# CSC 3150 Assignment 3 Report

Tianhao Shi, 120090472

Nov, 2022

## Environment



- CUDA Version: 11.6



Also tested on the school's cluster

# Execution

## Method 1: Use makefile

```
cd Assignment_3_120090472/source
make #(compile+run)
# make build #(build only)
```
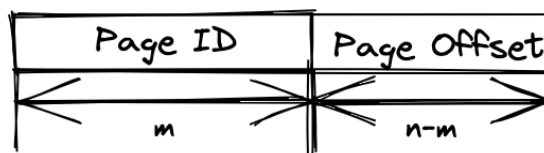
## Method 2: Batch file

```
cd Assignment_3_120090472/source
bash ./slurm.sh
```

# Design of Key Components

## Inverted Page Table

A page table is a one-to-one mapping betwen the virtual address and the physical address. In an inverted table, the index is the physical address (frames) and the virtual page id is stored as elements of the table. During address translation, the CPU issues a virtual address. The address is divided into two parts, **page id** and **page offset**, using the following rule:

For a virtual address with $n$ bits, if the page table entry is of size $2^m$, the upper $m$ bits are used as page id, and the rest $n$-$m$ bits are used as page offset.



## LRU

LRU, or "Least Recently Used", is a page replacement policy that OS use to decide which page to swap out in the inverted page table. This policy always choose, as the name suggests, the least recently used page to swap out of RAM (and page table). In my design, the underlying data structure is a **doubly linked list**. I store two pointers in the VirtualMemory struct, which tracks the *head* and *tail* of the linked list. The *head* points to the newest item in queue, while the tail points to the oldest. The update of LRU is not directly relevant to page hit/page miss. Instead, LRU is update when:

1. Page hit, but the requested page is not the latest used page.
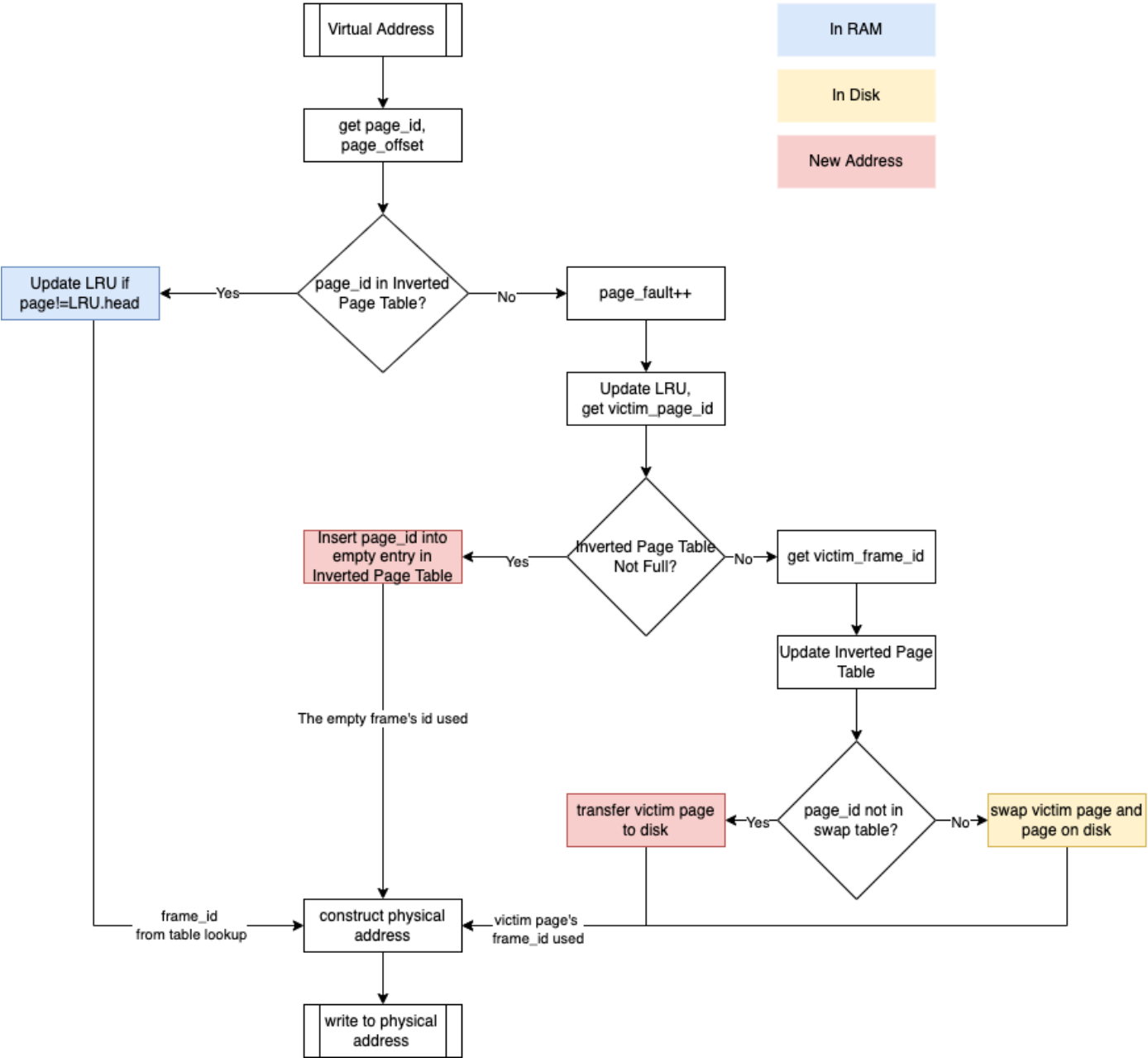2. Page miss, which means the page is previously not in the LRU list.

However, there are minor differences between the two scenarios. Namely, in the page hit scenario, all we need to do is traverse the LRU list and **perform a linear search**. The found node is moved to the head of the linked list.

A page miss indicates the requested page is not in the RAM. Therefore, there is no need to search. We can simply **append the node to the head** and update the list size. Notice that there is a limit for LRU size. When overflow occurs, the extraneous node is the LRU node and is removed from the linked list.

## Swap table

Our virtual memory size is 160KB, but the RAM size is only 32KB. This means for any input larger than the RAM size, there must be some data that resides in disk. To find the physical address of a page in disk, a swap table is used. The structure of swap table is similar to the inverted page table. The disk address is also partitioned into 32 byte pages, giving the swap table 4096 entries. Whenever swapping happens, the page id of the queried virtual address is searched in the swap table.

## Program Flow

The `vm_write` and `vm_read` function are almost identical. Here, I use `vm_write` as an example to illustrate the flow of my program.

For any virtual address, it is pointing to one of the following places:

1. A physical address in RAM
2. A physical address in Disk
3. Unallocated

If a virtual address is pointing to RAM, then the page_id must be stored in the inverted page table. Therefore, if the return value of querying page table is not -1, then we can construct the physical address using the corresponding frame_id.

If the page_id is not stored in RAM, **a page fault occurs**, and we must increment the page_fault counter.

If the page table is not full, then we can always insert the page to the first vacant place (there is no point swapping if we can store data in RAM).

However, if the page table is full, we **must move victim a page from RAM to disk**. The selection of this 'victim' is conducted using the previously mentioned **LRU** algorithm. If the swap table contains the requested page_id, then **the desired page is moved from disk to RAM**.

`vm_read` is almost identical except it reads from the final physical address.

# Bonus Design

Due to the limitted time, I implemented the bonus using the second version described in the additional note. The code is almost identical compared to the code in the first task.

However, notibaly the following things are added:

1. In `main.cu`, we need to launch four threads. Therefore, we call `mykernel <<<1, 4, INVERT_PAGE_TABLE_SIZE >>> (input_size);`
2. It seems resolving race condition is not very straightforward (hard to implement `mutex` used in assignment 2). However, the CUDA library do provide a `__syncthreads()` function. This function serves as a barrier, blocking threads until every thread reches this function. Using this function, the threads are executed consequtively, in the order of 0->1->2->3. Each time, the vm is initialized. However, the page fault pointer is not cleaned and is continuously updating, allowing the program to count the overall fault number.

```
1   if (threadIdx.x == 0) {
2     vm_init(&vm, data, storage, pt, &pagefault_num, PAGE_SIZE,
3       INVERT_PAGE_TABLE_SIZE, STORAGE_SIZE / PAGE_SIZE, PHYSICAL_MEM_SIZE, STORAGE_SIZE,
4       PHYSICAL_MEM_SIZE / PAGE_SIZE);
5     user_program(&vm, input, results, input_size);
6   }
7   __syncthreads();
8   if (threadIdx.x == 1) {
9     vm_init(&vm, data, storage, pt, &pagefault_num, PAGE_SIZE,
10      INVERT_PAGE_TABLE_SIZE, STORAGE_SIZE / PAGE_SIZE, PHYSICAL_MEM_SIZE, STORAGE_SIZE,
11      PHYSICAL_MEM_SIZE / PAGE_SIZE);
12    user_program(&vm, input, results, input_size);
13  }
14  __syncthreads();
15  if (threadIdx.x == 2) {
16    vm_init(&vm, data, storage, pt, &pagefault_num, PAGE_SIZE,
17      INVERT_PAGE_TABLE_SIZE, STORAGE_SIZE / PAGE_SIZE, PHYSICAL_MEM_SIZE, STORAGE_SIZE,
18      PHYSICAL_MEM_SIZE / PAGE_SIZE);
19    user_program(&vm, input, results, input_size);
20  }
21  __syncthreads();
22  if (threadIdx.x == 3) {
23    vm_init(&vm, data, storage, pt, &pagefault_num, PAGE_SIZE,
24      INVERT_PAGE_TABLE_SIZE, STORAGE_SIZE / PAGE_SIZE, PHYSICAL_MEM_SIZE, STORAGE_SIZE,
25      PHYSICAL_MEM_SIZE / PAGE_SIZE);
26    user_program(&vm, input, results, input_size);
27  }
28  __syncthreads();
```
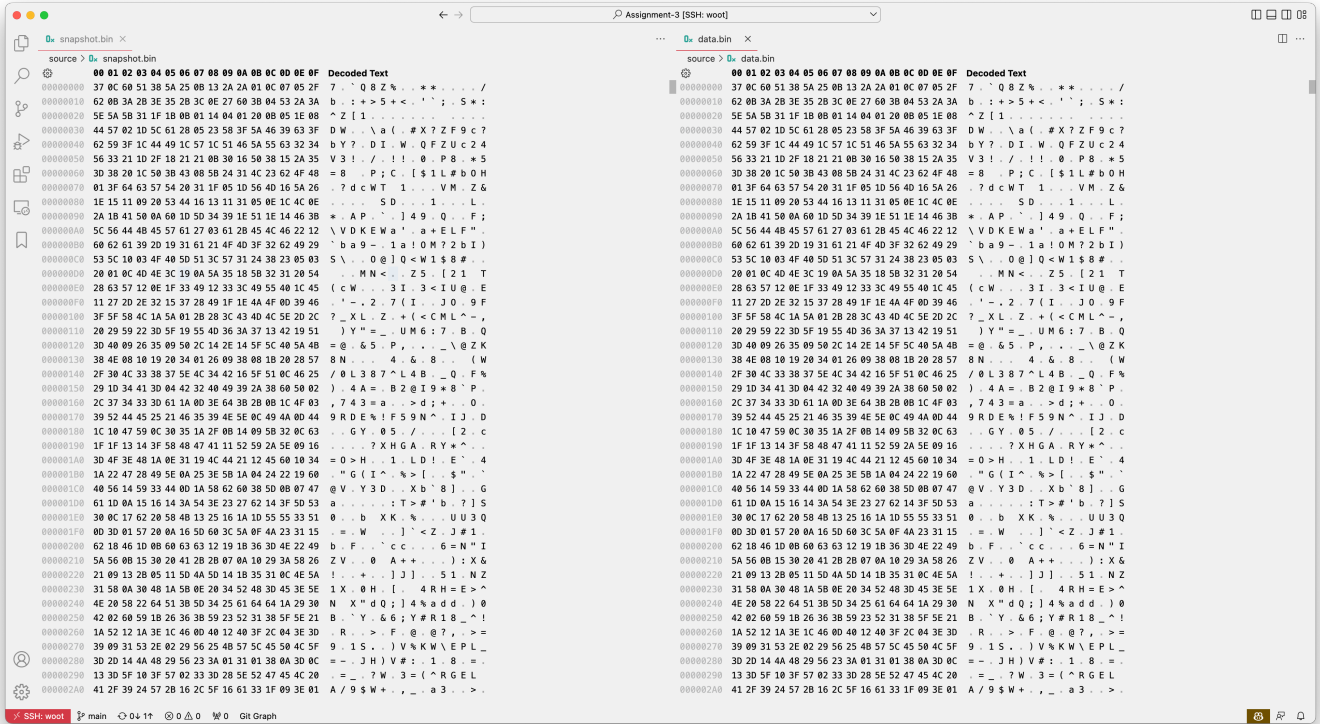
## Output and Analysis:

**For testcase 1, the output of my program is 8193:**

```
❯ make
nvcc --relocatable-device-code=true -g -G main.cu user_program.cu virtual_memory.cu -o main
./main
pagefault number is 8193
```

snapshot output

Cmp command:

```
❯ cmp -l data.bin snapshot.bin | gawk '{printf "%08X %02X %02X\n", $1, strtonum(0$2), strtonum(0$3)}'
```

Analysis:

The user program of testcase 1 does the following things:

```
1. Write [0, 132k) of the data.bin to VM starting at address 0 in ascending order.
1. Read virtual address [132k-1, 96k-1]
1. Take a snapshot of the virtual memory with offset 0
```

During these procedures, the first step contributes 132k/32=4k page faults, as they are all new pages. The second step only contributes to 1 page fault, because it reads all the content stored in the RAM (page hit), plus one element in the disk. The final step also contributes 4k page faults, because it reads from the beginning of VM. This means the pages loaded in step 2 are all replaced starting from address 0.

Threfore, totally we have $8 \times 2^{10} + 1 = 8193$ page faults.

A visualization is as followed

**For testcase 2, the output of my program is 9125**

```
❯ mv user_program.c user_program.c.back
❯ mv user_program2.c user_program.c
❯ make
nvcc --relocatable-device-code=true -g -G main.cu user_program.cu virtual_memory.cu -o main
./main
pagefault number is 9125
```

Analysis:

In this program the following things are done:

```
1. Write [0, 132k) of the data.bin to VM starting at address 32k in ascending order
1. Write [32k, 64k-1] of the data.bin to VM stasrting at address 0 in ascending order
1. Take a snapshot of the virtual memory with offset 32k
```

The first step populates the VM with 132K data, resulting in 4k page faults. The second step causes 1k-1 page faults, and the snapshot causes 4k page faults. Therefore, in total we have $9 * 2^{10} - 1 = 9215$ page faults.

**Bonus**

Only testcase 1 is used to test the bonus case. Since my code can be view as reading the .bin file four times, the total sum is 4*8193=32772 times.

```
586789    [Thread 3 Read] VA: 131063 Pid:4095 Poff:23 [HIT!!:32772] Fid: 1023 Val 51 PA 32759
586790    [Thread 3 Read] VA: 131064 Pid:4095 Poff:24 [HIT!!:32772] Fid: 1023 Val 48 PA 32760
586791    [Thread 3 Read] VA: 131065 Pid:4095 Poff:25 [HIT!!:32772] Fid: 1023 Val 63 PA 32761
586792    [Thread 3 Read] VA: 131066 Pid:4095 Poff:26 [HIT!!:32772] Fid: 1023 Val 45 PA 32762
586793    [Thread 3 Read] VA: 131067 Pid:4095 Poff:27 [HIT!!:32772] Fid: 1023 Val 97 PA 32763
586794    [Thread 3 Read] VA: 131068 Pid:4095 Poff:28 [HIT!!:32772] Fid: 1023 Val 61 PA 32764
586795    [Thread 3 Read] VA: 131069 Pid:4095 Poff:29 [HIT!!:32772] Fid: 1023 Val 23 PA 32765
586796    [Thread 3 Read] VA: 131070 Pid:4095 Poff:30 [HIT!!:32772] Fid: 1023 Val 68 PA 32766
586797    [Thread 3 Read] VA: 131071 Pid:4095 Poff:31 [HIT!!:32772] Fid: 1023 Val 93 PA 32767
586798    pagefault number is 32772
586799
```

# Problems Encountered & Solutions

## Implementation of LRU

At first I tried to use array as the underlying datastructure for implementing LRU. Because the vm_init function seems to suggest us to use the second 1k of inverted page table as the LRU array. This idea sounds easy, as the top element is the most recently used one. However, whenever we need to update the LRU array, we always need O(N) to move elements one by one. Moreover, when I tried to debug my programs, it seems that CUDA-GDB is having some trouble displaying value of the members' values of `struct VirtualMemory`. In the end I choose to replace the data structure with doubly linked list, which allowed the code to do insertion in O(1) time.

## Debugging Mechanism

Considering the large input size and heterogeneous nature of CUDA, finding an efficient debugging program is very difficult. I came up with three ways to debug my program:

1. Create the `tasks.json` and `launch.json` by referring to VSCode and Nsight studio doccumentation. I uploaded these two files on [github](). This method allows me to debug program using VSCode GUI in a step-wise manner. Moreover, at anytime I can examine the value of expressions to get accurate information. However, this method has two serious drawbacks. First, although VSCode attaches the program to the Cuda-GDB, the running speed of program is significantly slower than other debugging methods. Second, I find that the content of `printf` is not updated immediatly. I later found out that, according to [this post]() this is a feature of CUDA and there is nothing I can do about it. Therefore I also

use the following methods to facilitate the debugging procedure.

2. Using CUDA-GDB directly from command line. Although I am pretty new to GDB debugging and I don't know many of the commands and features, CUDA-GDB allows me to locate bugs more specifically (e.g, if SEGFAULT happens, I know which line triggers this).

3. Use rich `printf` in source code and run the executable file from the command line. Redirect stdout and stderr to a file and look for desired information in the file.

```
1    input size: 131072
2    [Write]Val: 55 VA: 32768 Pid:1024 Poff:0 [FAULT:1]Trying to put 1024 in list.Before update:Current LRU size:1Fid: 0 PA 0. LRU size:1.
3    [Write]Val: 12 VA: 32769 Pid:1024 Poff:1 [HIT!!:1] Fid: 0 PA 1. LRU size:1.
4    [Write]Val: 96 VA: 32770 Pid:1024 Poff:2 [HIT!!:1] Fid: 0 PA 2. LRU size:1.
5    [Write]Val: 81 VA: 32771 Pid:1024 Poff:3 [HIT!!:1] Fid: 0 PA 3. LRU size:1.
6    [Write]Val: 56 VA: 32772 Pid:1024 Poff:4 [HIT!!:1] Fid: 0 PA 4. LRU size:1.
7    [Write]Val: 90 VA: 32773 Pid:1024 Poff:5 [HIT!!:1] Fid: 0 PA 5. LRU size:1.
8    [Write]Val: 37 VA: 32774 Pid:1024 Poff:6 [HIT!!:1] Fid: 0 PA 6. LRU size:1.
9    [Write]Val: 11 VA: 32775 Pid:1024 Poff:7 [HIT!!:1] Fid: 0 PA 7. LRU size:1.
0    [Write]Val: 19 VA: 32776 Pid:1024 Poff:8 [HIT!!:1] Fid: 0 PA 8. LRU size:1.
1    [Write]Val: 42 VA: 32777 Pid:1024 Poff:9 [HIT!!:1] Fid: 0 PA 9. LRU size:1.
2    [Write]Val: 42 VA: 32778 Pid:1024 Poff:10 [HIT!!:1] Fid: 0 PA 10. LRU size:1.
3    [Write]Val: 1 VA: 32779 Pid:1024 Poff:11 [HIT!!:1] Fid: 0 PA 11. LRU size:1.
4    [Write]Val: 12 VA: 32780 Pid:1024 Poff:12 [HIT!!:1] Fid: 0 PA 12. LRU size:1.
5    [Write]Val: 7 VA: 32781 Pid:1024 Poff:13 [HIT!!:1] Fid: 0 PA 13. LRU size:1.
6    [Write]Val: 5 VA: 32782 Pid:1024 Poff:14 [HIT!!:1] Fid: 0 PA 14. LRU size:1.
7    [Write]Val: 47 VA: 32783 Pid:1024 Poff:15 [HIT!!:1] Fid: 0 PA 15. LRU size:1.
8    [Write]Val: 98 VA: 32784 Pid:1024 Poff:16 [HIT!!:1] Fid: 0 PA 16. LRU size:1.
9    [Write]Val: 11 VA: 32785 Pid:1024 Poff:17 [HIT!!:1] Fid: 0 PA 17. LRU size:1.
0    [Write]Val: 58 VA: 32786 Pid:1024 Poff:18 [HIT!!:1] Fid: 0 PA 18. LRU size:1.
1    [Write]Val: 43 VA: 32787 Pid:1024 Poff:19 [HIT!!:1] Fid: 0 PA 19. LRU size:1.
2    [Write]Val: 62 VA: 32788 Pid:1024 Poff:20 [HIT!!:1] Fid: 0 PA 20. LRU size:1.
3    [Write]Val: 53 VA: 32789 Pid:1024 Poff:21 [HIT!!:1] Fid: 0 PA 21. LRU size:1.
4    [Write]Val: 43 VA: 32790 Pid:1024 Poff:22 [HIT!!:1] Fid: 0 PA 22. LRU size:1.
5    [Write]Val: 60 VA: 32791 Pid:1024 Poff:23 [HIT!!:1] Fid: 0 PA 23. LRU size:1.
6    [Write]Val: 14 VA: 32792 Pid:1024 Poff:24 [HIT!!:1] Fid: 0 PA 24. LRU size:1.
7    [Write]Val: 39 VA: 32793 Pid:1024 Poff:25 [HIT!!:1] Fid: 0 PA 25. LRU size:1.
8    [Write]Val: 96 VA: 32794 Pid:1024 Poff:26 [HIT!!:1] Fid: 0 PA 26. LRU size:1.
9    [Write]Val: 59 VA: 32795 Pid:1024 Poff:27 [HIT!!:1] Fid: 0 PA 27. LRU size:1.
0    [Write]Val: 4 VA: 32796 Pid:1024 Poff:28 [HIT!!:1] Fid: 0 PA 28. LRU size:1.
1    [Write]Val: 83 VA: 32797 Pid:1024 Poff:29 [HIT!!:1] Fid: 0 PA 29. LRU size:1.
2    [Write]Val: 42 VA: 32798 Pid:1024 Poff:30 [HIT!!:1] Fid: 0 PA 30. LRU size:1.
3    [Write]Val: 58 VA: 32799 Pid:1024 Poff:31 [HIT!!:1] Fid: 0 PA 31. LRU size:1.
4    [Write]Val: 94 VA: 32800 Pid:1025 Poff:0 [FAULT:2]Trying to put 1025 in list.Before update:Current LRU size:2Fid: 1 PA 32. LRU size:2.
5    [Write]Val: 90 VA: 32801 Pid:1025 Poff:1 [HIT!!:2] Fid: 1 PA 33. LRU size:2.
6    [Write]Val: 91 VA: 32802 Pid:1025 Poff:2 [HIT!!:2] Fid: 1 PA 34. LRU size:2.
7    [Write]Val: 49 VA: 32803 Pid:1025 Poff:3 [HIT!!:2] Fid: 1 PA 35. LRU size:2.
```

# Learning Outcome

This process allowed me to get a more thorugh  and in-depth understanding of the paging policy, as well as the implementation and usage of virtual memory. It also allows me to gain some first hand-on experience on using CUDA. It is a very challenging project in terms of:

1. I have to come up how to implement a flawless paging policy from scratch
2. The input size of the binary file manes it hard to expose issues
3. Previously I have little experience in using GDB debugger (now I know how to use some of the basic features using CLI commands)
4. The CUDA-GDB itself has bugs, causing me to rely heavily on printf statements (on the bright side, I learned how to search for useful information in the log)

However, when I saw that I am able to handle 13000+ data without any mistake, I feel very proud of myself.

# Miscellaneous

I kept all the `printf` commands in my code, so it may print out a lot of information. Also, I included the outut.txt as a proof that my program DO run (hard to fake a file with ten thousad lines).