The following excerpts are taken from the User Manual for Intel's 8086 family architecture released in the late 1970s. In this problem, we will try to reason about the 8086 architecture to virtualize CPU and memory. We have condensed the 208-page manual to mere 11 pages for this problem. Chapter numbers and section numbers are kept consistent with the manual. So don't be surprised when you see Section 2.2 without seeing Section 2.1. *Occasional text in italics are additions made by the instructor for the problem.*

# Chapter 1: Introduction

**Microprocessors**:

At the core of the product line are three microprocessors *(including 8086)* that share these characteristics:

- Standard operating speed is 5 MHz (200ns cycle time); a selected 8 MHz version of 8086 CPU is also available.
- Chips are housed in reliable 40-pin packages.
- Up to 1 megabyte of memory can be addressed, along with a separate 64k bytes I/O space.

# Chapter 2: The 8086 and 8088 Central Processing Unit

Both CPUs are substantially more powerful than any microprocessor previously offered by Intel. Actual performance, of course, varies from application to application, but comparisons to the industry standard 2-MHz 8080A are instructive. The 8088 is from four to six times more powerful than the 8080A; the 8086 provides seven to ten times the 8080A's performance.

The high performance of the 8086 and 8088 is realized by combining a 16-bit internal data path with a pipelined architecture that allows instructions to be prefetched during spare bus cycles. Also contributing to performance is a compact instruction format that enables more instructions to be fetched in a given amount of time.

## 2.2  Process Architecture

Microprocessors generally execute a program by repeatedly cycling through the steps shown below (this description is somewhat simplified):

1. Fetch the next instruction from memory.
2. Read an operand (if required by the instruction).
3. Execute the instruction.
4. Write the result (if required by the instruction).

In previous CPUs, most of these steps have been performed serially, or with only a single bus cycle fetch overlap. The architecture of the 8086 and 8088 CPUs, while performing the

same steps, allocates them to two separate processing units within the CPU. The execution unit (EU) executes instructions; the bus interface unit (BIU) fetches instructions, reads operands, and writes results.

The two units can operate independently of one another and are able, under most circumstances, to extensively overlap instruction fetch with execution. The result is that, in most cases, the time normally required to fetch instructions "disappears" because the EU executes instructions that have already been fetched by the BIU.

**Execution Unit (EU)**

The execution units of the 8086 and 8088 are identical (figure 2-6). A 16-bit arithmetic/logic unit (ALU) in the EU maintains the CPU status and control flags, and manipulates the general registers and instruction operands. All registers and data paths in the EU are 16 bits wide for fast internal transfers.
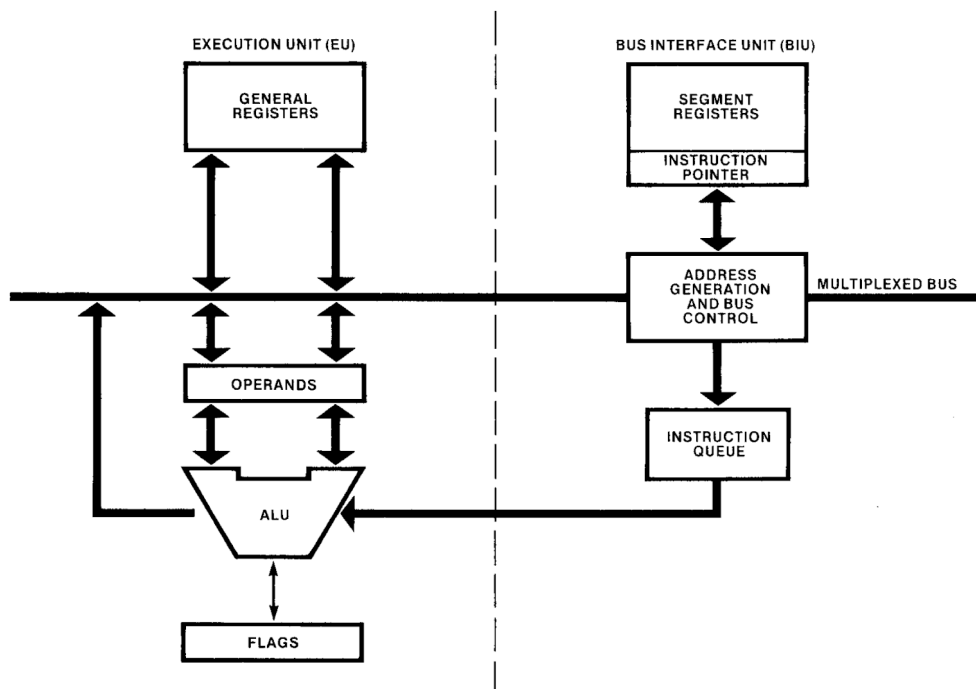


Figure 2-6. Execution and Bus Interface Units (EU and BIU)

The EU has no connection to the system bus, the "'outside world." It obtains instructions from a queue maintained by the BIU. Likewise, when an instruction requires access to memory or to a peripheral device, the EU requests the BIU to obtain or store the data. All addresses manipulated by the EU are 16 bits wide. The BIU, however, performs an address relocation that gives the EU access to the full megabyte of memory space (see section 2.3).

**Bus Interface Unit (BIU)**

The BIU performs all bus operations for the EU. Data is transferred between the CPU and memory or I/O devices upon demand from the EU. Sections 2.3 and 2.4 describe the

interaction of the BIU with memory and I/O devices *(we will not see 2.4 for this problem since we're only interested in CPU/memory virtualization)*.

**General registers**

Both CPUs have the same complement of eight 16-bit general registers (figure 2-7). The general registers are subdivided into two sets of four registers each: the data registers (sometimes called the H & L group for ""high" and "low"), and the pointer and index registers (sometimes called the P & I group).
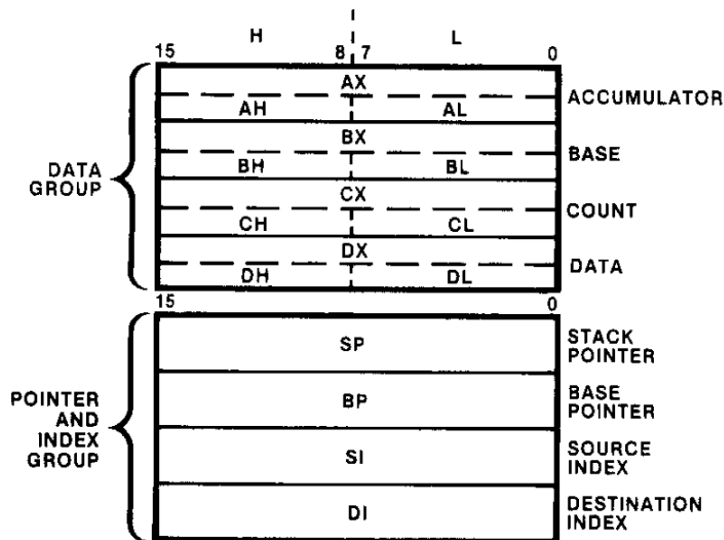


Figure 2-7. General Registers

The data registers are unique in that their upper (high) and lower halves are separately addressable. This means that each data register can be used interchangeably as a 16-bit register, or as two 8-bit registers. The other CPU registers always are accessed as 16-bit units only. The data registers can be used without constraint in most arithmetic and logic operations. In addition, some instructions use certain registers implicitly (see table 2-1) thus allowing compact yet powerful encoding.

Table 2-1: Implicit use of general registers

| REGISTER | OPERATIONS |
|---|---|
| AX | Word Multiply, Word Divide, Word I/O |
| AL | Byte Multiply, Byte Divide, Byte I/O, Translate, Decimal Arithmetic |
| AH | Byte Multiply, Byte Divide |
| BX | Translate |
| CX | String Operations, Loops |

| CL | Variable Shift and Rotate |
|----|---------------------------|
| DX | Word Multiply, Word Divide, Indirect I/O |
| SP | Stack Operations |
| SI | String Operations |
| DI | String Operations |

**Segment Registers**

The megabyte of 8086 and 8088 memory space is divided into logical segments of up to 64k bytes each. (Memory segmentation is described in section 2.3.) The CPU has direct access to four segments at a time; their base addresses (starting locations) are contained in the segment registers (see figure 2-8), The CS register points to the current code segment; instructions are fetched from this segment. The SS register points to the current stack segment; stack operations are performed on locations in this segment. The DS register points to the current data segment; it generally contains program variables. The ES register points to the current extra segment, which also is typically used for data storage.
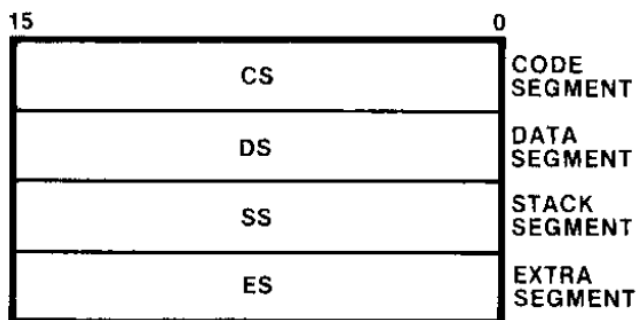


# Figure 2-8. Segment Registers

The segment registers are accessible to programs and can be manipulated with several instructions. Good programming practice and consideration of compatibility with future Intel hardware and software products dictate that the segment registers be used in a disciplined fashion.

**Instruction Pointer**

The 16-bit instruction pointer (IP) is analogous to the program counter (PC) in the 8080/8085 CPUs. The instruction pointer is updated by the BIU so that it contains the offset (distance in bytes) of the next instruction from the beginning of the current code segment; i.e., IP points to the next instruction. During normal execution, IP contains the offset of the next instruction to be fetched by the BIU; whenever IP is saved on the stack, however, it first is automatically adjusted to point to the next instruction to be executed. Programs do not have direct access

to the instruction pointer, but instructions cause it to change and to be saved on and restored from the stack.

**Flags**

The 8086 and 8088 have six 1-bit status flags (figure 2-9) that the EU posts to reflect certain properties of the result of an arithmetic or logic operation. A group of instructions is available that allows a program to alter its execution depending on the state of these flags, that is, on the result of a prior operation.
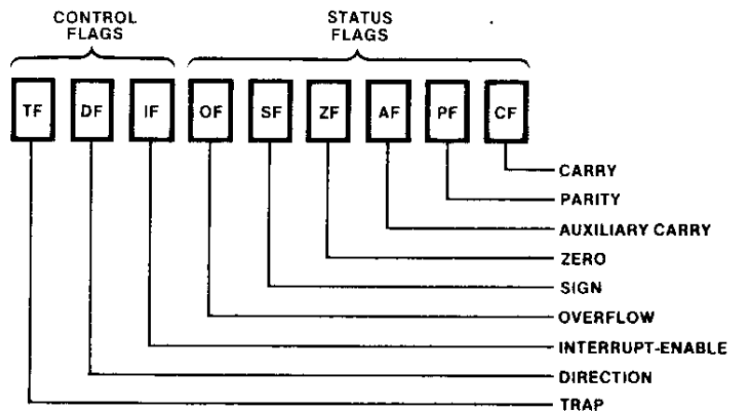


## Figure 2-9. Flags

Different instructions affect the status flags differently; in general, however, the flags reflect the following conditions:

1. If AF (the auxiliary carry flag) is set, there has been a carry out of the low nibble into the high nibble or a borrow from the high nibble into the low nibble of an 8-bit quantity (low-order byte of a 16-bit quantity). This flag is used by decimal arithmetic instructions.

2. If CF (the carry flag) is set, there has been a carry out of, or a borrow into, the high-order bit of the result (8- or 16-bit). The flag is used by instructions that add and subtract multibyte numbers.

3. If OF (the overflow flag) is set, an arithmetic overflow has occurred; that is, a significant digit has been lost because the size of the result exceeded the capacity of its destination location. An Interrupt On Overflow instruction is available that will generate an interrupt in this situation *(We will not worry about this instruction for this problem)*.

4. If SF (the sign flag) is set, the high-order bit of the result is a 1. Since negative binary numbers are represented in the 8086 and 8088 in standard two's complement notation, SF indicates the sign of the result (0 = positive, 1 = negative).

5. If PF (the parity flag) is set, the result has even parity, an even number of 1-bits. This flag can be used to check for data transmission errors.

6. If ZF (the zero flag) is set, the result of the operation is 0.

Three additional control flags (figure 2-9) can be set and cleared by programs to alter processor operations:

1. Setting DF (the direction flag) causes string instructions to auto-decrement; that is, to process strings from high addresses to low addresses, or from "right to left." Clearing DF causes string instructions to autoincrement, or to process strings from "left to right." *(We will not worry about string instructions for this problem).*

2. Setting IF (the interrupt-enable flag) allows the CPU to recognize external (maskable) interrupt requests. Clearing IF disables these interrupts. IF has no effect on either non-maskable external or internally generated interrupts.

3. Setting TF (the trap flag) puts the processor into a single-step mode for debugging. In this mode, the CPU automatically generates an internal interrupt after each instruction, allowing a program to be inspected as it executes instruction by instruction. Section 2.10 contains an example showing the use of TF in a single-step and breakpoint routine. *(We will not be seeing Section 2.10 details for this problem).*

## 2.3 Memory

This section describes how memory is functionally organized and used.

From a storage point of view, the 8086 and 8088 memory spaces are organized as identical arrays of 8-bit bytes. Instructions, byte data and word data may be freely stored at any byte address without regard for alignment thereby saving memory space by allowing code to be densely packed in memory. Instruction alignment does not materially affect the performance of either processor.

**Segmentation**

8086 and 8088 programs "view" the megabyte of memory space as a group of segments that are defined by the application. A segment is a logical unit of memory that may be up to 64k bytes long. Each segment is made up of contiguous memory locations and is an independent, separately-addressable unit. Every segment is assigned (by software) a base address, which is its starting location in the memory space. All segments begin on 16-byte memory boundaries. There are no other restrictions on segment locations; segments may be adjacent, disjoint, partially overlapped, or fully overlapped (see Figure 2-15). A physical memory location may be mapped into (contained in) one or more logical segments.
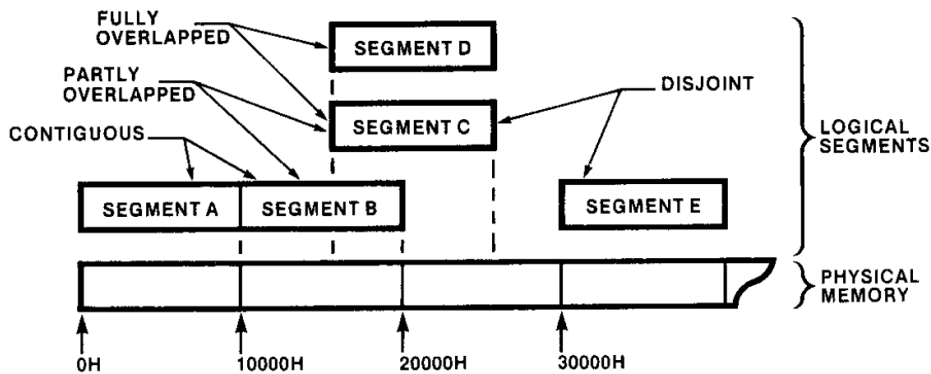
**Figure 2-15. Segment Locations in Physical Memory**

The segment registers point to (contain the base address values of) the four currently addressable segments (see figure 2-16). Programs obtain access to code and data in other segments by changing the segment registers to point to the desired segments.

Every application will define and use segments differently. The currently addressable segments provide a generous work space: 64k bytes for code, a 64k byte stack and 128k bytes of data storage. Many applications can be written to simply initialize the segment registers and then forget them. Larger applications should be designed with careful consideration given to segment definition.
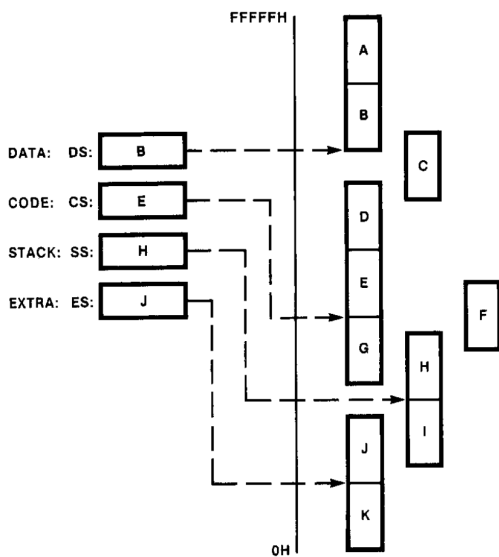


Figure 2-16. Currently Addressable Segments

## Physical Address Generation

It is useful to think of every memory location as having two kinds of addresses, physical and logical. A physical address is the 20-bit value that uniquely identifies each byte location in the megabyte memory space. Physical addresses may range from 0H through FFFFFH. All exchanges between the CPU and memory components use this physical address.

Programs deal with logical, rather than physical addresses and allow code to be developed without prior knowledge of where the code is to be located in memory and facilitate dynamic management of memory resources. A logical address consists of a segment base value and an offset value. For any given memory location, the segment base value locates the first byte of the containing segment and the offset value is the distance, in bytes, of the target location from the beginning of the segment. Segment base and offset values are unsigned 16-bit quantities; the lowest-addressed byte in a segment has an offset of 0. Many different logical addresses can map to the same physical location as shown in figure 2-17. In figure 2-17, physical memory location 2C3H is contained in two different overlapping segments, one beginning at 2B0H and the other at 2C0H.
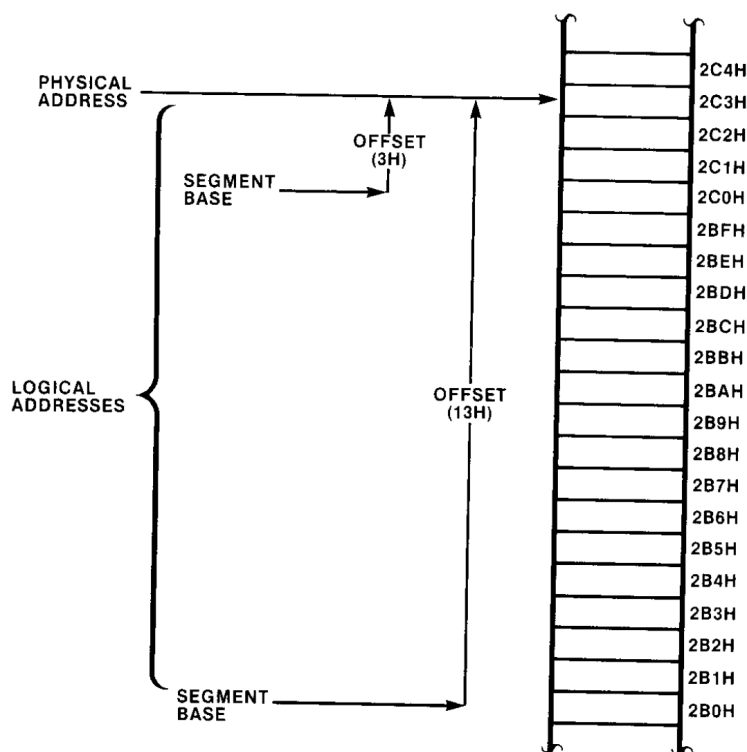


Figure 2-17. Logical and Physical Addresses

Whenever the BIU accesses memory—to fetch an instruction or to obtain or store a variable—it generates a physical address from a logical address. This is done by shifting the segment base value four-bit positions and adding the offset as illustrated in Figure 2-18. Note that this addition process provides for modulo 64k addressing (addresses wrap around from the end of a segment to the beginning of the same segment).
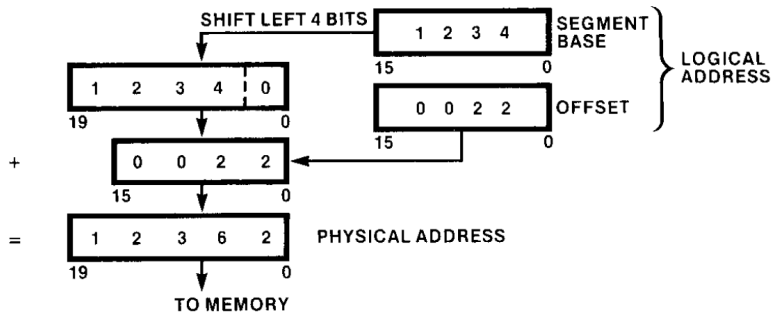
Figure 2-18. Physical Address Generation

The BIU obtains the logical address of a memory location from different sources depending on the type of reference that is being made (see table 2-2). Instructions always are fetched from the current code segment; IP contains the offset of the target instruction from the beginning of the code segment. Stack instructions always operate on the current stack segment; SP contains the offset of the top of the stack. Variables (memory operands) are assumed to reside in the current data segment[1].

Table 2-2 Logical Address Sources *(Simplified)*

| Type of memory reference | Segment base | OFFSET |
|---|---|---|
| Instruction Fetch | CS | IP |
| Stack Operation | SS | SP |
| Variable | DS | Address in instruction |
| BP Used as base register | SS | Address in instruction |

When register BP, the base pointer register, is designated as a base register in an instruction, the variable is assumed to reside in the current stack segment. Register BP thus provides a convenient way to address data on the stack.

**Dynamically Relocatable Code**

The segmented memory structure of the 8086 and 8088 makes it possible to write programs that are position-independent, or dynamically relocatable. Dynamic relocation allows a multiprogramming or multitasking system to make particularly effective use of available memory. Inactive programs can be written to disk and the space they occupied allocated to other programs. If a disk-resident program is needed later, it can be read back into any available memory location and restarted. Similarly, if a program needs a large contiguous block of storage, and the total amount is available only in nonadjacent fragments,

---

[1] *8086 programs can instruct the BIU to access a variable in one of the other currently addressable segments. For this problem, we will not worry about this to have simpler instruction.*

other program segments can be compacted to free up a continuous space. This process is shown graphically in Figure 2-19.

In order to be dynamically relocatable, a program must not load or alter its segment registers and must not transfer directly to a location outside the current code segment. In other words, all offsets in the program must be relative to fixed values contained in the segment registers. This allows the program to be moved anywhere in memory as long as the segment registers are updated to point to the new base addresses.
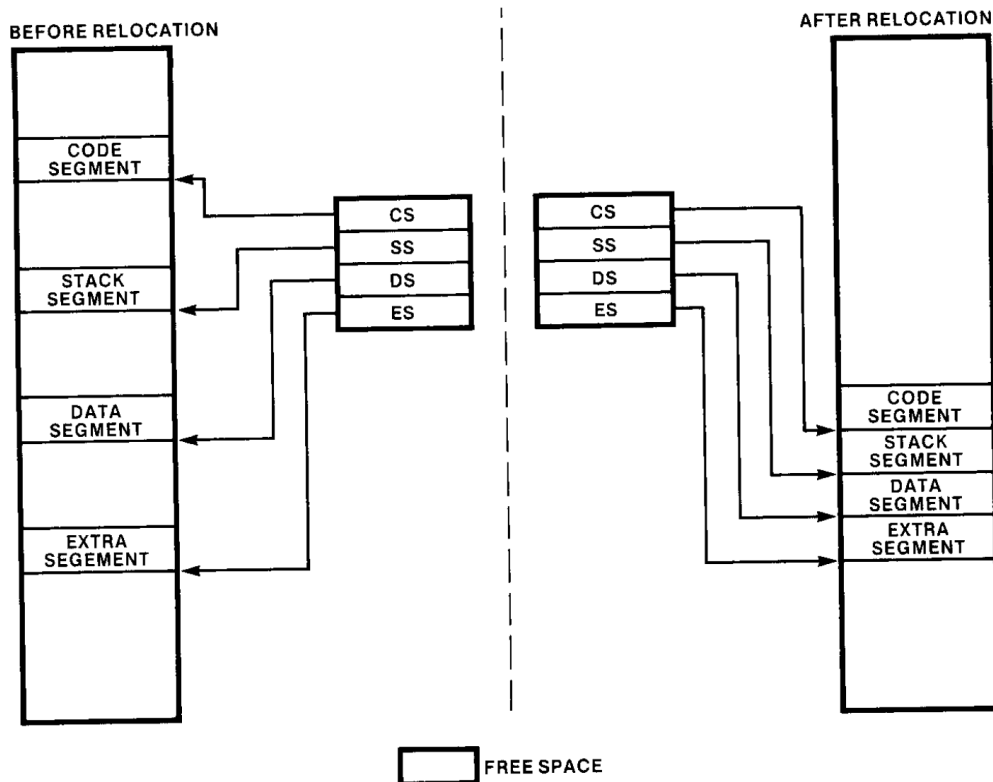


Figure 2-19. Dynamic Code Relocation

## Stack implementation

Stacks in the 8086 and 8088 are implemented in memory and are located by the stack segment register (SS) and the stack pointer register (SP). A system may have an unlimited number of stacks, and a stack may be up to 64k bytes long, the maximum length of a segment. (An attempt to expand a stack beyond 64k bytes overwrites the beginning of the stack.) One stack is directly addressable at a time; this is the current stack, often referred to simply as "the" stack. SS contains the base address of the current stack and SP points to the top of the stack (TOS). In other words, SP contains the offset of the top of the stack from the stack segment's base address. Note, however, that the stack's base address (contained in SS) is not the "bottom" of the stack.

8086 and 8088 stacks are 16 bits wide; instructions that operate on a stack add and remove stack items one word at a time. An item is pushed onto the stack (see figure 2-20) by decrementing SP by 2 and writing the item at the new TOS. An item is popped off the stack by copying it from TOS and then incrementing SP by 2. In other words, the stack grows down in memory toward its base address. Stack operations never move items on the stack, nor do they erase them. The top of the stack changes only as a result of updating the stack pointer.
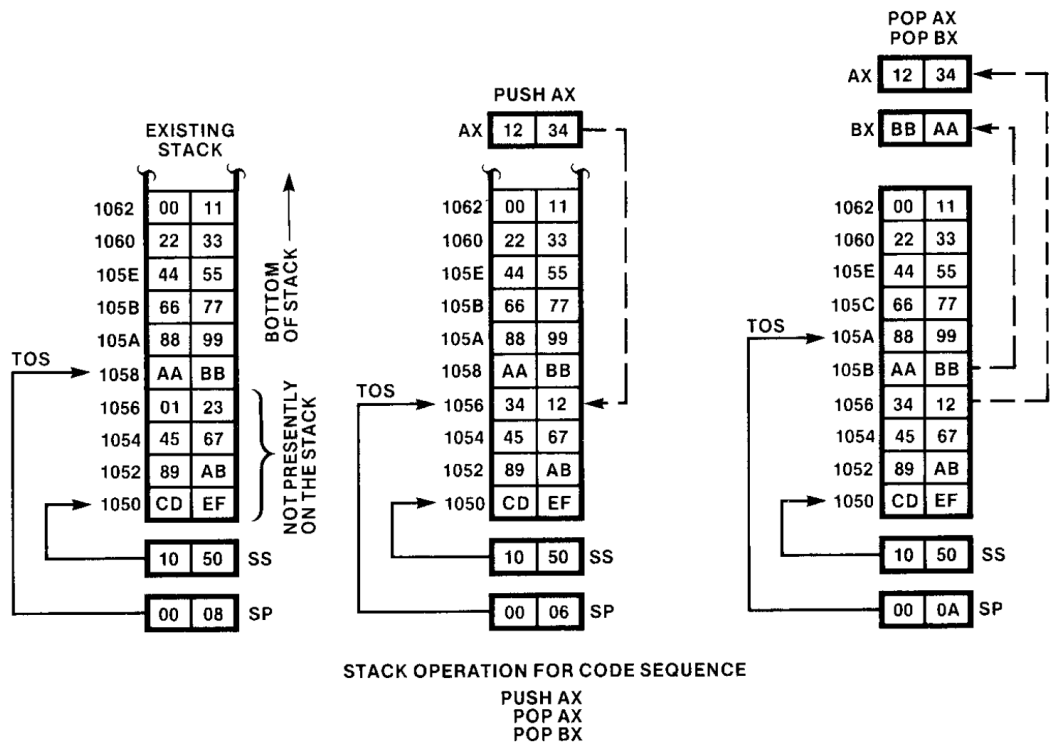


Figure 2-20. Stack Operation



U = UNDEFINED; VALUE IS INDETERMINATE
O = OVERFLOW FLAG
D = DIRECTION FLAG
I = INTERRUPT ENABLE FLAG
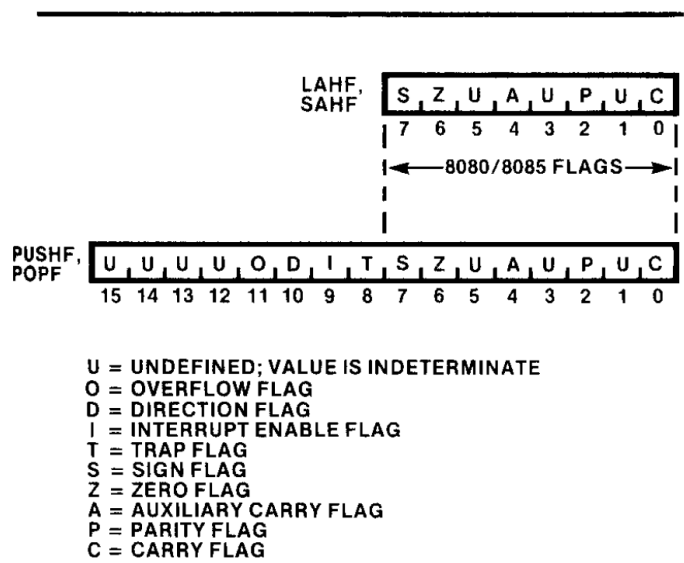T = TRAP FLAG
S = SIGN FLAG
Z = ZERO FLAG
A = AUXILIARY CARRY FLAG
P = PARITY FLAG
C = CARRY FLAG

Figure 2-32. Flag Storage Formats