

Gnu Assembler

The Gnu assembler(Gas) is an open source Assembler. The assembler is included as standard in the gcc distribution and is available for download to run under Windows. It provides support for the instructionset of the host CPU. Cross assembler versions are also available.

The assembler is invoked by the command `as`.

`as` provides support for the object module format used in whichever gcc distribution it comes with. If you are programming in assembler, Gas provides a more complete range of instructions, in association with better portability between operating systems than competing assemblers.

It is beyond the scope of this course to provide a complete guide to assembler programming

for the Intel processor family. Readers wanting a general background in assembler programming should consult appropriate text books [] in conjunction with the processor reference manuals published by Intel[?][?] and AMD [?].

General instruction syntax

Assembler programs take the form of a sequence of lines with one machine instruction per line. The instructions themselves take the form of an optional label, an operation code name conditionally followed by up to three comma separated operands. For example:

```
l1: cltq           # 0 operand instr
    pop %rax      # 1 operand instr
    mov %rax,%r9  # 2 operand instr
    imulq $3,%r9,%r10 # 3 operand instr
```

As shown above, a comment can be placed on an assembler line, with the comment distinguished from the instruction by a leading `#`. The label, if present is separated from the operation code name by a colon.

In the `as` assembler, unlike in the original Intel assembler, the direction of assignment in an instruction is always from left to right*, so that

```
mov %rax,%rcx
```

means

$$\%rax \rightarrow \%rcx$$

and

```
addss %xmm0,%xmm3
```

means

$$(\%xmm0 + \%xmm3) \rightarrow \%xmm3$$

*This is a result of `as` having originated as a Motorola assembler that was converted to recognise Intel opcodes. Motorola follow a left to right assignment convention.

Operand forms

Operands to instructions can be constants, register names or memory locations.

Constants

Constants are values known at assembly time, and take the form of numbers, labels, characters or arithmetic expressions whose components are themselves constants.

The most important constant values are numbers. A constant must be prefixed by a \$

```
mov $7,%rax # load 7 into register rax
```

Only integer constants can occur in instructions. Floating point constants can occur in data declarations but not in instructions.

Floating point constants are also supported as operands to store allocation directives (see section ??):

```
.double 3.14156  
.float 9.2e3
```

It is important to realise that due to limitations of the AMD and Intel instruction-sets, floating point constants can not be directly used as operands to instructions. Any floating point constants used in an algorithm have to be assembled into a distinct area of memory and loaded into registers from there.

Labels

Constants can also take the form of labels. As the assembler program is processed, the assembler allocates an integer value to each label.

We can load a register with the address referred to by a label by including the label as a constant operand:

```
mov  $ sourcebuf,%rsi
```

You must not confuse this with

```
mov  sourcebuf,%rsi
```

which loads the register with the word in memory labeled by `sourcebuf`.

Constant expressions

Suppose there exists a data-structures for which one has a base address label, it is often convenient to be able to refer to fields within this structure in terms of their offset from the start of the structure. Consider the example of a vector of 4 single precision floating point values at a location with label `myvec`. The actual address at which `myvec` will be placed is determined by Gas, we do not know it. We may know that we want the address of the 3rd element of the vector:

```
mov  $ myvec + $3 *$4, %esi
```

will place the address of this word into the `esi` register.

Constant expressions

Gas allows one to place arithmetic expressions whose sub-expressions are constants wherever a constant can occur. The arithmetic operators are written C style as shown below.

operator	means	operator	means
	bitwise or	+	add
^	bitwise xor	-	subtract
&	bitwise and	*	multiply
<<	shift left	/	signed division
>>	shift right		
%	remainder		

Registers

Operands can be register names. The available register names are shown in table ???. In the binary operation codes interpreted by the CPU, registers are identified using 3-bit integers. Depending on the operation code, these 3 bit fields are interpreted as the different categories of register shown in table ???.

Register names should be preceded by % in the assembler syntax to distinguish them from label names.

You should be aware that in the Intel architecture a number of registers are aliased to the same state vectors, thus for example the `eax`, `ax`, `al`, `ah` registers all share bits. More insidiously the floating point registers `ST0..ST7` not only share state with the MMX registers,

but their mapping to these registers is dynamic and variable.

The first 8 registers have names which they inherit from the Pentium processors, the AMD opteron added additional registers to the general register set and simd set.

	general registers				fpu	mmx	SIMD
#	8 reg	16 reg	32 reg	64 reg	80 stack	64 reg	128 reg
	Aliased				Aliased		
0	al	ax	eax	rax	st0	mm0	xmm0
1	cl	bx	ecx	rcx	st1	mm1	xmm1
2	dl	cx	edx	rdx	st2	mm2	xmm2
3	bl	bx	ebx	rbx	st3	mm3	xmm3
4	ah	sp	esp	rsp	st4	mm4	xmm4
5	ch	bp	ebp	rbp	st5	mm5	xmm5
6	dh	si	esi	rsi	st6	mm6	xmm6
7	bh	di	edi	rdi	st7	mm7	xmm7
8..15	r8..r15						xmm8..xmm15

Memory Locations

Memory locations are syntactically represented by labels, addresses or locations pointed to by registers: `100` , `myvec` , `0(%rsi)` all represent memory locations.

The address expressions, unlike constant expressions, can contain components whose values are not known until program execution. The final example above refers to the memory location addressed by the value in the `%rsi` register, and as such, depends on the history of prior computations affecting that register.

Here are the register addressing forms

- $n(\text{reg})$ address $\text{reg} + n$,
 - eg `8(%rdx)`
- $n(\text{reg},m)$ address $n + \text{reg} * m$
 - eg `16(%rax,2)` , m must be 2,4,8
- $n(\text{reg1},\text{reg2})$ address $n + \text{reg1} + \text{reg2}$,
 - eg `20(%rax,%rsi)`
- $n(\text{reg1},\text{reg2},m)$ address $n + \text{reg1} + \text{reg2} * m$
 - eg `100(rbx,%r10,8)` , m must be 2,4 or 8

Sectioning

Programs running under Linux have their memory divided into 4 sections:

- text** is the section of memory containing operation codes to be executed. It is typically mapped as read only by the paging system.
- data** is the section of memory containing initialised global variables, which can be altered following the start of the program.
- bss** is the section containing uninitialised global variables.

`stack` is the section in which dynamically allocated local variables of subroutines are located.

The `.section` directive is used by assembler programmers to specify into which section of memory they want subsequent lines of code to be assembled. For example in the listing shown in algorithm 1 we divide the program into three sections: a `text` section containing `myfunc`, a `bss` section containing 64 undefined bytes and a `data` section containing a vector of 4 integers.

The label `myfuncbase` can be used with *negative* offsets to access locations within the `bss`, whilst the label `myfuncglobal` can be used with *positive* offsets to access elements of the vector in the `data` section.

Algorithm 1 Examples of the use of section and data reservation directives

```
.section .text
    .global myfunc
myfunc:enter $128,$0
# body of function goes here
    leave
    ret $0

.section .bss
    .align 16
    .space 64 # reserve 64 bytes

myfuncBase:
.section .data
myfuncglobal: # reserve 4 by 32-bit integers
    .int 1
    .int 2
    .int 3
    .int 5
```

Data reservation

Data must be reserved in distinct ways in the different sections. In the data section, the data definition directives `.byte`, `.word`, `.int`, and `.quad` are used to define bytes, words(16bit), doublewords(32bit) and quad words (64bit). The directive must be followed by a constant expression. When defining bytes or words the constant must be an integer. `.float` and `.double` may be used to define floating point constants as shown previously.

In the bss section the directive `.space` is used to reserve a specified number of bytes, but no value is associated with these bytes.

Stack data

Data can be allocated in the stack section by use of the `enter` operation code name. This takes the form:

```
enter space, level
```

It should be used as the first operation code name of a function. The `level` parameter is only of relevance in block structured languages and should be set to 0 for assembler programming. The `space` parameter specifies the number of bytes to be reserved for the private use of the function. Once the `enter` instruction has executed, the data can be accessed at *negative* offsets from the `rbp` register.

Releasing stack space dynamically

The last two instructions in a function should, as shown in algorithm be

```
leave  
ret $0
```

The combined effect of these is to free the space reserved on the stack by `enter`, and pop the return address from the stack. The parameter to the operation code name `ret` is used to specify how many bytes of function parameters should be discarded from the stack. If one is interfacing to C this should always be set to 0.

Label qualification

The default scope of a label is the assembler source file containing the line it prefixes. But labels can be used to mark the start of functions that are to be called from C or other high level languages. To indicate that they have scope beyond the current assembler file, the `.global` directive should be used as shown in algorithm 1.

The converse case, where an assembler file calls a function exported by a C program is handled by the `.extern` directive:

```
.extern printreal
```

```
call printreal
```

in the above example we assume that `printreal` is a C function called from assembler.

Linking and object file formats

There are 4 object file formats that are commonly used on Linux and Windows systems as shown in table ???. This lists the name of the format, its file extension - which is often ambiguous and the combination of operating system and compiler that makes use of it. A flag provided to Gas specifies which format it should use. We will only go into the use of the gcc compiler, since this is portable between Windows and Linux.

Let us assume we have a C program called `c2asm.c` and an assembler file `asmfromc.asm`. Suppose we wish to combine these into a single executable module `c2asm`. We issue the following commands at the console:

```
as -o asmfromc.o asmfromc.asm
```

```
gcc -oc2asm c2asm.c asmfromc.o
```