

Assignment-2

Roll No	Name
2018007	Aditya Singh Rathore

Part 1 (SDT computation on GPU)

- (a) Write a CUDA version of the SDT computation using shared memory.
- (b) Document your approach to the problem.
- (c) Perform computations on CPU and GPU with image of sizes (256, 512, 1024, and 2048). Tabulate CPU and GPU (kernel and overall) timing results, plot speedups (kernel and overall), and report the `MSE` error in each case.

Note : Most optimised approach explained here. [Appendix, Kernel 5]

- The idea is compute minimum distances of each pixel from edge in one kernel and compute SDT from those minimum distances in second kernel.
- The number of edge pixels is calculated on the GPU.
 - Rather than atomic add to global variable for each thread, we have a block level variable and all threads add count atomically to that variable.
 - Once it has counted edge pixels in the given block, we can add the block level count to global count (once per block).
- Edge pixels are computed on the CPU. We cannot have multiple pixels elements writing to edge array, using atomic ops its a serial operation.
- Once we calculate minimum distance for each pixel, SDT calculation is straight forward.
- For computing minimum distance, we have to do the following:

```
1  for pixel in image:
2      d_min = inf
3      a, b = pixel coordinates
4      for edge in edges:
5          x, y = edge coordinates
6          d = (a-x)^2 + (b-y)^2
7          d_min = min(d_min, d)
8      min_dist[pixel] = d_min
```

- If we see calculating d , we can expand it further as:

```
1  d = (a-x)^2 + (b-y)^2
2  d = a^2 + b^2 + x^2 + y^2 - 2*(ax + by)
```

- We are computing a lot of those terms again and again. We can pre-compute them once and use them again as follows

```
1  struct edge{
2      x, y, sqr
3  }
```

```

4  Edges = struct edge[number of edges]
5  for edge in edges:
6      x, y = edge coordinates
7      Edges[edge].x = x
8      Edges[edge].y = y
9      Edges[edge].sqr = x^2 + y^2
10
11  for pixel in image:
12      d_min = inf
13      a, b = pixel coordinates
14      sqr = a^2 + b^2
15      for edge in Edges:
16          d = sqr + edge.sqr - 2*(a*edge.x + b*edge.y)
17          d_min = min(d_min, d)
18      min_dist[pixel] = d_min

```

- We can see that each pixel will access each edge. We can store edges in shared memory for faster access.
- But, size of edges can be very large. So we will divide the edges into chunks of 1024 and call the kernel again and again with each chunk to get the global minimum distance. (see optimization why 1024)

```

1  Edges = | c1 | c2 | c3 | ... | cn-1 | cn |
2  sz_edge = len(Edges)
3  size of c1 ... cn-1 = 1024
4  size of cn = sz_edge % 1024

```

- Calculating minimum distance for each pixel using chunks

```

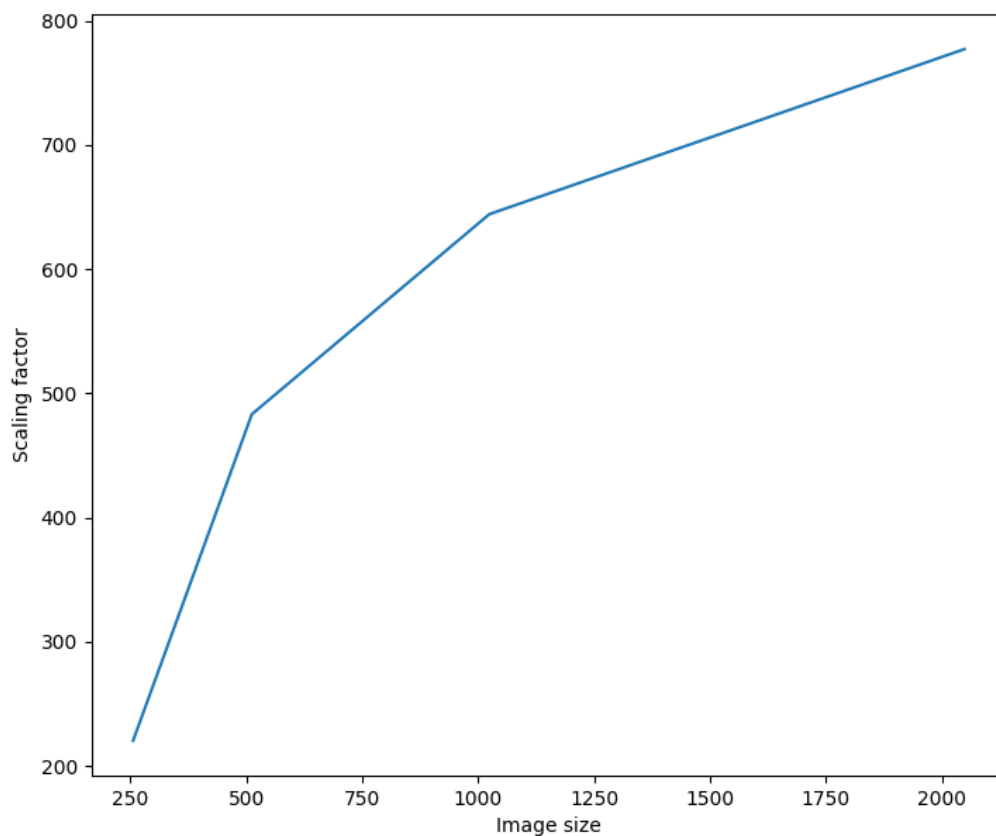
1  __global__ void compute_dist(float* min_dist, struct Edge*
    global_edges,int start)
2  {
3      extern __shared__ struct Edge edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5      edges[threadIdx.x] = global_edges[start + threadIdx.x];
6      __syncthreads();
7      if (i >= Sz[0]) return;
8      int x = i%Width[0];
9      int y = i/Width[0];
10     float sqr = x*x + y*y;
11     float min = min_dist[i];
12     float dist2;
13     for(int k = 0; k < blockDim.x; k++){
14         dist2 = edges[k].sqr + sqr - 2*(x*edges[k].x + y*edges[k].y);
15         if(dist2 < min) min = dist2;
16     }
17     min_dist[i] = min;
18 }

```

Results

Image Size	CPU time(Inter i7 8th generation, 16 GB Ram)	GPU time (IIITD server, NVIDIA1080) [Kernel-5]	Speed up	MSE
256	275.818 ms	1.25133 ms	220.419	0
512	4474.67 ms	9.26115 ms	483.16	0
1024	71157 ms	110.456 ms	644.2	0
2048	757062 ms	974.297 ms	777.03	0

- Scaling Factor



Part 2

How will you modify your approach to use constant memory instead of shared memory? Explain why using constant memory instead of shared memory is a good/bad choice in this case.

Note : Done for the most optimised approach. [Appendix, Kernel 5]

- We can store the computed edges in constant memory.
- But the constant memory is limited, we cannot do it for all images. Therefore it is a bad choice to use constant memory in this case.
- I was able to compute for image size `256x256` where number of edges was 2881.
 - I fixed number of edges beforehand.
- Instead of writing to constant memory, we can simply read from constant memory.

- For 256x256 image on NVIDIA GTX 1050 Ti, constant memory did not improve the performance.
 - Shared memory : 3.15299 ms
 - Constant memory : 3.45699 ms

Part 3 (Kernel Analysis)

(a) Analyze your CUDA kernel in terms of efficiency using nvprof / nvvp tool.

(b) Identify bottlenecks in your kernel.

- The main bottleneck is the computation of minimum distance for each pixel.
- I went through 5 different kernels [see Appendix]. Here are the bottlenecks identified in each kernel.

Kernel-1

- This was the most basic.
- Each pixel iterates over each thread to calculate minimum distance. All of this is fetched from the global memory.
- We can store edges in shared memory of each block.
- Didn't need a profiler to see this.
- Timings for this Kernel (NVIDIA GTX 1050 Ti)

Size	CPU Time	GPU Time	Speedup	MSE
256	275.818 ms	7.05322	39.10	0
512	4474.67 ms	107.066 ms	41.79	0
1024	71157 ms	1462.88 ms	48.64	0
2048	757062 ms	15873.7 ms	47.69	0

Kernel-2

- Read Appendix for idea
- Generated by improving Kernel-1
- Using profiler, I was able to identify that global load efficiency and shared memory efficiency were poor.

Results	
<p>Low Global Memory Load Efficiency (kernels accounting for 100% of compute have low efficiency (22.7% avg))</p> <p>Global load efficiency indicates how well the application's global loads are using device memory bandwidth. The efficiency is the number of bytes requested divided by the number of bytes that were transferred from device memory to satisfy those requests. Because device memory transfers bytes in blocks, the alignment and access pattern of a given load determines how many blocks must be transferred and thus determines the efficiency of that load. Low efficiency indicates that one or more global memory loads have a poor access pattern or alignment. Select this result to highlight kernels with low global load efficiency. More...</p>	
<p>Low Shared Memory Efficiency (kernels accounting for 100% of compute have low efficiency (3.1% avg))</p> <p>Shared memory efficiency indicates how well the application's shared memory accesses are using the available shared memory bandwidth. The efficiency is the number of shared loads and stores divided by the number of shared memory transactions required to perform those loads and stores. The alignment and access pattern of a given shared memory access determines how many transfers are required and thus determines the efficiency of that access. Low efficiency indicates that one or more shared memory accesses have a poor access pattern or alignment. Select this result to highlight kernels with low shared memory efficiency. More...</p>	
Queued	n/a
Submitted	n/a
Start	1.57616 s (1,576,159,6
End	1.60742 s (1,607,421,5
Duration	31.26194 ms (31,261,5
Stream	Default
Grid Size	[1025, 1, 1]
Block Size	[1024, 1, 1]
Registers/Thread	29
Shared Memory/Block	4 KiB
Launch Type	Normal
▼ Efficiency	
Global Load Efficiency	22.2%
Global Store Efficiency	100%
Shared Efficiency	3.1%
Warp Execution Efficiency	99.8%
Not-Predicated-Off Warp Execution Efficiency	89.9%

Problem

- There problem was that single thread was copying memory to shared memory [Appendix]

```
1  __global__ void compute_dist(float* min_dist, int* global_edges, int
2  height, int width, int start, int chunk_size)
3  {
4      ...
5      if(threadIdx.x == 0){
6          //Leader thread will write chunk to memory
7          for(int j = 0; j < chunk_size; j++){
8              edges[j] = global_edges[j+start];
9          }
10     __syncthreads();
```

- The timings for this kernel (NVIDIA GTX 1050 Ti)

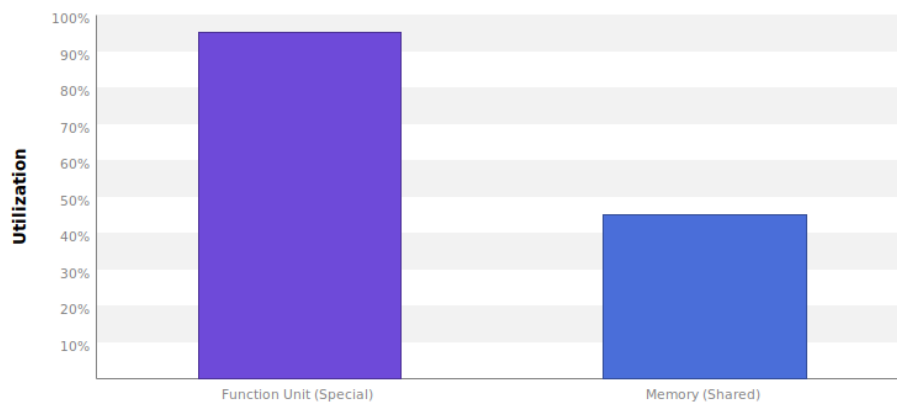
Size	CPU Time	GPU Time	Speedup	MSE
256	275.818 ms	6.42819	42.90	0
512	4474.67 ms	95.9526	46.63	0
1024	71157 ms	1399.2	50.85	0
2048	757062 ms	14311.8	52.89	0

Kernel-3

- Generated by improving Kernel-2
- See Appendix for Kernel
- Using profiler we identified computations as bottleneck

i Kernel Performance Is Bound By Compute

For device "GeForce GTX 1050 Ti" the kernel's memory utilization is significantly lower than its compute utilization. These utilization levels indicate that the performance of the kernel is most likely being limited by computation on the SMs.



Problem

- For computing minimum distance, we have to do the following:

```

1  for pixel in image:
2      d_min = inf
3      a, b = pixel coordinates
4      for edge in edges:
5          x, y = edge coordinates
6          d = (a-x)^2 + (b-y)^2
7          d_min = min(d_min, d)
8      min_dist[pixel] = d_min

```

- If we see calculating d , we can expand it further as:

```

1  d = (a-x)^2 + (b-y)^2
2  d = a^2 + b^2 + x^2 + y^2 - 2*(ax + by)

```

- We can compute $x^2 + y^2$ for each edge once.
- We can compute $a^2 + b^2$ for each pixel outside loop once.

Kernel-4

- Created by optimising Kernel-3

Problems

- We are still calculating x, y and $x^2 + y^2$ in each block for each edge.

```

1  __global__ void compute_dist(float* min_dist, int* global_edges,int start)
2  {
3      extern __shared__ int edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5      int edge = global_edges[threadIdx.x+start];
6      int _x = edge % Width[0];
7      int _y = edge / Width[0];
8      edges[3*threadIdx.x] = _x;
9      edges[3*threadIdx.x+1] = _y;
10     edges[3*threadIdx.x+2] = _x*_x + _y*_y;
11     __syncthreads();

```

- Timings for this kernel (NVIDIA GTX 1050 Ti)

Size	CPU Time	GPU Time	Speedup	MSE
256	275.818 ms	3.36058	82.08	0
512	4474.67 ms	44.8349	99.81	0
1024	71157 ms	560.283	127	0
2048	757062 ms	5474.75	138.2	0

Kernel-5

- Created by Optimising Kernel-4
- See Appendix
- Timings for this Kernel [NVIDIA GTX 1050 Ti]

Size	CPU Time	GPU Time	Speedup	MSE
256	275.818 ms	3.14413 ms	87.7	0
512	4474.67 ms	39.7117 ms	112.67	0
1024	71157 ms	527.827 ms	134.81	0
2048	757062 ms	5165.55 ms	146.57	0

Problems

- Kernel Calls could be paralleled.
- Shared Memory efficiency is still not 100%

Part 4 (Kernel Optimization)

- Rewrite a better version of your kernel based on your analysis (for example: improve occupancy, math performance, register/shmem usage, bank conflicts, coalescing, etc.).
- Document your optimization strategies. Compare and plot new speedups (kernel

and overall).

Kernel-1 -> Kernel-2

- Compute minimum distance for each pixel in different kernel.
- Compute `sdt` for each pixel in a different Kernel .
- For computing minimum distance, we will store the edges in shared memory because each pixel reads all edges, each block can store in shared memory which will be faster.
- We cannot store entire edge array in shared memory as it gets very large, so we will divide it into chunks, say 25 and call the kernel multiple times.
- For each pixel, we will initialise minimum distance to `FLT_MAX`.
- We will find the minimum distance with each chunk, first with chunk-1. Then with chunk-2 and so on. This way we will have the global minimum distance.

Kernel-2 -> Kernel-3

- The problem was low global load efficiency because a single thread was copying from global memory to shared memory.
- We fixed the chunk size = block size = 1024.
- Each thread in block copies from global memory to local memory

```
1  __global__ void compute_dist(float* min_dist, int* global_edges,int start)
2  {
3      extern __shared__ int edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5
6      edges[threadIdx.x] = global_edges[threadIdx.x+start];
7      __syncthreads();
8      ...
9  }
```

Results	
Low Shared Memory Efficiency [kernels accounting for 100% of compute have low efficiency (3.2% avg)]	
Shared memory efficiency indicates how well the application's shared memory accesses are using the available shared memory bandwidth. The efficiency is the number of shared loads and stores divided by the number of shared memory transactions required to perform those loads and stores. The alignment and access pattern of a given shared memory access determines how many transfers are required and thus determines the efficiency of that access. Low efficiency indicates that one or more shared memory accesses have a poor access pattern or alignment. Select this result to highlight kernels with low shared memory efficiency. More...	
Queued	n/a
Submitted	n/a
Start	1.49147 s (1,491,)
End	1.52124 s (1,521,)
Duration	29.77486 ms (29,)
Stream	Default
Grid Size	[1025,1,1]
Block Size	[1024,1,1]
Registers/Thread	28
Shared Memory/Block	4 KiB
Launch Type	Normal
Efficiency	
Global Load Efficiency	100%
Global Store Efficiency	100%
Shared Efficiency	3.2%
Warp Execution Efficiency	100%
Not-Predicated-Off Warp Execution Efficiency	90.1%
Occupancy	

- Global load efficiency is 100%

Kernel-3 -> Kernel-4

- Kernel is computation intensive.
- Simplified calculation by calculating edge variables once per block and pixel variable once

```
1  __global__ void compute_dist(float* min_dist, int* global_edges,int start)
2  {
3      ...
4      // Once per block
```



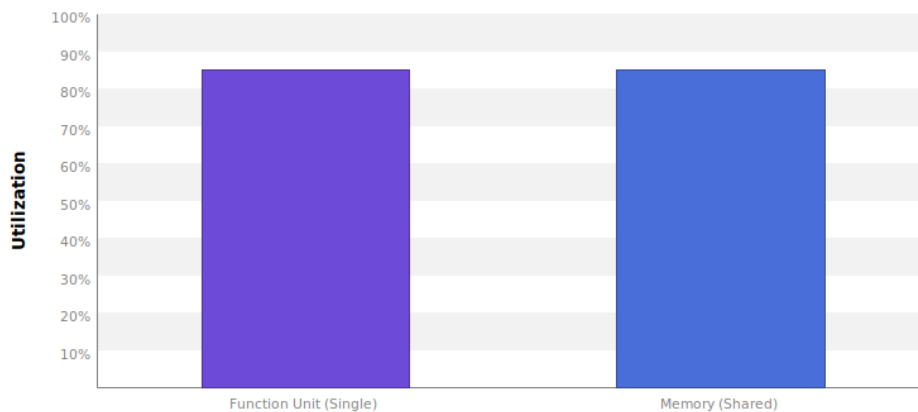
```

5     int edge = global_edges[threadIdx.x+start];
6     int _x = edge % Width[0];
7     int _y = edge / Width[0];
8     edges[3*threadIdx.x] = _x;
9     edges[3*threadIdx.x+1] = _y;
10    edges[3*threadIdx.x+2] = _x*_x + _y*_y;
11
12    ...
13    //Once per pixel
14    int x = i%Width[0];
15    int y = i/Width[0];
16    float sqr = x*x + y*y;
17    ...
18    for(int k = 0; k < blockDim.x; k++){
19        dist2 = edges[3*k+2] + sqr -2*(x*edges[3*k] + y*edges[3*k+1]);
20        if(dist2 < min) min = dist2;
21    }
22    ...
23 }

```

i High Compute And Memory Utilization

The kernel is utilizing greater than 80% of the available compute and memory performance of device "GeForce GTX 1050 Ti". These utilization levels indicate that additional performance improvement may be difficult to achieved for the kernel.



Kernel-4 -> Kernel-5

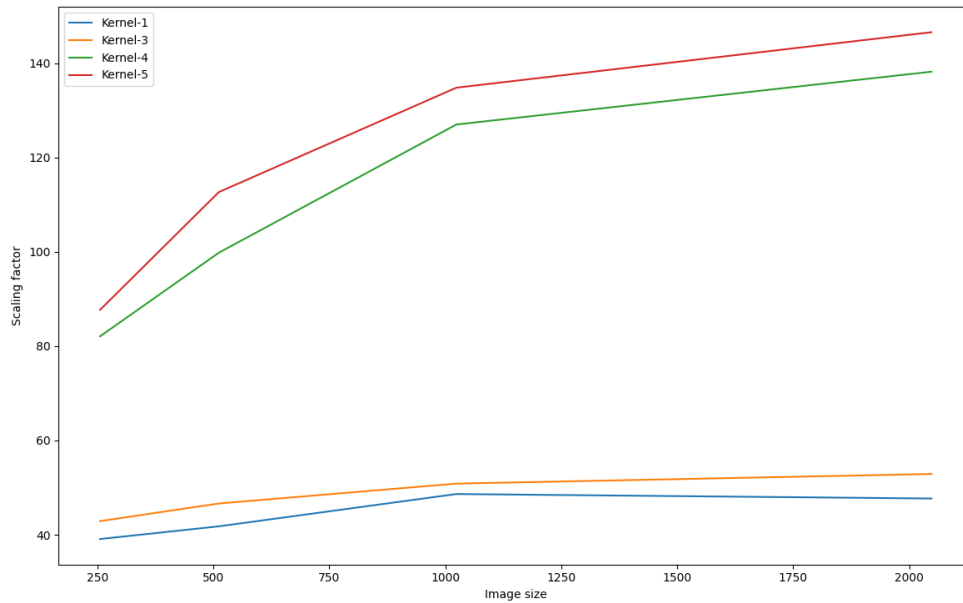
- Added Kernel to do edge computations only once/

Results

- On NVIDIA GTX 1050 Ti

Size	Kernel-1	Kernel-3	Kernel-4	Kernel-5
256	7.05322	6.42819	3.36058	3.14413
512	107.066	95.9526	44.8349	39.7117
1024	1462.88	1399.2	560.283	527.827
2048	15873.7	14311.8	5474.75	5165.55

- Scaling Factor comparison



Appendix

Kernel-1

- Compute Edges on the CPU.
- All Computations in one kernel.
- Each thread will compute SDT for each pixel.

```

1  __global__ void compute_sdt(float* sdt, unsigned char* bitmap, int* edges, int
height, int width, int edge_size)
2  {
3      int x = threadIdx.x + blockIdx.x*blockDim.x;
4      int y = threadIdx.y + blockIdx.y*blockDim.y;
5      if ( x < width && y < height){
6          float min_dist = FLT_MAX;
7          for(int k=0; k<edge_size; k++)
8          {
9              float _x = edges[k] % width;
10             float _y = edges[k] / width;
11             float dx = _x - x;
12             float dy = _y - y;
13             float dist2 = dx*dx + dy*dy;
14             if(dist2 < min_dist) min_dist = dist2;
15         }
16         float sign = (bitmap[x + y*width] >= 127)? 1.0f : -1.0f;
17         sdt[x + y*width] = sign * sqrtf(min_dist);
18     }
19 }

```

Kernel-2

- Compute minimum distance for each pixel in different kernel.
- Compute `sdt` for each pixel in a different Kernel .
- For computing minimum distance, we will store the edges in shared memory because each pixel reads all edges, each block can store in shared memory which will be faster.
- We cannot store entire edge array in shared memory as it gets very large, so we will divide it into chunks, say 25 and call the kernel multiple times.
- For each pixel, we will initialise minimum distance to `FLT_MAX`.
- We will find the minimum distance with each chunk, first with chunk-1. Then with chunk-2 and so on. This way we will have the global minimum distance.

Kernel

```

1  __global__ void compute_dist(float* min_dist, int* global_edges, int height, int
width, int start, int chunk_size)
2  {
3      extern __shared__ int edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5      if(threadIdx.x == 0){
6          //Leader thread will write chunk to memory
7          for(int j = 0; j < chunk_size; j++){
8              edges[j] = global_edges[j+start];
9          }
10     }
11     __syncthreads();
12     if (i >= height*width) return;
13
14     int x = i%width;
15     int y = i/width;
16     float min = min_dist[i];
17     float _x, _y, dx, dy, dist2;
18     for(int k = 0; k < chunk_size; k++){
19         _x = edges[k] % width;
20         _y = edges[k] / width;
21         dx = _x - x;
22         dy = _y - y;
23         dist2 = dx*dx + dy*dy;
24         if(dist2 < min) min = dist2;
25     }
26     min_dist[i] = min;
27 }
28
29 __global__ void compute_sdt(unsigned char* bitmap, float* min_dist, float* sdt, int
sz){
30     int i = threadIdx.x + blockDim.x*blockIdx.x;
31     if(i >= sz) return;
32     float sign = (bitmap[i] >= 127)? 1.0f : -1.0f;
33     sdt[i] = sign * sqrtf(min_dist[i]);
34 }

```

Calling

```

1  const auto num_chunks = 25;
2  const auto chunk_size = sz_edge/num_chunks;
3  const auto last_chunk = sz_edge%num_chunks;
4  const auto block_size = 256;

```

```

5  const auto grid_size = (sz/block_size) + 1;
6  for(int i = 0; i < num_chunks; i++){
7      compute_dist<<< grid_size, block_size, chunk_size*sizeof(int)>>>
      (d_min_dist.arr(),
8       d_edges.arr(),
9       height, width,
10      i*chunk_size,
11      chunk_size);
12  }
13  if (last_chunk != 0){
14      compute_dist<<< grid_size, block_size, last_chunk*sizeof(int)>>>
      (d_min_dist.arr(),
15       d_edges.arr(),
16       height, width,
17       num_chunks*chunk_size, last_chunk);
18  }

```

Kernel-3

- Same as Kernel-2.
- Global load store efficiency is improved.
- Edges is divided into equal sized chunks of size 1024.
- Instead of first thread in block writing to shared memory, all threads write to shared memory (chunk size = block size).

Kernel

```

1  __global__ void compute_dist(float* min_dist, int* global_edges,int start)
2  {
3      extern __shared__ int edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5
6      edges[threadIdx.x] = global_edges[threadIdx.x+start];
7      __syncthreads();
8
9      int sz = Sz[0];
10     if (i >= sz) return;
11
12     int width = Width[0];
13
14     int x = i%width;
15     int y = i/width;
16     float min = min_dist[i];
17     float _x, _y, dx, dy, dist2;
18     for(int k = 0; k < blockDim.x; k++){
19         int edge = edges[k];
20         _x = edge % width;
21         _y = edge / width;
22         dx = _x - x;
23         dy = _y - y;
24         dist2 = dx*dx + dy*dy;
25         if(dist2 < min) min = dist2;
26     }
27     min_dist[i] = min;
28 }

```

Launch

```
1  const auto chunk_size = 1024;
2  const auto num_chunks = sz_edge/chunk_size;
3  const auto last_chunk = sz_edge%chunk_size;
4
5  const auto block_size = chunk_size;
6  const auto grid_size = (sz/block_size) + 1;
7  const auto grid_last_chunk = (sz/last_chunk) + 1;
8  for(int i = 0; i < num_chunks; i++){
9      compute_dist<<< grid_size, block_size, chunk_size*sizeof(int)>>>
10     (d_min_dist.arr(),
11      d_edges.arr(),
12      i*chunk_size);
13 }
14 if (last_chunk != 0){
15     compute_dist<<< grid_last_chunk, last_chunk, last_chunk*sizeof(int)>>>
16     (d_min_dist.arr(),
17      d_edges.arr(),
18      num_chunks*chunk_size);
19 }
```

Kernel-4

- When computing `dist2`, each thread will be computing some computations that can be done only once at block level (explained in Kernel Optimization).

Kernel

```
1  __global__ void compute_dist(float* min_dist, int* global_edges,int start)
2  {
3      extern __shared__ int edges[];
4      int i = threadIdx.x + blockIdx.x*blockDim.x;
5      int edge = global_edges[threadIdx.x+start];
6      int _x = edge % Width[0];
7      int _y = edge / Width[0];
8      edges[3*threadIdx.x] = _x;
9      edges[3*threadIdx.x+1] = _y;
10     edges[3*threadIdx.x+2] = _x*_x + _y*_y;
11     __syncthreads();
12
13     int sz = Sz[0];
14     if (i >= sz) return;
15
16     int x = i%Width[0];
17     int y = i/Width[0];
18     float sqr = x*x + y*y;
19     float min = min_dist[i];
20     float dist2;
21     for(int k = 0; k < blockDim.x; k++){
22         dist2 = edges[3*k+2] + sqr -2*(x*edges[3*k] + y*edges[3*k+1]);
23         if(dist2 < min) min = dist2;
24     }
25
26     min_dist[i] = min;
27 }
```

Kernel-5

- We can count number of edge pixels on the GPU.
- We are calculating `x`, `y` and `sqr` for each edge in every block.
- We can still do better by computing them globally and creating an `array of structs` for each edge.

Kernel

```
1 struct Edge{
2     float x, y, sqr;
3 };
```

- Count the number of pixels.

```
1 __global__ void count_edge_pixels(unsigned char* bitmap, int* sz_edges){
2     __shared__ int num_edges;
3     int i = threadIdx.x + blockDim.x*blockIdx.x;
4     if(i >= Sz[0]) return;
5
6     if (threadIdx.x == 0){
7         num_edges = 0;
8     }
9     __syncthreads();
10
11     if (bitmap[i] == 255){
12         atomicAdd(&num_edges, 1);
13     }
14     __syncthreads();
15     if (threadIdx.x == 0){
16         atomicAdd(sz_edges, num_edges);
17     }
18 }
```

- Calculate `x`, `y` and `sqr` for each edge index calculated on the CPU.

```
1 __global__ void compute_edge_pixels(struct Edge* edges, int* edge_indices){
2     int i = threadIdx.x + blockDim.x*blockIdx.x;
3     if(i >= Sz_edges[0]) return;
4     int edge = edge_indices[i];
5     int _x = edge % Width[0];
6     int _y = edge / Width[0];
7     edges[i].x = _x;
8     edges[i].y = _y;
9     edges[i].sqr = _x*_x + _y*_y;
10 }
```

- We create shared memory of `struct Edge`

```
1 __global__ void compute_dist(float* min_dist, struct Edge* global_edges, int start)
2 {
3     extern __shared__ struct Edge edges[];
4     int i = threadIdx.x + blockDim.x*blockIdx.x;
```

```
5     edges[threadIdx.x] = global_edges[start + threadIdx.x];
6     __syncthreads();
7     if (i >= Sz[0]) return;
8     int x = i%Width[0];
9     int y = i/Width[0];
10    float sqr = x*x + y*y;
11    float min = min_dist[i];
12    float dist2;
13    for(int k = 0; k < blockDim.x; k++){
14        dist2 = edges[k].sqr + sqr - 2*(x*edges[k].x + y*edges[k].y);
15        if(dist2 < min) min = dist2;
16    }
17    min_dist[i] = min;
18 }
```