# ROMhacking Int. 104

# Compression Scheme RE

*(Dynamic and Static approach)*

By Bunkai
(Last update: August 1st 2024)

# Index

# Introduction & Tools

*\* Warning: the Int. means Intermediate, this is not a beginner topic. However, the intention of the document is to be a reference for beginners, and a reminder for veterans.*

*\* RE stands for Reverse Engineering, so expect a lot of that.*

*For this document/tutorial, you will need to have the following tools:*

*Game's ROM:  Marmalade Boy (Japan).sgb with CRC 32: 0F3FF7DA*

*The Hex Editor of your choice. The screenshots are from crystal tile, but ImHex and HxD work too*

*A console emulator. Screenshots for this document are from BGB and Mesen2.*

*Your code editor of choice if you want to program your tools. Tools in this doc are code in python.*

*A lot of patience, effort, and a place to write your notes to review them when needed.*

# Preface

*This document attempts to be a learning reference to 'attack' compression in a game, from start to finish, with all the (minimum required) theory followed by practical examples. All of the steps and explanations are applied to an actual project with actual screenshots of the project.*

*References will be given when required, you can use them to dive deeper. But rest assured, you will only need the tools used in each step to follow along. You can just read this paper without any extra material if your only goal is to learn about the process.*

*You must know that each section of this document can be read standalone. If you want to read only one section, don't be afraid to do so, however,  you will benefit from reading all.*

*\* Note: The actual extraction/reinsertion from/to the ROM won't be detailed here. If you need a hint for that, just mind the offset to extract or reinsert bytes and copy them in or out.*

*The end goal is to understand and be able to use and export the knowledge to do this process, not to parrot. Remember: Curiosity leads to knowledge, be curious*

# Static Analysis

I am going to follow, and copy, the great explanations from abridgewater in his: *Rom Hacking 201: A non-ASM approach to compressed data*, which you can find here:
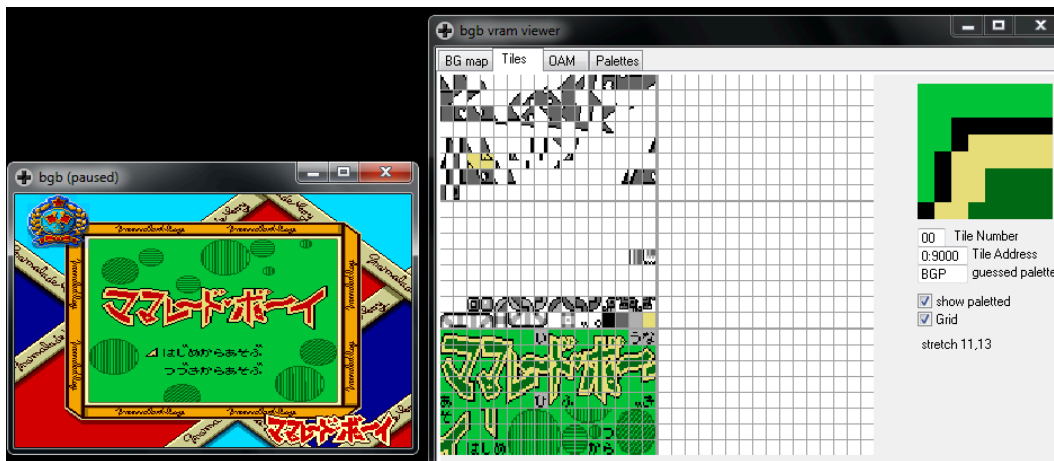https://www.romhacking.net/documents/878/
However, I will be applying them to the "marmalade boy (japan).sgb" game in order to learn and "break" his compression routine. (I will also add some aid from the use of a, dynamic, debugger in some parts for explanation purposes)
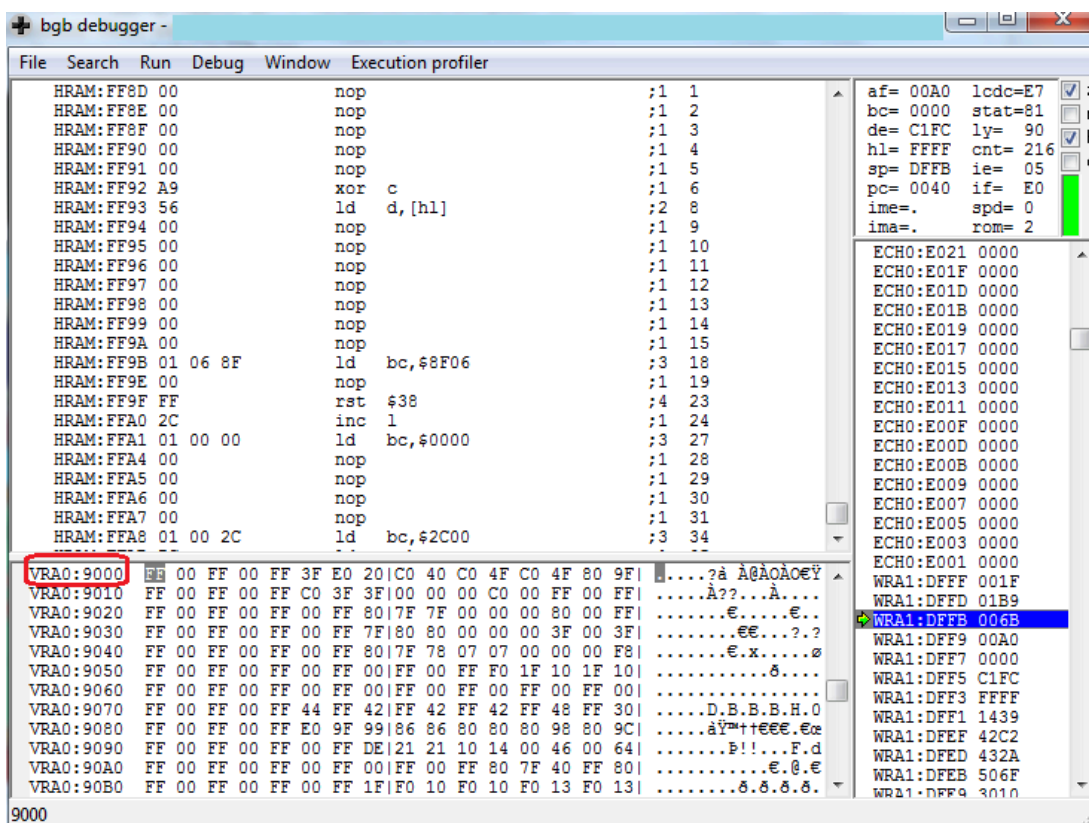
* Investigating the Title Screen

Step one, find our decompressed data. Start the ROM in an emulator that has debugging support, get to the title screen, and pause the emulation. Next, use the tile viewer to find the address of some of the decompressed tiles in VRAM.

Here you can see an example of one Tile and its address 0:9000



Now, point your memory viewer at that (0:9000) address of the V-RAM to see the decompressed data in arrays of bytes. See below for a reference:

That very same RAM block is copied below for quick reference:

```
        00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F
-------------------------------------------------------------------------------------

VRA0:9000: FF 00 FF 00 FF 3F E0 20 C0 40 C0 4F C0 4F 80 9F
VRA0:9010: FF 00 FF 00 FF C0 3F 3F 00 00 00 C0 00 FF 00 FF
VRA0:9020: FF 00 FF 00 FF 00 FF 80 7F 7F 00 00 00 80 00 FF
VRA0:9030: FF 00 FF 00 FF 00 FF 7F 80 80 00 00 00 3F 00 3F
VRA0:9040: FF 00 FF 00 FF 00 FF 80 7F 78 07 07 00 00 00 F8
VRA0:9050: FF 00 FF 00 FF 00 FF 00 FF 00 FF F0 1F 10 1F 10
VRA0:9060: FF 00 FF 00 FF 00 FF 00 FF 00 FF 00 FF 00 FF 00
VRA0:9070: FF 00 FF 00 FF 44 FF 42 FF 42 FF 42 FF 48 FF 30
VRA0:9080: FF 00 FF 00 FF E0 9F 99 86 86 80 80 80 98 80 9C
VRA0:9090: FF 00 FF 00 FF 00 FF DE 21 21 10 14 00 46 00 64
VRA0:90A0: FF 00 FF 00 FF 00 FF 00 FF 00 FF 80 7F 40 FF 80
VRA0:90B0: FF 00 FF 00 FF 00 FF 1F F0 10 F0 10 F0
```

Pro Tip for this game: In mesen2's real time debugger, if you go to the VRAM $1400 address. You can do some screen edits like the one in the image below, with the あ tile edited.



Explanation of this: take the あ bytes from the screenshot ($1400: FF 00 FF 20 FF FC FF 20 FF 7C FF AA FF 92 FF 64) to https://www.huderlem.com/demos/gameboy2bpp.html



The changes done here are explained as follows, each 16 bytes (a full tile). If you separate that in blocks of 2:

"FF 64" each of these pairs is a row of the tile being the 2nd half "64" the tile palette and the 1st half "FF" the drawn part of the raw.

This also suggests that the font could be the same used in the script.
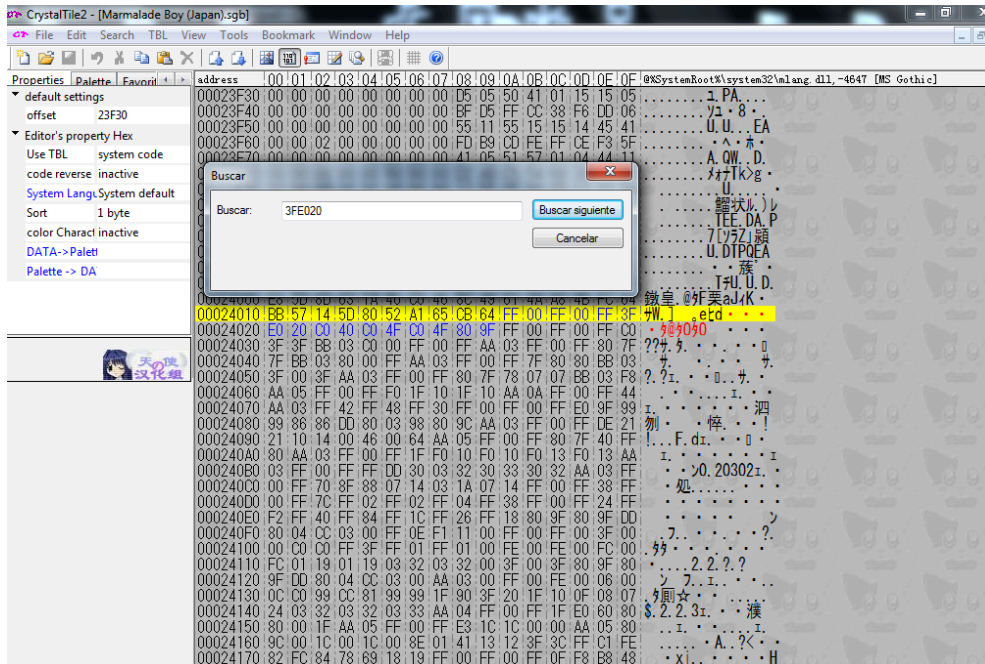(As proved later, this is not the case)

Back to our analisis now, this RAM block has a very regular pattern in the beginning, (it's very common and repeated):

9000: 'FF 00 FF 00 FF'
So, we'll skip ahead to a significant break in the pattern:  At $9005, "3F E0 20"
We want this pattern to make a relative search in the ROM so we can find the related data.

Tip: You will find false positives, just remember to keep searching until you find the correct match. Searched match displayed below:



Jackpot! We note that the pattern from VRAM $9005, "0e f1 3f c0" appears here at $2401A as "FF 00 FF 00 FF 3F E0 20 C0 40 C0 4F C0 4F 80 9F", and the remaining of the line matches too, further confirming that this is the data that we're looking for. (If the full block is shown, everything is in plain, non compressed. If it shows tidbits with other hexes in-between then you have probably found it but it is compressed or have some extra unrelated data for the game's code to work)

Besides, you can do a corruption test of that address (this means editing those hexes, and saving the file to run it in the emulator and see if the image displayed changes). Since the display changes, it proves to be the right location too.



Pantalla original          Pantalla corrupta

We'll now grab a bit of this ROM data for deeper analysis. Copied below for quick reference:

```
         00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F
       ------------------------------------------------------------------------
24010:                                       FF 00 FF 00 FF 3F
24020: E0 20 C0 40 C0 4F C0 4F 80 9F FF 00 FF 00 FF C0
24030: 3F 3F BB 03 C0 00 FF 00 FF AA 03 FF 00 FF 80 7F
24040: 7F BB 03 80 00 FF AA 03 FF 00 FF 7F 80 80 BB 03
24050: 3F 00 3F AA 03 FF 00 FF 80 7F 78 07 07 BB 03 F8
24060: AA 05 FF 00 FF F0 1F 10 1F 10 AA 0A FF 00 FF 44
24070: AA 03 FF 42 FF 48 FF 30 FF 00 FF 00 FF E0 9F 99
24080: 99 86 86 DD 80 03 98 80 9C AA 03 FF 00 FF DE 21
24090: 21 10 14 00 46 00 64 AA 05 FF 00 FF 80 7F 40 FF
240A0: 80 AA 03 FF 00 FF 1F F0 10 F0 10 F0 13 F0 13 AA
240B0: 03 FF 00 FF FF DD 30 03 32 30 33 30 32 AA 03 FF
```

At this point we have a chunk of decompressed data and a couple of places where the decompressed data matches up with data in the ROM, meaning that we've found approximately where the compressed data sits. We don't yet know any details of the compression scheme other than that it includes literal data at times, and we don't know where the compressed data starts or where it ends.

To find either end of the compressed data, we want to look forward or backwards through the ROM until we find the compressed data that corresponds with the start or the end of the decompressed data. This, in turn, requires that we work out at least the basics of the overall compression scheme.

And the compressed data includes literal data, so what we do is to line it up with the decompressed data and look for patterns. What's in the compressed data that's not in the decompressed data, and what's in the decompressed data that's not in the compressed data? Our starting point is the first match we found, and we work our way forward from there.

VRAM $9000:  FF 00 FF 00 FF 3F E0 20 C0 40 C0 4F C0 4F 80 9F
ROM  $2401A: FF 00 FF 00 FF 3F E0 20 C0 40 C0 4F C0 4F 80 9F
> These two are equal so we can't take any extra info from the comparison

VRAM $9010:  FF 00 FF 00 FF C0 3F 3F 00 00 00 C0 00 FF 00 FF
ROM  $2402A: FF 00 FF 00 FF C0 3F 3F BB 03 C0 00 FF 00 FF AA
> From these two, we can see a pattern where 'BB 03' means take the first Byte and repeat three times '00 00 00'. no control code at the end, just the next byte to display

VRAM $9020:  FF 00 FF 00 FF 00 FF 80 7F 7F 00 00 00 80 00 FF
ROM  $2403A: 03 FF 00 FF 80 7F 7F BB 03 80 00

> We can see a similar pattern again here (taking the 'AA' from the previous row), we can see 'AA 03 FF 00' results in 'FF 00 FF 00 FF 00' where the 03 again marks how many times we repeat the bytes, and we could guess that the AA means it's a two Bytes string this time 'FF 00'


VRAM $9030:  FF 00 FF 00 FF 00 FF 7F 80 80 00 00 00 3F 00 3F
ROM  $2403B: FF AA 03 FF 00 FF 7F 80 80 BB 03 3F 00 3F
> we can see the 'AA 03 FF 00' and the 'BB 03' patterns again here.

VRAM $9040:  FF 00 FF 00 FF 00 FF 80 7F 78 07 07 00 00 00 F8
ROM  $24049: AA 03 FF 00 FF 80 7F 78 07 07 BB 03 F8
> we can see the 'AA 03 FF 00' and the 'BB 03' patterns again here.

VRAM $9050: FF 00 FF 00 FF 00 FF 00 FF 00 FF F0 1F 10 1F 10
ROM $2404D: AA 05 FF 00 FF F0 1F 10 1F 10
> we can see the 'AA 03 FF 00' but now with more bytes repeated, which change it to 'AA 05 FF 00' further confirming our previous guesses

We can now code this pattern into a routine to test. Using the rom hex block as the input file, should give the vram hex block as the output

From here on, the process is to keep referencing the pattern while also looking for how does the decompressor know when to stop decompressing? There are two ways for the decompressor to know when to stop. Either it knows how long the decompressed data will be and stops when it fills that space, or there is a special end marker that it looks for in the compressed data.

The end marker is the more likely of the two, and the only space in the encoding scheme that it can fit is as a back reference with a particular value. This is an easy enough thing to determine: Walk forward through the compressed data until you reach the end. If it just stops without any unusual backreference ( An item in a regular expression equivalent to the text matched by an earlier pattern in the expression), it's length-limited, and the length needs to be found. If there /is/ an unusual backreference, then that's the end marker.

[To be brutally honest, I didn't find this one through static, I got "lucky" because I had the routine disassembled from an old attempt to find the text (that I had archived until I could use it). I just didn't notice this was for the images until recognizing the same patterns in the static analysis. However, I will copy the reference of the end of this ROM's block below for those who want to do it]

ROM  $246A0: B6 92 6D 24 DB 49 B6 92 6D 24 DB 48 B7 81 7E 25
ROM  $246B0: DA 49 B6 93 6C 23 DC 47 B8 8F 70 3F C0 FF 00 **EE**
ROM  $246C0: 80 80 DD C0 04 A0 E0 90 F0 88 F8 88 F8 84 FC 82

After all the chewed info, we should be able to recreate the decompression routine. To test your tool at the beginning, use both (ROM and V-RAM) blocks you've been analyzing in the previous steps. Extract the block from the rom and decompress running your code,

check if the result matches with the vram block to see the accuracy until you have your tool ironed.

The process to create the tools will be talked about later in a future section. But, this one about the static analysis can be considered as finished.

# Dynamic Analysis

I am going to follow, and copy, the great explanations from YasaSheep in her: *ROM Hacking 202 - Reverse engineering with the debugger*, which you can find here: https://www.romhacking.net/documents/877/

However, I will be applying them to the "marmalade boy (japan).sgb" game in order to learn and "break" his compression routine.

I want to add here my whole gratitude to Phonymike for all his help documenting the ASM routine of the Marmalade game. (many of those code comments are done by him)
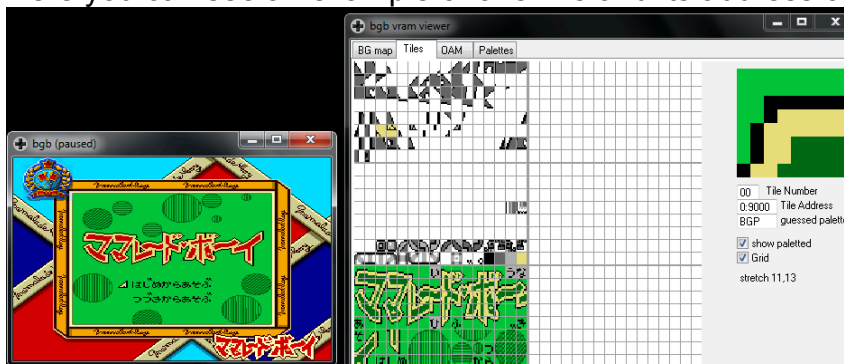
* Note: This 202 doc. uses the previous one as a base to find (have) the first addresses and a few routine logistics. However, in order to make both methods as standalone and complete as possible, I will reuse/copy the first page from the Static method (Because it's the same beginning in both cases, the difference starts when you branch out between analyzing the data in the hex values or the code in the debugger.)

* Extra note: This method was done at first following the code to find the text, we recognized the routine as a compression form but it obviously wasn't the text; As a result, the screenshots of this method weren't saved, only the ASM and comments were archived for future use. It hasn't been used again until I recognized the patterns in the static analysis (like I mentioned before). For those reasons, the process is mostly recreated and won't throw as many tips as one would write if it was done from scratch; but, I believe it will be "good enough" to explain the method. Besides, it gives a nice way to have a method comparison with the same data, and makes the "paper" most complete.
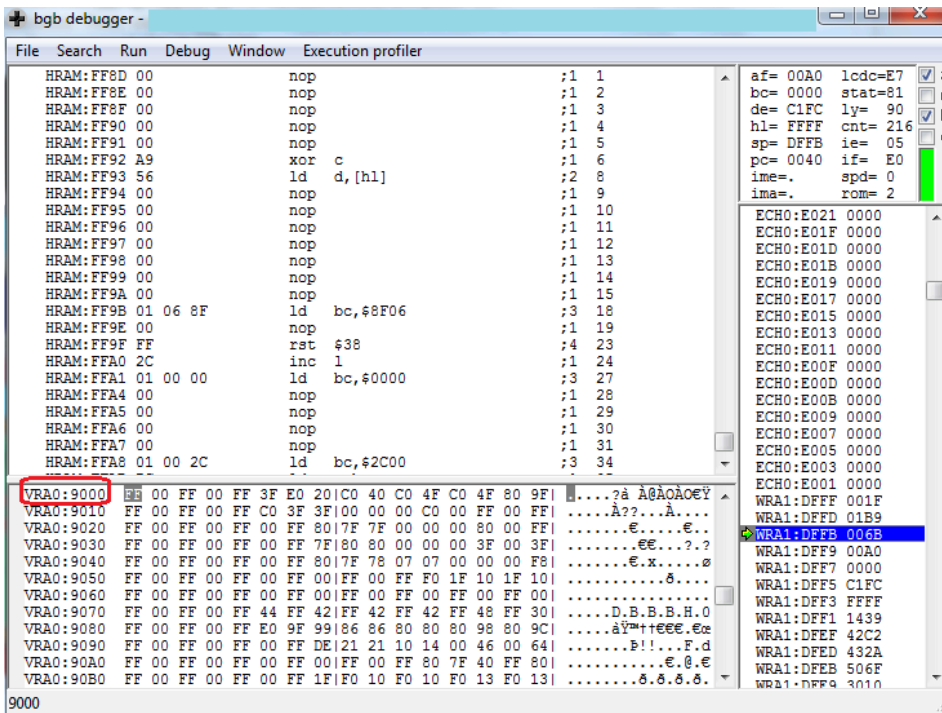
* Investigating the Title Screen

Step one, find our decompressed data. Start the ROM in an emulator that has debugging support, get to the title screen, and pause the emulation. Next, use the tile viewer to find the address of some of the decompressed tiles in VRAM.

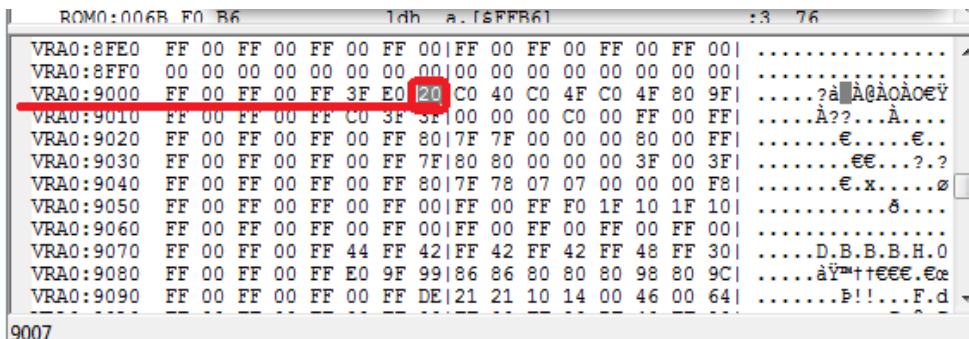Here you can see an example of one Tile and its address 0:9000

With that, point your memory viewer at that (0:9000) address of the V-RAM to see the decompressed data in arrays of bytes. See below for a reference:



You need to read those values until you find one that is unique enough so it stops when watched, during the execution running, but not so common to stop at every instruction.

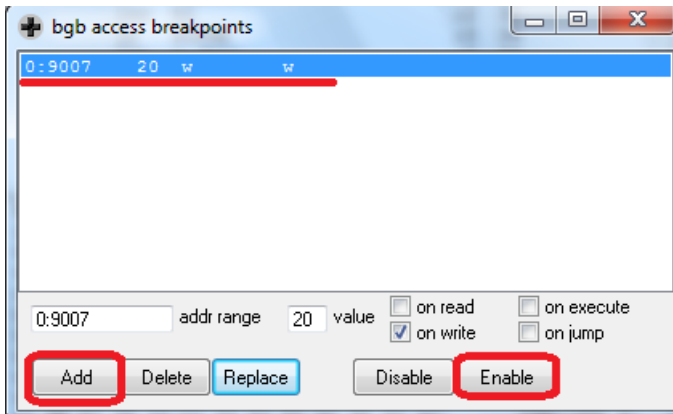I went with that '20' in the VRAM Address $9007.

VRAM0: 9000        FF 00 FF 00 FF 3F E0 20 C0 40 C0 4F C0 4F 80 9F



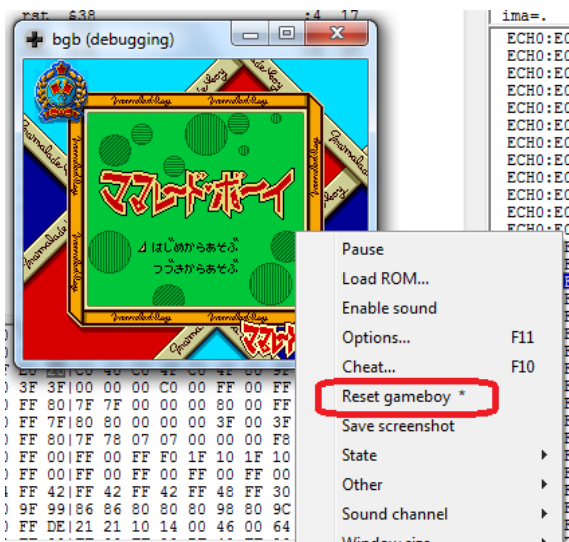Now make a breakpoint in the debugger, so it tells you when that happens.

-   right-click and select "Set access breakpoint"

You want to see when that address is written on with that value, so make sure "on write" is the only checkbox. For now. every other checkbox will make it stop for things you are not interested in, which will be useless effort.
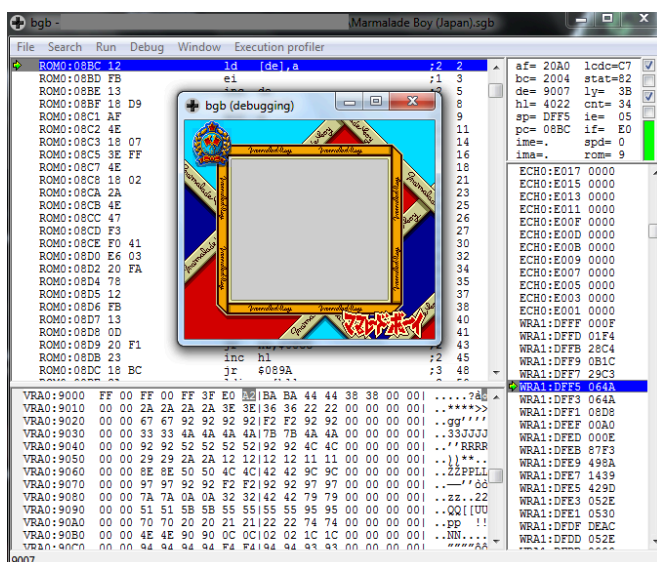


The debugger doesn't have a rewind, so you want to reset the emulation to make the game go back and redo the operations until it reaches the breakpoint.

Right-click the game screen and select "Reset gameboy



After that, the game will load and pause right before the value is updated to the one we are expecting. So, don't panic. You can see the screenshot of that moment below:
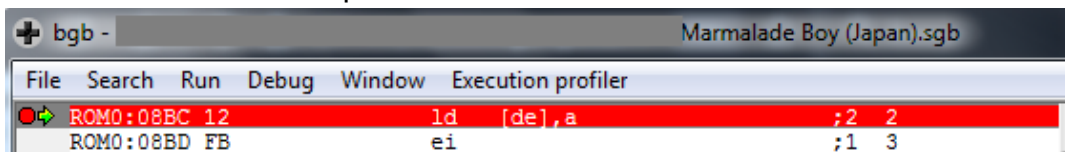


The values displayed in the emu are:

ROM0:08BC 12 ...

VRA0:9000
FF 00 FF 00 FF 3F E0 A2 BA BA 44 44 38 38 00

This is the operation that writes the value to the memory position we asked the emulator to watch.          **ROM0:08BC 12     ld [de],a**

This operation is "Load A into the address pointed to by DE". Essentially, as expected, it writes A (which contains 0x20).
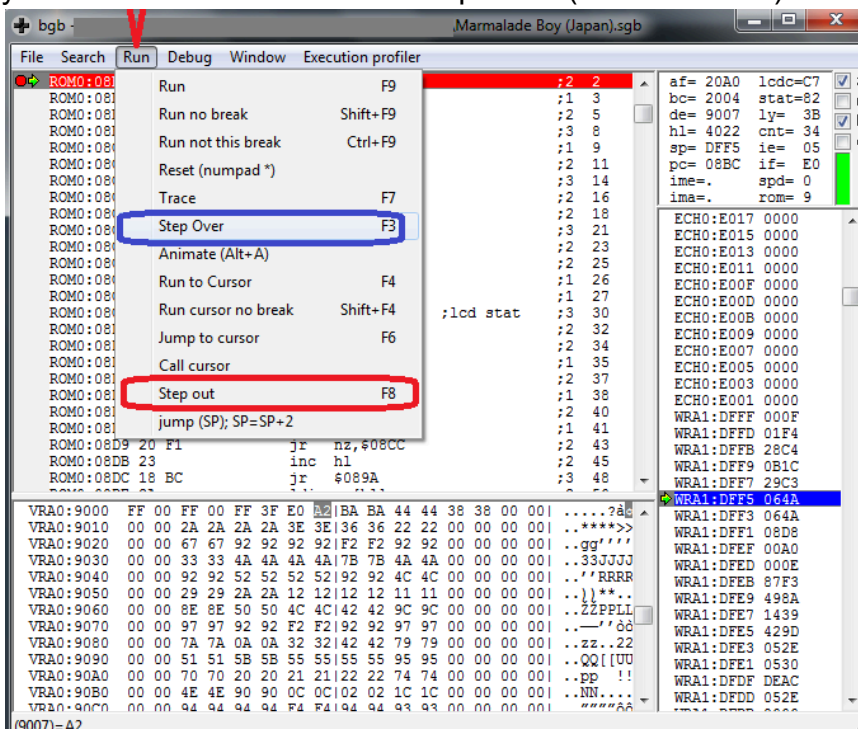
Tip: you can set a Breakpoint here then disable it from the "Debug → breakpoints" menu in order to save it as a bookmark). Just double click on it, the line will be highlighted in red to show the breakpoint.



However, typically the function we end up in from doing this isn't the one we need, so we want to go higher. BGB has a nice function for this!

Select "Run → Step out" (circled in red) to go up to the calling function. We will end up on the address after the call to this function.

If for any reason that doesn't work or you want to read all instructions from that function, you can do a Select "Run → Step over" (circled in blue)



If we step over, we want to find the jump that makes the routine go back to its caller, if we step over, we'll land at the same place. If we did it right, for our example, you will be right at **ROM0:089A** which, lucky for this project, is the beginning of our compression routine.

*(hl) = data pointer*
*ROM0:089A 2A ldi a,(hl)          ;get next byte, then inc data pointer*

\* Note: In many cases, that routine could be a subroutine or part of the patterns in the compression routine. For those cases you will have to spend more time with "step into" and "trace" until you go back to the main compression routine. Meanwhile you can comment and write up every operation the routine does, to use that info later for creating your own decompression tool.

In any case, from this point on, you set a breakpoint here (restart the emulation), and start to analyze the whole routine:

*The important thing to do is watch the registers and how they change as well as what memory is being accessed. These are the byte changes that we want to document and know, since this is what the routine does (to decompress) and where (to what is being decompressed or loaded)*

- To keep the same sample pattern for the whole document, as a way to compare and find consistency through the shared data between sections, I will copy the `BB` pattern below with comments.

**; Pattern found: BB 05 - will write 00 five times**
**ROM0:089B FE BB cp a,BB**
**ROM0:089D 28 22 jr z,08C1          ;write value 00, next byte is count**
**; …**
**; After the jump, the pattern is expanded into what the code does with the data**
**;BB routine**
**ROM0:08C1 AF xor a          ;a = 00**
**ROM0:08C2 4E ld c,(hl)          ;c = count**
**ROM0:08C3 18 07 jr 08CC          ;wait for H-Blank and write byte**

A bit of extra explanation here to know what we want to look for in a more generic approach:

*- First, we locate the jump, which is the place where the subroutine for the pattern will be "called"*

*- Second, we note what the code does right before, because that's the data it prepares and the data that wants to work with, in the subroutine.*
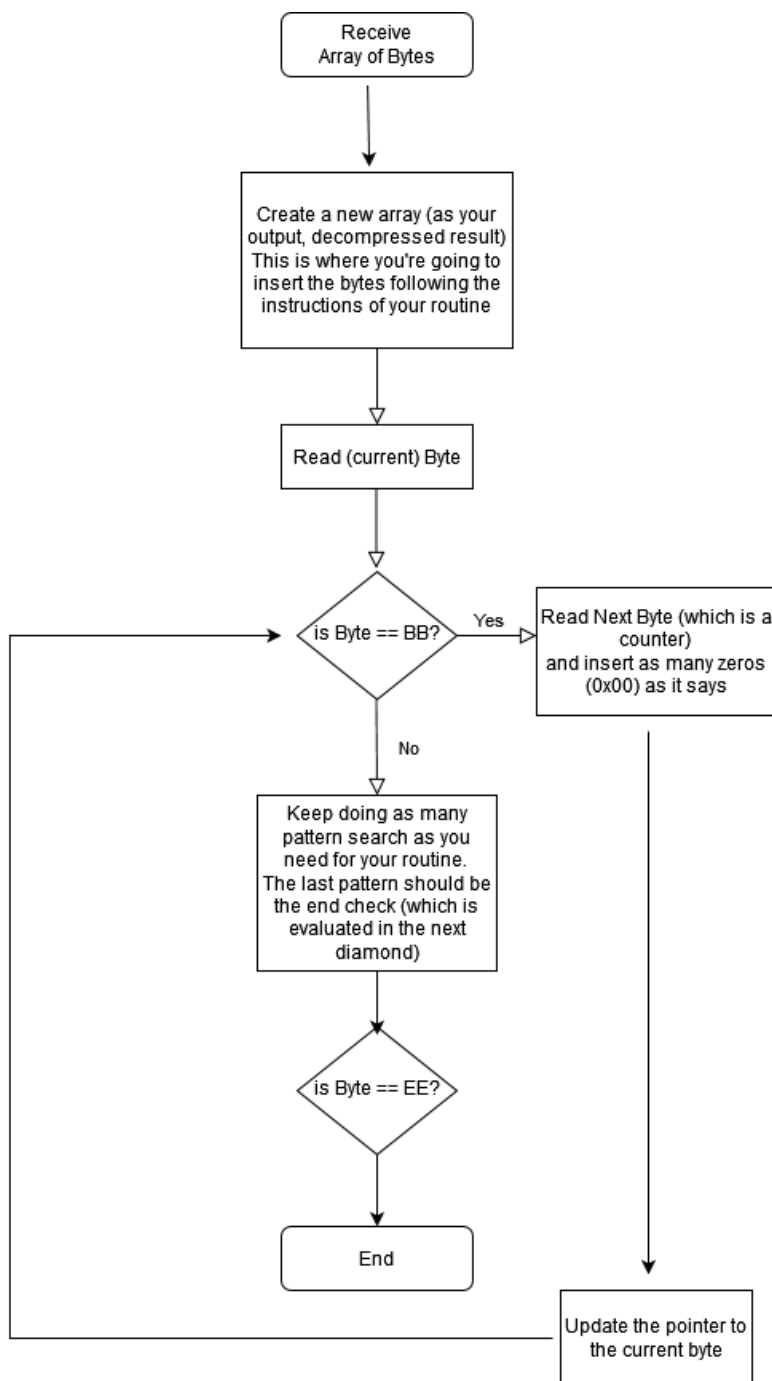
*- Third and last, we follow what the subroutine does by reading its operations with the received data, and try to guess what it can mean at a higher level (this last bit is probably the harder and the one you need to comment and practice the most).*

# Creating the tools - Decompression

This section uses the data from the previous analysis. The complexity here resides on how good your analysis is and which method you used, but the idea is the same: Create a script to launch that routine at will without the game emulator running. This means that you are going to take each step of the compression algorithm and make it into code in your coding language of choice. I will be using python here, but you can do the same with any other.

Below, the main flow is explain but only the first pattern is focused on, so you have a general idea without being overwhelmed, and at the same time, you have one pattern to check against:

*The flowchart based on the Static Analysis is shown below:*



Using this method, you need to develop your own flowchart. Hence why I am showing you one for this pattern, as example

VRAM $9010:  FF 00 FF 00 FF C0 3F 3F 00 00 00 C0 00 FF 00 FF

ROM  $2402A: FF 00 FF 00 FF C0 3F 3F BB 03 C0 00 FF 00 FF AA

> From these two, we can see a pattern where 'BB 03' means take the first Byte and repeat three times '00 00 00'. no control code at the end, just the next byte to display

*The Idea and code based on the Dynamic Analysis is shown below:*

BB 05 - will write 00 five times
ROM0:089B FE BB cp a,BB
ROM0:089D 28 22 jr z,08C1              ;write value 00, next byte is count

ROM0:08C1 AF xor a             ;a = 00
ROM0:08C2 4E ld c,(hl)         ;c = count
ROM0:08C3 18 07 jr 08CC         ;wait for H-Blank and write byte
        *Do the other patterns until reaching the end's control code and finish*

\* Note: You have seen the approach depending on which analysis method you want. As for Pros and Cons:

-   With the static method, you will have to somehow make yourself flowcharts of how the data evolves and why.

-   With the dynamic method you will have to port what the code does from one language to another.

If you followed the previous 'BB' patterns and made them into python, you will have something like this

```
def decompress(data):
        output = []            # create an output file
        i = 0
        while i < len(data):           # loop through the compressed  file and
                byte = data[i]
                i += 1
                if byte == 0xBB:     # check your patterns
                        count = data[i]
                        i += 1
                        output.extend([0x00] * count)    # Insert your bytes accordingly

# The code for the other patterns that would follow is omitted.
# This is to help visualize and highlight some key points at the end of the routine
# Those keys are written (and coded) below

                elif byte == 0xEE:  # Control code, Nobody wants an endless loop
                        break
                elif byte == 0x99:   # Be careful with overlooking some control patterns.
                        value = data[i]      # They might be scarce, so pay attention.
                        i += 1
                        output.append(value)
                else:                # Don't forget to copy the bytes that belong to no
                        output.append(byte)      # pattern.   They are part of the file too.
        return output
```

# Creating the tools - (Re)compression

This part of the document might be the one closest to dev than RE. However you can still do some steps of the analyzing methods to create and check your tool.

Warning: In contrast to the decompression method, you might have "limit cases" that won't be seen until you test the tool. The compression method has to be "exactly the same" not better and not worse, any difference will make the game read the inserted (compressed) image wrong.

- You will need a flowchart, based on the patterns you used for decompression, but this time  you have to code the tool to create the control bytes instead of finding them. See next page if you want to see an example for this project

- You have to take into account what can be considered a pattern and what must be passed as a plain Byte (one that is copied 'as is').

- You must not forget unique Bytes. An example in our marmalade routine would be the '0x99'. When you decompressed those, you probably saw them repeated and passed them as non important, but now you must be careful because if you see a '0x99', you have to insert an extra '0x99' to indicate to the game's decompression routine not to take it as a pattern warning for the next byte read.

> *# Check unique Bytes inside the loop through the file.*
> *if input_data[i] in [0xBB, 0xCC, 0xDD, 0xAA, 0x99, 0xEE]:*
> 　　*compressed.extend([0x99, input_data[i]])*
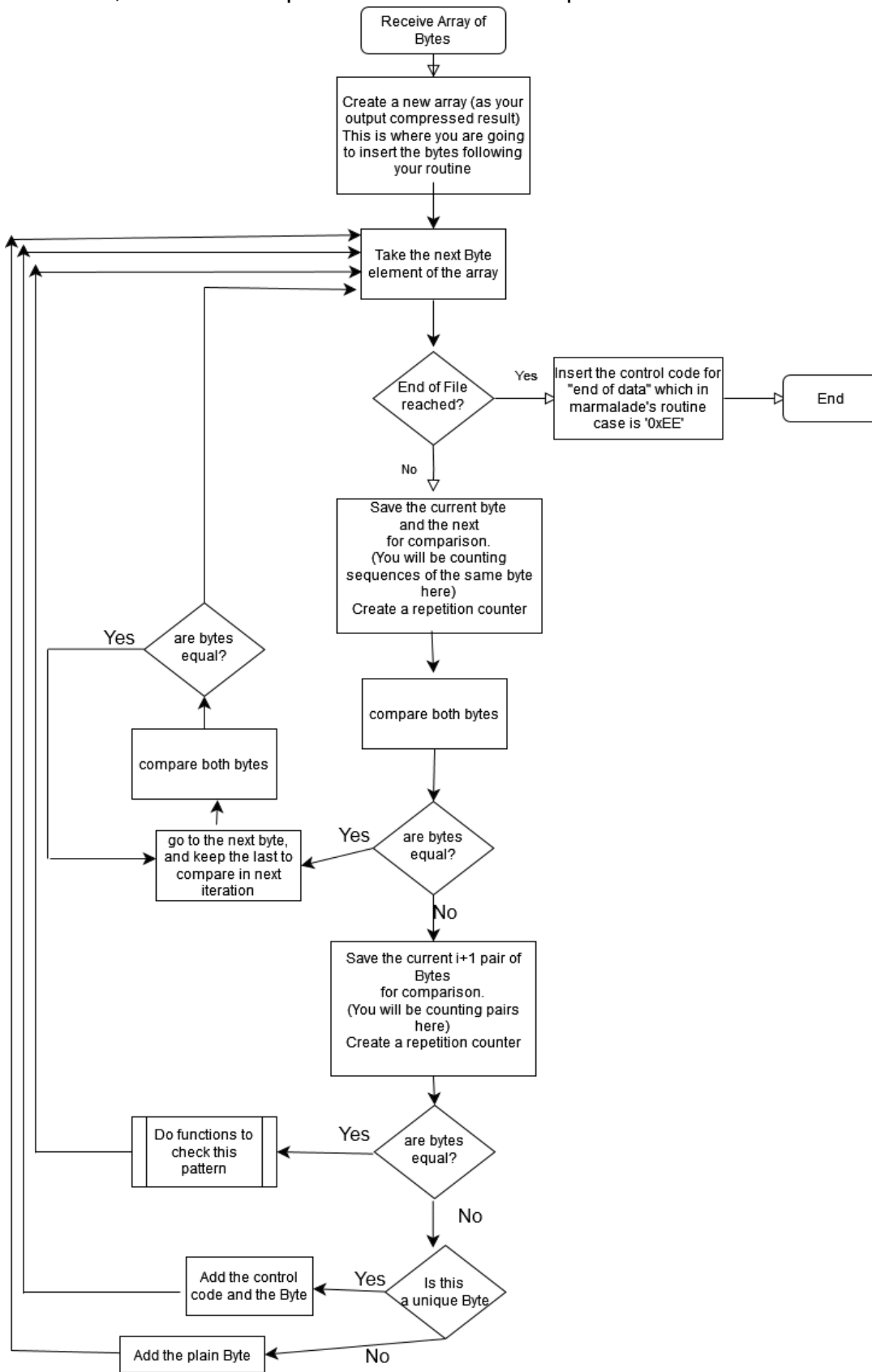> *else:*
> 　　*compressed.append(input_data[i])*

- Remember to add the "END" control code, in our marmalade boy's routine was '0xEE', but in some cases can be the length of the image file at the beginning or at the end of the file. If you don't do it, you may find errors like endless files or RAM overloads.

> *# After looping through the elements of your array/file, before the exit*
> *compressed.append(0xEE)  # Add the "END" control code*
> *return compressed*

- Your compression routine and its notes will come handy to develop a counterpart.

* Hint: You can use the original compressed image from the rom to compare with your compressed one side by side to find errors in the algorithm you have coded.

Flowchart, based on the patterns used for decompression:

Receive Array of Bytes

Create a new array (as your output compressed result) This is where you are going to insert the bytes following your routine

Take the next Byte element of the array

End of File reached? — Yes → Insert the control code for "end of data" which in marmalade's routine case is '0xEE' → End

No

Save the current byte and the next for comparison. (You will be counting sequences of the same byte here) Create a repetition counter

compare both bytes

are bytes equal? — Yes → go to the next byte, and keep the last to compare in next iteration → compare both bytes → are bytes equal? — Yes

No

Save the current i+1 pair of Bytes for comparison. (You will be counting pairs here) Create a repetition counter

are bytes equal? — Yes → Do functions to check this pattern

No

Is this a unique Byte — Yes → Add the control code and the Byte

No → Add the plain Byte

Snippet of the code for beginning of the routine and for the "Same-Byte" Sequence pattern finding:

```
def compress(input_data):
    compressed = bytearray()
    i = 0
    while i < len(input_data):
    # Same-Byte Sequence pattern finding
```

* Note: It uses a counter for 3 because 2 wouldn't save any space. Example: '00 00' would become 'BB 02', both cases require 2 bytes. [This can also be considered a special or limit case since some routines just transform the '00 00' into 'BB 02' either way].
Flowchart example for this compression routine

```
        if i + 1 < len(input_data):
        count = 1
        while i + count < len(input_data) and input_data[i] == input_data[i + count]:
            count += 1

        if count >= 3: # A sequence has been found
```

```
# The code for the other patterns that would follow is omitted.
# This is to help visualize and highlight some key points at the end of the routine
# Those keys are written (and coded) below

    # Check unique Bytes inside the loop through the file

    if input_data[i] in [0xBB, 0xCC, 0xDD, 0xAA, 0x99, 0xEE]:
            compressed.extend([0x99, input_data[i]])
    else:
            compressed.append(input_data[i])

    # After looping through the elements of your array/file, before the exit

    compressed.append(0xEE)  # Add the "END" control code
    return compressed
```

# Credits and Special Thanks