

Lab 4 report (Tianqi Gao, Weihang Chu)

```

GaoLiaoLiao — tgx588@joker:release — ssh — 90x38
lab1* lab2* lab3* lab4*
[tgx588@joker release]$ ./lab4

**=====**
Processing 16777216 elements...
Host CPU Processing time: 46.275002 (ms)
CUDA Processing time: 5.518000 (ms)
Speedup: 8.386190X
Test PASSED
[tgx588@joker release]$ ./lab4
[tgx588@joker src]$ ls
lab1 lab2 lab3 lab4 matrix1.txt matrix2.txt project
[tgx588@joker src]$ rm -r lab4
[tgx588@joker release]$ ./lab4
[tgx588@joker src]$ ls
lab1 lab2 lab3 matrix1.txt matrix2.txt project
[tgx588@joker src]$

**=====**
Processing 16777216 elements...
Host CPU Processing time: 46.356998 (ms)
CUDA Processing time: 5.510000 (ms)
Speedup: 8.413248X
Test PASSED
[tgx588@joker release]$ ./lab4
[tgx588@joker src]$ ls
lab1 lab2 lab3 lab4 matrix1.txt matrix2.txt project
[tgx588@joker src]$ cd l
lab1/ lab2/ lab3/ lab4/
[tgx588@joker src]$ cd lab4/

**=====**
Processing 16777216 elements...
Host CPU Processing time: 48.148998 (ms)
CUDA Processing time: 5.514000 (ms)
Speedup: 8.732136X
Test PASSED
[tgx588@joker release]$ ./lab4
[tgx588@joker lab4]$ vi scan_l
scan_largearray.cu scan_largearray_kernel.cu
[tgx588@joker lab4]$ vi scan_largearray.cu
[tgx588@joker lab4]$ make clean
[tgx588@joker lab4]$ make

**=====**
Processing 16777216 elements...
Host CPU Processing time: 46.219002 (ms)
CUDA Processing time: 5.523000 (ms)
Speedup: 8.368460X
Test PASSED
[tgx588@joker release]$

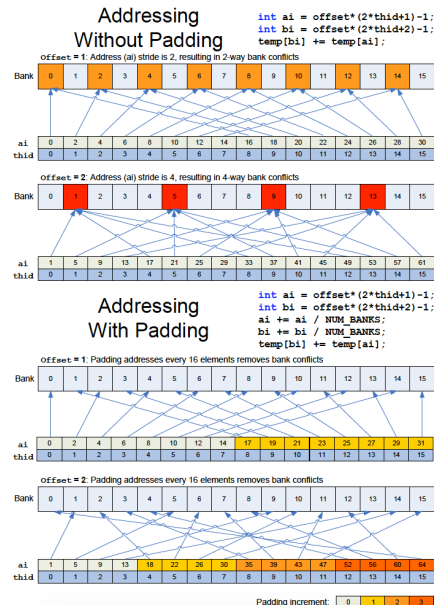
```

For the result, when we run our program without arguments and we processing 16777216 elements. The CPU processing time is 49ms, but the CuDA processing time is 5.543 ms. tH SPEEDUP IS 8.854772x. So the CUDA is more better.

For the handle of arrays not a power of two in size, we process the remain with a scan kernel modified to handle non power of 2 arrays. This kernel pads the shared memory array used out to the next higher power of two and initializes this extra memory to zero while loading in the data from device memory.

For the handle of the shared memory bank conflicts, we just need to pad the 2D array to a width that is not evenly divisible by number of shared memory banks. For this problem, we can achieve this by adding a variable amount of padding to each shared memory array index we compute. We just add to the index the value of the index divided by the number of shared memory banks

Below is a picture I found on the website which has the same idea of handle bank conflict.



In this lab, we also use the shared memory to speed up. What's more, in the beginning, we just implement the Hillis and Steele's algorithm which is not very good. Then, we use the Blelloch's algorithm, which can speed up the computation.

For the FLOPS rate in CPU, because the complexity of this algorithm is $O(n)$, so the $\text{FLOPS} = 16777216 / (49/1000) = 342392163$. For the FLOPS rate in CUDA, because the complexity of this algorithm is also $2n$, so the $\text{FLOPS} = 2 * 16777216 / (5.543/1000) = 6053478621.685$. for the CPU, the theoretical performance limits is 20G, for the CUDA it is 1t. For the bottlenecks, The CPU is write/read. For the CUDA, it is the number of threads.