, [esi-4]	<i>; Move 4 bytes at memory address ESI+(-4) into EAX</i>
mov [esi+eax], cl	<i>; Move the contents of CL into the byte at address ; ESI+EAX</i>
mov edx, [esi+4*ebx]	<i>; Move the 4 bytes of data at address ESI+4*EBX into ; EDX</i>

Listing 1.4: Invalid x86 addressing modes

mov eax, [ebx-ecx]	<i>; Can only add register values</i>
mov [eax+esi+edi], ebx	<i>; At most 2 registers in address computation</i>

1.3.3 Size Directives

In general, the intended size of the of the data item at a given memory address can be inferred from the assembly code instruction in which it is referenced. For example, in all of the above instructions, the size of the memory regions could be inferred from the size of the register operand – when we were loading a 32-bit register, the assembler could infer that the region of memory we were referring to was 4 bytes wide. When we were storing the value of a one byte register to memory, the assembler could infer that we wanted the address to refer to a single byte in memory. However, in some cases the size of a referred-to memory region is ambiguous. Consider the instruction `mov [ebx], 2`.

Should this instruction move the value 2 into the single byte at address EBX? Perhaps it should move the 32-bit integer representation of 2 into the 4-bytes starting at address EBX. Since either is a valid possible interpretation, the assembler must be explicitly directed as to which is correct. The size directives `BYTE PTR`, `WORD PTR`, and `DWORD PTR` serve this purpose. For examples, see Listing 1.5.

Listing 1.5: x86 size directive usage

mov BYTE PTR [ebx], 2	<i>; Move 2 into the single byte at memory ; location EBX</i>
mov WORD PTR [ebx], 2	<i>; Move the 16-bit integer representation of 2 ; into the 2 bytes starting at address EBX</i>
mov DWORD PTR [ebx], 2	<i>; Move the 32-bit integer representation of 2 ; into the 4 bytes starting at address EBX</i>

1.4 Instructions

Machine instructions generally fall into three categories: data movement, arithmetic/logic, and control-flow. In this section, we will look at important examples of x86 instructions from each category. This section should not be considered an exhaustive list of x86 instructions, but rather a useful subset.

In this section, we will use the following notation:

- `<reg32>` - means any 32-bit register described in Section 2, for example, ESI.
- `<reg16>` - means any 16-bit register described in Section 2, for example, BX.
- `<reg8>` - means any 8-bit register described in Section 2, for example AL.
- `<reg>` - means any of the above.
- `<mem>` - will refer to a memory address, as described in Section 1.3.2, for example `[EAX]`, or `[var+4]`, or `DWORD PTR [EAX+EBX]`.
- `<con32>` - means any 32-bit constant.
- `<con16>` - means any 16-bit constant.
- `<con8>` - means any 8-bit constant.
- `<con>` - means any of the above sized constants.

1.4.1 Data Movement Instructions

Instruction: mov

Syntax:

```
mov <reg>, <reg>
mov <reg>, <mem>
mov <mem>, <reg>
mov <reg>, <const>
mov <mem>, <const>
```

Semantics: The mov instruction moves the data item referred to by its second operand (i.e. register contents, memory contents, or a constant value) into the location referred to by its first operand (i.e. a register or memory). While register-to-register moves are possible, direct memory-to-memory moves are not. In cases where memory transfers are desired, the source memory contents must first be loaded into a register, then can be stored to the destination memory address.

Examples:

```
mov eax, ebx           ; transfer ebx to eax
mov BYTE PTR [var], 5  ; store the value 5 into the byte at
                        ; memory location ``var``
```

Instruction: push

Syntax:

```
push <reg32>
push <mem>
push <con32>
```

Semantics: The push instruction places its operand onto the top of the hardware supported stack in memory. Specifically, push first decrements ESP by 4, then places its operand into the contents of the 32-bit location at address `[ESP]`. ESP (the stack pointer) is decremented by push since the x86 stack grows down – i.e. the stack grows from high addresses to lower addresses.

Examples:

```
push eax      ; push the contents of eax onto the stack
push [var]    ; push the 4 bytes at address ``var`` onto stack
```


Instruction: imul

Syntax:

```
imul <reg32>, <reg32>
imul <reg32>, <mem>
imul <reg32>, <reg32>, <con>
imul <reg32>, <mem>, <con>
```

Semantics: The imul instruction has two basic formats: two-operand (first two syntax listings above) and three-operand (last two syntax listings above). The two-operand form multiplies its two operands together and stores the result in the first operand. The result (i.e. first) operand must be a register. The three operand form multiplies its second and third operands together and stores the result in its first operand. Again, the result operand must be a register. Furthermore, the third operand is restricted to being a constant value.

Examples:

```
imul eax, [var]          ; multiply the contents of EAX by the
                        ; 32-bit contents of the memory location
                        ; 'var'. Store the result in EAX.
imul esi, edi, 25        ; multiply the contents of EDI by 25.
                        ; Store the result in ESI.
```

Instruction: idiv

Syntax:

```
idiv <reg32>
idiv <mem>
```

Semantics: The idiv instruction is used to divide the contents of the 64 bit integer EDX:EAX (constructed by viewing EDX as the most significant four bytes and EAX as the least significant four bytes) by the specified operand value. The quotient result of the division is stored into EAX, while the remainder is placed in EDX. This instruction must be used with care. Before executing the instruction, the appropriate value to be divided must be placed into EDX and EAX. Clearly, this value is overwritten when the idiv instruction is executed.

Examples:

```
idiv ebx                ; divide the contents of EDX:EAX by the
                        ; contents of EBX. Place the quotient
                        ; in EAX and the remainder in EDX.
idiv DWORD PTR [var]    ; same as above, but divide by the
                        ; 32-bit value stored at memory
                        ; location 'var'.
```

Instruction: and, or, xor

Syntax:

```
and <reg>, <reg>      or <reg>, <reg>      xor <reg>, <reg>
and <reg>, <mem>      or <reg>, <mem>      xor <reg>, <mem>
and <mem>, <reg>      or <mem>, <reg>      xor <mem>, <reg>
and <reg>, <con>      or <reg>, <con>      xor <reg>, <con>
and <mem>, <con>      or <mem>, <con>      xor <mem>, <con>
```

Semantics: These instructions perform the specified logical operation (logical bitwise and, or, and exclusive or, respectively) on their operands, placing the result in the first operand location.

Examples:

```
and eax, 0fH           ; clear all but the last 4 bits of EAX.
xor edx, edx           ; set the contents of EDX to zero.
```

Instruction: not

Syntax: not <reg>
 not <mem>

Semantics: Performs the logical negation of the operand contents (i.e. flips all bit values).

Examples: not BYTE PTR [var] ; negate all bits in the byte at the
 ; memory location ``var``.

Instruction: neg

Syntax: neg <reg>
 neg <mem>

Semantics: Performs the arithmetic (i.e. two's complement) negation of the operand contents.

Examples: neg eax ; negate the contents of EAX.
 neg [var] ; negate the contents of ``var``

Instruction: shl, shr

Syntax: shl <reg>, <con8> shr <reg>, <con8>
 shl <mem>, <con8> shr <mem>, <con8>
 shl <reg>, cl shr <reg>, cl
 shl <mem>, cl shr <mem>, cl

Semantics: These instructions shift the bits in their first operand's contents left and right (shl and shr, respectively), padding the resulting empty bit positions with zeros. The shifted operand can be shifted up to 31 places. The number of bits to shift is specified by the second operand, which can be either an 8-bit constant or the register CL. In either case, shifts counts of greater than 31 are performed modulo 32.

Examples: shl eax 5 ; shift the contents of eax left by 5 bit
 ; positions
 shr [var] 3 ; shift the contents of ``var`` right by 3
 ; bit positions

1.4.3 Control Flow Instructions

In this section, we will refer to labeled locations in the program text as <label>. Labels can be inserted anywhere in x86 assembly code text by entering a label name followed by a colon. For example, consider the code fragment in Listing 1.6. The second instruction in this code fragment is labeled "begin". Elsewhere in the code, we can refer to the memory location that this instruction is located at in memory using the more convenient symbolic name "begin" instead of having to refer to the memory address as an integer.

Listing 1.6: x86 labeled code location

	mov esi, [ebp+8]
begin:	xor ecx, ecx
	mov eax, [esi]

Instruction: jmp

Syntax: `jmp <label>`

Semantics: Transfers program control flow to the instruction at the memory location indicated by the operand.

Examples: `jmp begin` ; jumps to the ``begin'' label

Instruction: jCC

Syntax: `je <label>` ; Jump when equal
`jne <label>` ; Jump when not equal
`jz <label>` ; Jump when last result was zero
`jg <label>` ; Jump when greater than
`jge <label>` ; Jump when greater than or equal to
`jle <label>` ; Jump when less than
`jlt <label>` ; Jump when less than or equal to

Semantics: These instructions are conditional jumps that are based on the status of a set of **condition codes** that are stored in a special register called the *machine status word*. The contents of the machine status word include information about the last arithmetic operation performed. For example, one bit of this word indicates if the last result was zero. Another indicates if the last result was negative. Based on these condition codes, a number of conditional jumps can be performed. For example, the `jz` instruction performs a jump to the specified operand label if the result of the last arithmetic operation (e.g. `add`, `sub`, etc.) was zero. Otherwise, control proceeds to the next instruction in sequence after the `jz`. These conditional jumps are the underlying support needed to implement high-level language features such as “if” statements and loops (e.g. “while” and “for”).

A number of the conditional branches are given names that are intuitively based on the last operation performed being a special compare instruction, `cmp` (see below). For example, conditional branches such as `jle` and `jne` are based on first performing a `cmp` operation on the desired operands.

Examples: `cmp eax, ebx` ; if the contents of `eax` are less than or
`jle done` ; equal to the contents of `EBX`, jump to the
; code location labeled ``done''.

Instruction: cmp

Syntax: `cmp <reg>, <reg>`
`cmp <reg>, <mem>`
`cmp <mem>, <reg>`
`cmp <reg>, <con>`
`cmp <mem>, <con>`

Semantics: Compares the two specified operands, setting the condition codes in the machine status word appropriately. In fact, this instruction is equivalent to the `sub` instruction, except the result of the subtraction is discarded.

Examples: `cmp DWORD PTR [var], 10` ; if the 4 bytes stored at memory
`jeq loop` ; location ``var'' equal the 4-byte
; integer value 10, then jump to the
; code location labeled `loop`

(i.e. integer-sized) memory region labeled “var”.

Next in each file, we find the declaration of the function named `returnTwo`. In the C++ file we have declared the function to be `extern “C”`. This declaration indicates that the C++ compiler should use C naming conventions when labeling the function `returnTwo` in the resulting object file that it produces. In fact, this naming convention means that the function `returnTwo` should map to the label `_returnTwo` in the object code. In the assembly code, we have labeled the beginning of the subroutine `_returnTwo` using the `PROC` directive, and have declared the label `_returnTwo` to be public. Again, the result of these actions will be that the subroutine will map to the symbol `_returnTwo` in the object code that the assembler generates.

The function bodies are straight-forward. As we will see in more detail in Section 6, return values for functions are placed into EAX by convention, hence the instruction to move the contents of “var” into EAX in the assembly code.

Given these equivalent function definitions, use of either version of the function is the same. A sample call to the function `returnTwo` is depicted in Listing 1.9. This C++ code could be linked to either definition of the function and would produce the same results (note, we could not link to both definitions, or the linker would produce a “multiply defined symbol” error. The mechanics of program linking will be discussed in an associated document that relates to the specific programming environment that you will use to assemble and run programs.

Listing 1.8: C++ code to return 2

```
int var = 2;

extern "C" returnTwo();

int returnTwo() {
    return var;
}
```

Listing 1.9: Calling `returnTwo()` from C++

```
#include <iostream>
using namespace std;

extern "C" int returnTwo();

int main() {
    cout << "calling _returnTwo() _returned: "
         << returnTwo() << endl;
    return 0;
}
```


Chapter 2

x86 32-bit C Calling Convention

...

This chapter was derived from a document written by Adam Ferrari and later updated by Alan Batson, Mike Lack and Anita Jones

2.1 What is a Calling Convention?

In Section 1.5 we saw a simple example of a subroutine defined in x86 assembly language. In fact, this subroutine was quite simple – it did not modify any registers except EAX (which was needed to return the result), and it did not call any other subroutines. In practice, such simple function definitions are rarely useful. When more complex subroutines are combined in a single program, a number of complicating issues arise. For example, how are parameters passed to a subroutine? Can subroutines overwrite the values in a register, or does the caller expect the register contents to be preserved? Where should local variables in a subroutine be stored? How should results be returned from functions?

To allow separate programmers to share code and develop libraries for use by many programs, and to simplify the use of subroutines in general, programmers typically adopt a common *calling convention*. The calling convention is simply a set of rules that answers the above questions without ambiguity to simplify the definition and use of subroutines. For example, given a set of calling convention rules, a programmer need not examine the definition of a subroutine to determine how parameters should be passed to that subroutine. Furthermore, given a set of calling convention rules, high-level language compilers can be made to follow the rules, thus allowing hand-coded assembly language routines and high-level language routines to call one another.

In practice, even for a single processor instruction set, many calling conventions are possible. In this class we will examine and use one of the most important conventions: the C language calling convention. Understanding this convention will allow you to write assembly language subroutines that are safely callable from C and C++ code, and will also enable you to call C library functions from your assembly language code.

2.2 The C Calling Convention

The C calling convention is based heavily on the use of the hardware-supported stack. To understand the C calling convention, you should first make sure that you fully understand the push, pop, call, and ret instructions – these will be the basis for most of the rules. In this calling convention, subroutine parameters are passed on the stack. Registers are saved on the stack, and local variables used by subroutines are placed in memory on the stack. In fact, this stack-centric implementation of subroutines is not unique to the C language or the x86 architecture. The vast majority of high-level procedural languages implemented on most processors have used similar calling convention.

The calling convention is broken into two sets of rules. The first set of rules is employed by the caller of the subroutine, and the second set of rules is observed by the writer of the subroutine (the “callee”). It should be emphasized that mistakes in the observance of these rules quickly result in fatal program errors; thus meticulous care should be used when implementing the call convention in your own subroutines.

2.3 The Caller’s Rules

The caller should adhere to the following rules when invoking a subroutine:

1. Before calling a subroutine, the caller should save the contents of certain registers that are designated caller-saved. The caller-saved registers are EAX, ECX, EDX. If you want the contents of these registers to be preserved across the subroutine call, push them onto the stack.
2. To pass parameters to the subroutine, push them onto the stack before the call. The parameters should be pushed in inverted order (i.e. last parameter first) – since the stack grows down, the first parameter will be stored at the lowest address (this inversion of parameters was historically used to allow functions to be passed a variable number of parameters).
3. To call the subroutine, use the call instruction. This instruction places the return address on top of the parameters on the stack, and branches to the subroutine code.
4. After the subroutine returns, (i.e. immediately following the call instruction) the caller must remove the parameters from stack. This restores the stack to its state before the call was performed.
5. The caller can expect to find the return value of the subroutine in the register EAX.
6. The caller restores the contents of caller-saved registers (EAX, ECX, EDX) by popping them off of the stack. The caller can assume that no other registers were modified by the subroutine.

2.4 The Callee’s Rules

The definition of the subroutine should adhere to the following rules:

1. At the beginning of the subroutine, the function should push the value of EBP onto the stack, and then copy the value of ESP into EBP using the following instructions:

Listing 2.1: x86 callee code, part 1

```
push ebp
mov ebp, esp
```

The reason for this initial action is the maintenance of the base pointer, EBP. The base pointer is used by convention as a point of reference for finding parameters and local variables on the stack. Essentially, when any subroutine is executing, the base pointer is a “snapshot” of the stack pointer value from when the subroutine started executing. Parameters and local variables will always be located at known, constant offsets away from the base pointer value. We push the old base pointer value at the beginning of the subroutine so that we can later restore the appropriate base pointer value for the caller when the subroutine returns. Remember, the caller isn't expecting the subroutine to change the value of the base pointer. We then move the stack pointer into EBP to obtain our point of reference for accessing parameters and local variables.

2. Next, allocate local variables by making space on the stack. Recall, the stack grows down, so to make space on the top of the stack, the stack pointer should be decremented. The amount by which the stack pointer is decremented depends on the number of local variables needed. For example, if 3 local integers (4 bytes each) were required, the stack pointer would need to be decremented by 12 to make space for these local variables. I.e:

Listing 2.2: x86 callee code, part 2

```
sub esp, 12
```

As with parameters, local variables will be located at known offsets from the base pointer.

3. Next, the values of any registers that are designated callee-saved that will be used by the function must be saved. To save registers, push them onto the stack. The callee-saved registers are EBX, EDI and ESI (ESP and EBP will also be preserved by the call convention, but need not be pushed on the stack during this step).

After these three actions are performed, the actual operation of the subroutine may proceed. When the subroutine is ready to return, the call convention rules continue:

4. When the function is done, the return value for the function should be placed in EAX if it is not already there.
5. The function must restore the old values of any callee-saved registers (EBX, EDI and ESI) that were modified. The register contents are restored by popping them from the stack. Note, the registers should be popped in the inverse order that they were pushed.
6. Next, we deallocate local variables. The obvious way to do this might be to add the appropriate value to the stack pointer (since the space was allocated by subtracting the needed amount from the stack pointer). In practice, a less error-prone way to deallocate the variables is to move the value in the base pointer into the stack pointer, i.e.:

Listing 2.3: x86 callee code, part 3

```
mov esp, ebp
```

This trick works because the base pointer always contains the value that the stack pointer contained immediately prior to the allocation of the local variables.

7. Immediately before returning, we must restore the caller's base pointer value by popping EBP off the stack. Remember, the first thing we did on entry to the subroutine was to push the base pointer to save its old value.
8. Finally, we return to the caller by executing a ret instruction. This instruction will find and remove the appropriate return address from the stack.

It might be noted that the callee's rules fall cleanly into two halves that are basically mirror images of one another. The first half of the rules apply to the beginning of the function, and are therefor commonly

Listing 2.4: Example function call, caller's rules obeyed

```

; Want to call a function "myFunc" that takes three
; integer parameters. First parameter is in EAX.
; Second parameter is the constant 123. Third
; parameter is in memory location "var"

push [var]      ; Push last parameter first
push 123
push eax        ; Push first parameter last

call myFunc ; Call the function (assume C naming)

; On return, clean up the stack. We have 12 bytes
; (3 parameters * 4 bytes each) on the stack, and the
; stack grows down. Thus, to get rid of the parameters,
; we can simply add 12 to the stack pointer

add esp, 12

; The result produced by "myFunc" is now available for
; use in the register EAX. No other register values
; have changed

```

said to define the *prologue* to the function. The latter half of the rules apply to the end of the function, and are thus commonly said to define the *epilogue* of the function.

2.5 Calling Convention Example

The above rules may seem somewhat abstract on first examination. In practice, the rules become simple to use when they are well understood and familiar. To start the process of better understanding the call convention, we now examine a simple example of a subroutine call and a subroutine definition.

In Listing 2.4 a sample function call is depicted. Note how the caller pushes the parameters onto the stack in inverted order before the call. The call instruction is used to jump to the beginning of the subroutine in anticipation of the fact that the subroutine will use the ret instruction to return when the subroutine completes. When the subroutine returns, the parameters must be removed from the stack. A simple way to do this is to add the appropriate amount to the stack pointer (since the stack grows down). Finally, the result is available in EAX. Relative to the caller's rules, the callee's rules are somewhat more complex. An example subroutine implementation that obeys the callee's rules is depicted in ListingFigure 8. The subroutine prologue performs the standard actions of saving a snapshot of the stack pointer in EBP (the base pointer), allocating local variables by decrementing the stack pointer, and saving register values on the stack.

In the body of the subroutine we can now more clearly see the use of the base pointer illustrated. Both parameters and local variables are located at constant offsets from the base pointer for the duration of the

Listing 2.5: Example function definition, callee's rules obeyed

```

global myFunc

section .text

myFunc:
    ; *** Standard subroutine prologue ***
    push ebp          ; Save the old base pointer value.
    mov ebp, esp      ; Set the new base pointer value.
    sub esp, 4        ; Make room for one 4-byte local variable.
    push edi          ; Save the values of registers that the function
    push esi          ; will modify. This function uses EDI and ESI.
                    ; (no need to save EAX, EBP, or ESP)

    ; *** Subroutine Body ***
    mov eax, [ebp+8]   ; Put value of parameter 1 into EAX
    mov esi, [ebp+12]  ; Put value of parameter 2 into ESI
    mov edi, [ebp+16]  ; Put value of parameter 3 into EDI

    mov [ebp-4], edi   ; Put EDI into the local variable
    add [ebp-4], esi   ; Add ESI into the local variable
    add eax, [ebp-4]   ; Add the contents of the local variable
                    ; into EAX (final result)

    ; *** Standard subroutine epilogue ***
    pop esi           ; Recover register values
    pop edi
    mov esp, ebp      ; Deallocate local variables
    pop ebp           ; Restore the caller's base pointer value
    ret

```

subroutines execution. In particular, we notice that since parameters were placed onto the stack before the subroutine was called, they are always located below the base pointer (i.e. at higher addresses) on the stack. The first parameter to the subroutine can always be found at memory location `[EBP+8]`, the second at `[EBP+12]`, the third at `[EBP+16]`, and so on. Similarly, since local variables are allocated after the base pointer is set, they always reside above the base pointer (i.e. at lower addresses) on the stack. In particular, the first local variable is always located at `[EBP-4]`, the second at `[EBP-8]`, and so on. Understanding this conventional use of the base pointer allows us to quickly identify the use of local variables and parameters within a function body.

The function epilogue, as expected, is basically a mirror image of the function prologue. The caller's register values are recovered from the stack, the local variables are deallocated by resetting the stack pointer, the caller's base pointer value is recovered, and the `ret` instruction is used to return to the appropriate code location in the caller.

A good way to visualize the operation of the calling convention is to draw the contents of the nearby region of the stack during subroutine execution. Figure 2.1 depicts the contents of the stack during the

execution of the body of `myFunc` (depicted in Listing 2.5). Notice, lower addresses are depicted lower in the figure, and thus the “top” of the stack is the bottom-most cell. This corresponds visually to the intuitive statement that the x86 hardware stack “grows down.” The cells depicted in the stack are 32-bit wide memory locations, thus the memory addresses of the cells are 4 bytes apart. From this picture we see clearly why the first parameter resides at an offset of 8 bytes from the base pointer. Above the parameters on the stack (and below the base pointer), the call instruction placed the return address, thus leading to an extra 4 bytes of offset from the base pointer to the first parameter.

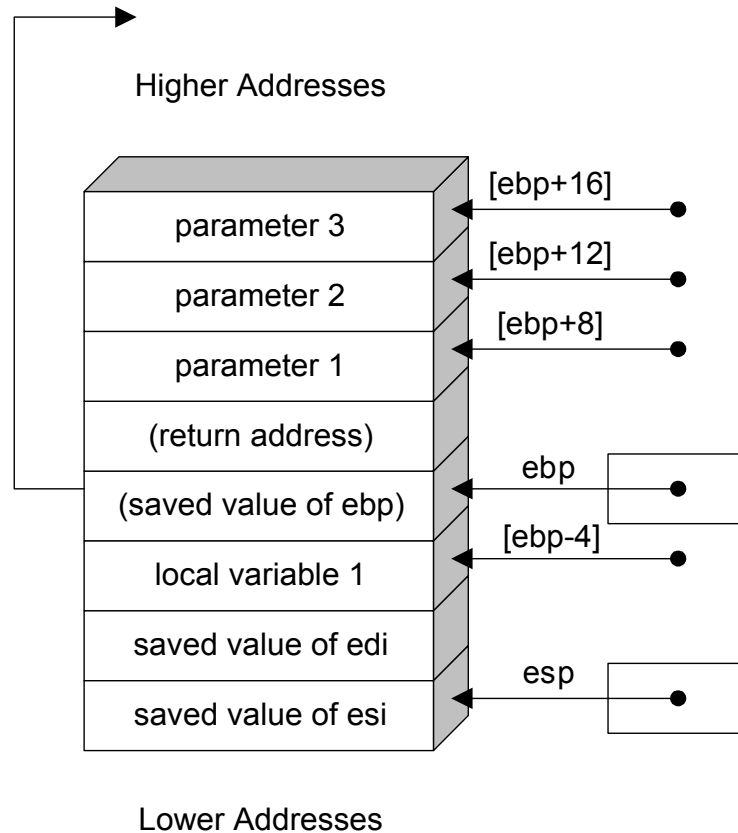


Figure 2.1: A picture of the stack in memory during the execution of the body of `myFunc`

The assembly code for `myFunc()` was shown above in Listing 2.5). The C++ code to call that subroutine is shown in Listing 2.6.

Listing 2.6: Example C++ code to invoke a 3-parameter x86 subroutine

```
#include <iostream>
using namespace std;

extern "C" int myFunc(int ,int ,int );

int main() {
    int x = 3;
    cout << "myFunc() returned:_"
         << myFunc(x,5,10) << endl;
    return 0;
}
```