

## CST 8152 – Compilers - Assignment #1

**Due Date:** prior or on **June 8<sup>th</sup>, 2020 (midnight)**

**Earnings:** 6% of your course grade (plus 1% bonus)

### **Purpose: Programming and Using Dynamic Structures (buffers) with C**

This is a review of and an exercise in C coding style, programming techniques, data types and structures, memory management, and simple file input/output. It will give you a better understanding of the type of internal data structures used by a simple compiler you will be building this semester. This assignment will be also an exercise in “*excessively defensive programming*”. You are to write functions that should be “overly” protected and should not abruptly terminate or “**crash**” at run-time due to invalid function parameters, erroneous internal calculations, or memory violations. To complete the assignment, you should fulfill the following two tasks:

### **Task 1: The Buffer Data Structure and Utility Functions**

Buffers are often used when developing compilers because of their efficiency (see page 111 of your textbook). You are to implement a buffer that can operate in three different modes: a “fixed-size” buffer, an “additive self-incrementing” buffer, and a “multiplicative self-incrementing” buffer. The buffer implementation is based on two associated data structures: a **Buffer Descriptor** (or *Buffer Handle*) and an array of characters (the actual character buffer). Both structures are to be created “on demand” at run time, that is, they are to be allocated dynamically. The **Buffer Descriptor** or **Buffer Handle** - the names suggest the purpose of this buffer control data structure - contains all the necessary information about the array of characters: a pointer to the beginning of the character array location in memory, the current capacity, the next character entry position, the increment factor, the operational mode and some additional parameters.

In this assignment you are to complete the coding for a “*buffer utility*”, which includes the buffer data structure and the associated functions, following strictly the given specifications. Use the data declarations and function prototypes given below. **Do not change the names and the data types of the functions and the variables. Any change will be regarded as a serious specification violation.** Write the associated code.

The following structure declaration must be used to implement the Buffer Descriptor:

```
typedef struct BufferDescriptor {
    char *cb_head;      /* pointer to the beginning of character array (character buffer) */
    short capacity;     /* current dynamic memory size (in bytes) allocated to character buffer */
    short addc_offset;  /* the offset (in chars) to the add-character location */
    short getc_offset;  /* the offset (in chars) to the get-character location */
    short markc_offset; /* the offset (in chars) to the mark location */
    char inc_factor;    /* character array increment factor*/
    char mode;          /* operational mode indicator*/
    unsigned short flags; /* contains character array reallocation flag and end-of-buffer flag */
} Buffer, *pBuffer;
```

Where:

**capacity** is the current total size (measured in bytes) of the memory allocated for the character array by **malloc()/realloc()** functions. In the text below it is referred also as current capacity. It is whatever value you have used in the call to **malloc()/realloc()** that allocates the storage pointed to by **cb\_head**.

**inc\_factor** is a buffer increment factor. It is used in the calculations of a new buffer **capacity** when the buffer needs to grow. The buffer needs to grow when it is full but still another character needs to be added to the buffer. The buffer is full when **addc\_offset** measured in bytes is equal to **capacity** and thus all the allocated memory has been used. The **inc\_factor** is only used when the buffer operates in one of the “self-incrementing” modes. In “additive self-incrementing” mode it is a positive integer number in the range of 1 to 255 and represents directly the increment (measured in characters) that must be added to the current capacity every time the buffer needs to grow. In “multiplicative self-incrementing” mode it is a positive integer number in the range of 1 to 100 and represents a percentage used to calculate the new capacity increment that must be added to the current capacity every time the buffer needs to grow.

**addc\_offset** is the distance (measured in chars) from the beginning of the character array (**cb\_head**) to the location where the next character is to be added to the existing buffer content. **addc\_offset** (measured in bytes) must never be larger than **capacity**, or else you are overrunning the buffer in memory and your program may crash at run-time or destroy data.

**getc\_offset** is the distance (measured in chars) from the beginning of the character array (**cb\_head**) to the location of the character which will be returned if the function **b\_getc()** is called. The value **getc\_offset** (measured in chars) must never be larger than **addc\_offset**, or else you are overrunning the buffer in memory and your program may get wrong data or crash at run-time. If the value of **getc\_offset** is equal to the value of **addc\_offset**, the buffer has reached the end of its current content.

**markc\_offset** is the distance (measured in chars) from the beginning of the character array (**cb\_head**) to the location of a *mark*. A *mark* is a location in the buffer, which indicates the position of a specific character (for example, the beginning of a word or a phrase).

**mode** is an operational mode indicator. It can be set to three different integer numbers: 1, 0, and -1. The number 0 indicates that the buffer operates in “fixed-size” mode; 1 indicates “additive self-incrementing” mode; and -1 indicates “multiplicative self-incrementing” mode. The mode is set when a new buffer is created and cannot be changed later.

**flags** is a field containing different flags and indicators. In cases when storage space must be as small as possible, the common approach is to pack several data items into single variable; one common use is a set of single-bit or multiple-bit flags or indicators in applications like compiler buffers, file buffers, and database fields. The flags are usually manipulated through different bitwise operations using a set of “masks.” Alternative technique is to use *bit-fields*. Using bit-fields allows individual fields to be manipulated in the same way as structure members are manipulated. Since almost everything about bit-fields is implementation dependent, this approach should be avoided if the portability is a concern. In this implementation, you are to use bitwise operations and masks (see *bitmask.c* example).

Each flag or indicator uses one or more bits of the **flags** field. The flags usually indicate that something happened during a routine operation (end of file, end of buffer, integer arithmetic sign overflow, and so on). Multiple-bit indicators can indicate, for example, the mode of the buffer (three different combinations – therefore 2-bits are needed). In this implementation the **flags** field has the following structure:

Bit	MSB 15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0 LSB
Content	1	1	1	1	1	1	1	1	1	1	1	1	1	x	x	1
Description	reserved for future use must be set to 1 and stay 1 all the time in this implementation													r_flag	eob flag	reserved

The LSB 1 of the **flags** field is a single-bit *end-of-buffer* (**eob**) flag. The **eob** bit is by default **0**, and when set to 1, it indicates that the end of the buffer content has been reached during the buffer read operation (**b\_getc()** function). If **eob** is set to **1**, the function **b\_getc()** should not be called before the **getc\_offset** is reset by another operation.

The LSB 2 of the **flags** field is a single-bit reallocation flag (**r\_flag**). The **r\_flag** bit is by default **0**, and when set to 1, it indicates that the location of the buffer character array in memory has been changed due to memory reallocation. This could happen when the buffer needs to expand or shrink. The flag can be used to avoid dangling pointers when pointers instead of offsets are used to access the information in the character buffer.

The rest of the bits are reserved for further use and must be set by default to **1**. When setting or resetting **eob** or **r\_flag** they must not be changed by the bitwise operation manipulating bit 2 and 1.

You are to implement the following set of buffer utility functions (operations). All function definitions must be stored in a file named **buffer.c**. Later they will be used by all other parts of the compiler when a temporary storage space is needed.

The **first implementation step** in all functions must be the validation (if possible and appropriate) of the function arguments. If an argument value is invalid, the function must return immediately an appropriate failure indicator.

### **Buffer \* b\_allocate (short init\_capacity, char inc\_factor, char o\_mode)**

This function creates a new buffer in memory (on the program heap). The function

- tries to allocate memory for one **Buffer** structure using **calloc()**;
- tries to allocates memory for one dynamic character buffer (character array) calling **malloc()** with the given initial capacity **init\_capacity**. The range of the parameter **init\_capacity** must be between **0** and the **MAXIMUM ALLOWED POSITIVE VALUE – 1** inclusive. The *maximum allowed positive value* is determined by the data type of the parameter **init\_capacity**. If the **init\_capacity** is **0**, the function tries to create a character buffer with default size **200** characters. If the **init\_capacity** is **0**, the function ignores the current value of the parameter **inc\_factor** and sets the buffer structure **inc\_factor** to **15** in mode **a** and **m** or to **0** in mode **f**. The pointer returned by **malloc()** is assigned to **cb\_head**;
- sets the Buffer structure operational mode indicator **mode** and the **inc\_factor**. If the **o\_mode** is the symbol **f**, the **mode** and the buffer **inc\_factor** are set to number **0**. If the **inc\_factor** parameter is **0** and **init\_capacity** is not **0** (see above), the **mode** and the buffer **inc\_factor** are set to **0**. If the **o\_mode** is **a** and **inc\_factor** is in the range of **1** to **255** inclusive, the **mode** is set to number **1** and the buffer **inc\_factor** is set to the value of **inc\_factor**. If the **o\_mode** is **m** and **inc\_factor** is in the range of **1** to **100** inclusive, the **mode** is set to number **-1** and the **inc\_factor** value is assigned to the buffer **inc\_factor**;

- copies the given *init\_capacity* value into the Buffer structure **capacity** variable;
- sets the **flags** field to its default value which is FFF9 hexadecimal.

Finally, on success, the function returns a pointer to the **Buffer** structure. It must return **NULL** pointer on any error which violates the constraints imposed upon the buffer parameters or prevents the creation of a working buffer. If run-time error occurs, the function must return immediately after the error is discovered. Check for all possible errors which can occur at run time. Do not allow “memory leaks”, “dangling” pointers, or “bad” parameters.

#### ***pBuffer b\_addc (pBuffer const pBD, char symbol)***

Using a bitwise operation the function resets the **flags** field *r\_flag* bit to 0 and tries to add the character **symbol** to the character array of the given **buffer** pointed by **pBD**. If the buffer is operational and it is not full, the symbol can be stored in the character buffer. In this case, the function adds the character to the content of the character buffer, increments **addc\_offset** by 1 and returns.

If the character buffer is already full, the function will try to resize the buffer by increasing the current capacity to a new capacity. How the capacity is increased depends on the current operational mode of the buffer.

If the operational mode is **0**, the function returns **NULL**.

If the operational mode is **1**, it tries to increase the current capacity of the buffer to a *new capacity* by adding **inc\_factor** (converted to bytes) to **capacity**. If the result from the operation is positive and does not exceed the **MAXIMUM ALLOWED POSITIVE VALUE – 1** (minus 1), the function proceeds. If the result from the operation is positive but exceeds the **MAXIMUM ALLOWED POSITIVE VALUE – 1** (minus 1), it assigns the **MAXIMUM ALLOWED POSITIVE VALUE – 1** to the *new capacity* and proceeds. The **MAXIMUM ALLOWED POSITIVE VALUE** is determined by the data type of the variable, which contains the buffer capacity.

If the result from the operation is negative, it returns **NULL**.

If the operational mode is **-1** it tries to increase the current capacity of the buffer to a *new capacity* in the following manner:

- If the current capacity can not be incremented anymore because it has already reached the maximum capacity of the buffer, the function returns **NULL**.

The function tries to increase the current capacity using the following formulae:

$$\begin{aligned} \text{available space} &= \text{maximum buffer capacity} - \text{current capacity} \\ \text{new increment} &= \text{available space} * \text{inc\_factor} / 100 \\ \text{new capacity} &= \text{current capacity} + \text{new increment} \end{aligned}$$

The *maximum buffer capacity* is the **MAXIMUM ALLOWED POSITIVE VALUE – 1**. If the *new capacity* has been incremented successfully, no further adjustment of the *new capacity* is required.

If as a result of the calculations, the *current capacity* cannot be incremented, but the *current capacity* is still smaller than the **MAXIMUM ALLOWED POSITIVE VALUE – 1**, then the **MAXIMUM ALLOWED POSITIVE VALUE – 1** is assigned to the *new capacity* and the function proceeds.

If the capacity increment in mode **1** or **-1** is successful, the function performs the following operations:

- the function tries to expand the character buffer calling **realloc()** with the *new capacity*. If the reallocation fails, the function returns NULL;
- if the location in memory of the character buffer has been changed by the reallocation, the function sets **r\_flag** bit to **1** using a bitwise operation;
- adds (appends) the character **symbol** to the buffer content;
- changes the value of **addc\_offset** by 1, and saves the newly calculated capacity value into **capacity** variable;
- the function returns a pointer to the **Buffer** structure.

The function must return NULL on any error. Some of the possible errors are indicated above but you must check for all possible errors that can occur at run-time. Do not allow “**memory leaks**”. Avoid creating “**dangling pointers**” and using “**bad**” parameters. The function **must not destroy** the buffer or the contents of the buffer even when an error occurs – it must simply return NULL leaving the existing buffer content intact. **A change in the project platform (16-bit, 32-bit or 64-bit) must not lead to improper behavior.**

#### ***int b\_clear (Buffer \* const pBD)***

The function retains the memory space currently allocated to the buffer, but re-initializes all appropriate data members of the given **Buffer** structure (buffer descriptor) so that the buffer will appear as just created to the client functions (for example, next call to **b\_addc()** will put the character at the beginning of the character buffer). The function does not need to clear the existing contents of the character buffer. If a run-time error is possible, the function should return **-1** in order to notify the calling function about the failure.

#### ***void b\_free (Buffer \* const pBD)***

The function de-allocates (frees) the memory occupied by the character buffer and the **Buffer** structure (buffer descriptor). The function should not cause abnormal behavior (crash).

#### ***int b\_isfull (Buffer \* const pBD)***

The function returns **1** if the character buffer is full; it returns **0** otherwise. If a run-time error is possible, the function should return **-1**.

#### ***short b\_addcoffset (Buffer \* const pBD)***

The function returns the current **addc\_offset**. If a run-time error is possible, the function should return **-1**.

#### ***short b\_capacity (Buffer \* const pBD)***

The function returns the current capacity of the character buffer. If a run-time error is possible, the function should return **-1**.

#### ***short b\_markc (pBuffer const pBD, short mark)***

The function sets **markc\_offset** to **mark**. The parameter **mark** must be within the current limit of the buffer (0 to **addc\_offset** inclusive). The function returns the currently set **markc\_offset**. If a run-time error is possible, the function should return **-1**.

### ***int b\_mode (Buffer \* const pBD)***

The function returns the value of **mode** to the calling function. If a run-time error is possible, the function should notify the calling function about the failure .

### ***size\_t b\_incfactor (Buffer \* const pBD)***

The function returns the non-negative value of **inc\_factor** to the calling function. If a run-time error is possible, the function should return **0x100**.

### ***int b\_load (FILE \* const fi, Buffer \* const pBD)***

The function loads (reads) an open input file specified by **fi** into a buffer specified by **pBD**. The function must use the standard function **fgetc(fi)** to read one character at a time and the function **b\_addc()** to add the character to the buffer. If the current character cannot be added to the buffer (**b\_addc()** returns NULL), the function returns the character to the file stream (file buffer) using **ungetc()** library function and then returns **-2** (use the defined **LOAD\_FAIL** constant). The operation is repeated until the standard macro **feof(fi)** detects end-of-file on the input file. The end-of-file character must not be added to the content of the buffer.

Only the standard macro **feof(fi)** must be used to detect end-of-file on the input file. Using other means to detect end-of-file on the input file will be considered a significant specification violation. If some other run-time errors are possible, the function should return **-1**. If the loading operation is successful, the function must return the number of characters added to the buffer.

### ***int b\_isempty (Buffer \* const pBD)***

If the **addc\_offset** is 0, the function returns 1; otherwise it returns 0. If a run-time error is possible, it should return **-1**.

### ***char b\_getc (Buffer \* const pBD)***

This function is used to read the buffer. The function performs the following steps:

- checks the argument for validity (possible run-time error). If it is not valid, it returns **-2**;
- if **getc\_offset** and **addc\_offset** are equal, using a bitwise operation it sets the **flags** field **eob** bit to **1** and returns number **0**; otherwise, using a bitwise operation it sets **eob** to **0**;
- returns the character located at **getc\_offset**. Before returning it increments **getc\_offset** by **1**.

### ***int b\_eob (Buffer \* const pBD)***

The function returns the **value** of the **flags** field determined only by the **eob** bit.

A bitwise operation must be used to return the value of the **flags** field. If a run-time error is possible, it should return **-1**.

### ***int b\_print (Buffer \* const pBD, char nl)***

This function is intended to be used for diagnostic purposes only. Using the **printf()** library function the function prints character by character the contents of the character buffer to the standard output (stdout). In a loop the function prints the content of the buffer calling **b\_getc()** and using **b\_eob()** to detect the end of the buffer content (using other means to detect the end of buffer content will be considered a significant specification violation). After the loop ends, it checks the **nl** and if it is not **0**, it prints a new line character. Finally, it returns the number of characters printed. The function returns **-1** on failure.

### ***Buffer \* b\_compact(Buffer \* const pBD, char symbol)***

For all operational modes of the buffer the function shrinks (or in some cases may expand) the buffer to a *new capacity*. The *new capacity* is the current limit plus a space for one more character. In other words the *new capacity* is ***addc\_offset + 1*** converted to bytes. The function uses ***realloc()*** to adjust the *new capacity*, and then updates all the necessary members of the buffer descriptor structure. Before returning a pointer to ***Buffer***, the function adds the ***symbol*** to the end of the character buffer (**do not** use ***b\_addc()***, use ***addc\_offset*** to add the symbol) and increments ***addc\_offset***. The function must return NULL if for some reason it cannot to perform the required operation. It must set the ***r\_flag*** bit appropriately.

### ***char b\_rflag (Buffer \* const pBD)***

The function returns the **value** of the ***flags*** field determined only by the ***r\_flag*** bit. A bitwise operation must be used to return the value of the ***flags*** field. If a run-time error is possible, it should return **-1**.

### ***short b\_retract (Buffer \* const pBD)***

The function decrements ***getc\_offset*** by 1. If a run-time error is possible, it should return **-1**; otherwise it returns ***getc\_offset***.

### ***short b\_reset (Buffer \* const pBD)***

The function sets ***getc\_offset*** to the value of the current ***markc\_offset***. If a run-time error is possible, it should return **-1**; otherwise it returns ***getc\_offset***.

### ***short b\_getcoffset (Buffer \* const pBD)***

The function returns ***getc\_offset*** to the calling function. If a run-time error is possible, it should return **-1**.

### ***int b\_rewind(Buffer \* const pBD)***

The function set the ***getc\_offset*** and ***markc\_offset*** to 0, so that the buffer can be reread again. If a run-time error is possible, it should return **-1**; otherwise it returns 0;

### ***char \* b\_location(Buffer \* const pBD, short loc\_offset)***

The function returns a pointer to the location of the character buffer indicated by ***loc\_offset***. ***loc\_offset*** is the distance (measured in chars) from the beginning of the character array (***cb\_head***). If a run-time error is possible, it should return **NULL**.

All constant definitions, data type and function declarations (prototypes) must be located in a header file named ***buffer.h***. You are allowed to use only named constants in your programs (except when incrementing something by 1 or setting a numeric value to 0). To name a constant you must use ***#define*** preprocessor directive (see ***buffer.h***). The incomplete ***buffer.h*** is posted on Brightspace (BS). All function definitions must be stored in a file named ***buffer.c***.

## **Task 2: Testing the Buffer**

To test your program you are to use the test harness program ***platy\_bt.c*** (do not modify it) and the input files ***ass1e.pls*** (an empty file), and ***ass1.pls***. The corresponding output files are ***ass1e.out*** and ***ass1ai.out*** (mode = 1), ***ass1mi.out*** (mode = -1), ***ass1fi.out*** (mode = 0). Those files are generated by my test program and contain plain ASCII text. They are available as part of the assignment postings on Brightspace (BS). You must create a standard console project named ***buffer*** with an executable target ***buffer*** (see ***Creating\_C\_Project*** document in Lab0). The project must contain only one header file (***buffer.h***) and two source files: ***buffer.c*** and ***platy\_bt.c***.

Here is a brief description of the program that is provided for you on Brightspace (BS). It simulates “normal” operating conditions for your buffer utility. The program (***platy\_bt.c***) main function takes to parameters from the command line: an input file name and a character (**f** - fixed-size, **a** – additive self-increment, or **m** – multiplicative self increment) specifying the buffer operational mode. It opens up a file with the specified name (for example, ***ass1.pls***), creates a buffer, and loads it with data from the file using the ***b\_load()*** function. Then the program prints the current capacity, the current size (limit), the current operational mode, the increment factor, the current mark, the flags, and the contents of the buffer. Finally, it compacts the buffer, and if the operation is successful, it prints the buffer contents again. Your program must not overflow any buffers in any operational mode, no matter how long the input file is. The provided main program will not test all your functions. You are strongly encouraged to test all your buffer functions with your own test files and modified main function.

### **Bonus Task: Implementing a Preprocessor Macro Definition and Expansion (1%)**

Implement ***b\_isfull()*** both as a function and a macro expansion (macro). Using conditional processing you must allow the user to choose between using the macro or the function in the compiled code. If **B\_ISFULL** name is defined the macro should be used in the compiled code. If the **B\_ISFULL** name is not defined or undefined, the function should be used in the compiled code. To receive credit for the bonus task your code must be well documented, tested, and working.

### ***SUBMIT THE FOLLOWING:***

**Digital Submission: Compress** into a **zip** file the following files: *platy\_bt.c*, *buffer.h*, *buffer.c*, *ass1.pls*, *ass1e.pls* and the corresponding *test output files* produced by your program. Include your additional input/output test files if you have any. Also include a Cover page and a Test Plan. The submission must follow the course submission standards. You will find the Assignment Submission Standard as well as the Assignment Marking Guide (***CST8152\_ASSAMG.pdf***) for the Compilers course on the Brightspace.

Upload the zip file on Brightspace. The file must be submitted prior or on the due date as indicated in the assignment. The name of the file must be **Your Last Name** followed by the last three digits of your student number followed by your lab section number. For example: Ranev345\_s11.zip.

### **Assignments will not be marked if there are not source files in the digital submission.**

Assignments could be late, but the lateness will affect negatively your mark: see the Course Outline and the Marking Guide. All assignments must be successfully completed to receive credit for the course, even if the assignments are late.



**Evaluation Note:** Make your functions as efficient as possible. These functions are called many times during the compilation process. The functions will be graded with respect to design, documentation, error checking, robustness, and efficiency. When evaluating and marking your assignment, I will use the standard project and **platy\_bt.c** and the test files posted on Brightspace. If your program compiles, runs, and produces correct output files, it will be considered a *working program*. Additionally, I will try my best to “crash” your functions using a modified main program, which will test all your functions including calling them with “invalid” parameters. I will use also some additional test files (for example, a large file). This can lead to fairly big reduction of your assignment mark (see **CST8152\_ASSAMG** and **cMarkingSheetA1** documents).

Enjoy the assignment. And do not forget that:

***“Writing a program is like painting. It is better to start on a new canvas.”*** Ancient P-Artist  
***“It is part of the nature of humans to begin with romance (buffer) and build to reality (compiler).”*** by Ray Bradbury

**#define buff·er** (bŭf·ər) *noun* (Microsoft Bookshelf)

1. Something that lessens or absorbs the shock of an impact.
2. One that protects by intercepting or moderating adverse pressures or influences: “A sense of humor . . . may have served as a buffer against the . . . shocks of disappointment” (James Russell Lowell).
3. Something that separates potentially antagonistic entities, as an area between two rival powers that serves to lessen the danger of conflict.
4. *Chemistry.* A substance that minimizes change in the acidity of a solution when an acid or base is added to the solution.
5. **Computer Science.** A device or memory area used to store data temporarily and deliver it at a rate different from that at which it was received.

S^R, CST8152 – Compilers, May 1<sup>st</sup>, 2020