

GPU Assignment - 2

Report by Kaustav Vats (2016048)

Assigning a word from a file to each thread in the grid. Checking for the pattern with offset 0, 1 & 2 on the same thread. 4 bytes word are assigned row wise in the grid [Threads in First row get 4 bytes in sequence].

Naive approach to do pattern matching on GPU would be to directly read 4 bytes from global memory. Which is very costly since we are trying to match word with various offsets, so we also required next 4 bytes. This means that each thread is doing two read operations from global memory.

Since each thread is reading extra 4 bytes, which are also read by next thread in the sequence. So i solved this problem using shared memory.

I created a 2D shared memory in which each block do at least 1 read operation from global memory and some of the boundary thread do 2 read operations.

```
const dim3 block_size(32, 32);  
const dim3 num_blocks(1024, 1024);  
__shared__ unsigned int sm_text[TOTAL][TOTAL+1];
```

In this Total is BlockDim.x, Right Boundary threads do Read operation twice.

Since the keywords used for pattern matching are used again and again by all threads, I created a 1D Shared Memory Array to store those elements.

```
__shared__ unsigned int sm_words[MAX_WORDS];
```

To increase the frequency of each word match, I'm doing atomic add to increment frequency for each keyword in global memory. Initially i was trying to do frequency increment in shared memory then in the end of each thread, transferring data from shared memory to global memory.

Note:-

Its mentioned on many of the blogs that atomicAdd is slower for shared memory as compared to global memory.

Timings\File Size	Small	Medium	Large
CPU Time	0.259128s	0.667559s	1.34634s
GPU Kernel Time	14.2845ms	14.8894ms	15.9051ms
GPU Kernel + Memory Transfer	15.7383ms	17.97ms	21.8767ms
Speedup	18.1405	44.8344	84.6481
Speedup with memory Transfer	16.4648	37.1486	61.5419

