

Greedy Triangulation via GPU

Randy Shoemaker

Department of Computer Science

College of William and Mary

Williamsburg, USA

rwshoemaker@email.wm.edu

Abstract—We present a GPU implementation of the greedy triangulation algorithm from [1]. The version from Shoemaker et al. utilized the Message Passing Interface (MPI) which parallelized the greedy triangulation algorithm to run on a cluster. We present a version utilizing CUDA to run on many streaming multiprocessors which achieves significant speed up over Shoemaker et al. Our GPU implementation and our experimental data can be found in our Github repository¹.

I. PROJECT NOTE

Professor Ren,

I included this section to label which portion to use as the project and which portion to use as “extra credit”. In my project proposal I only set out to produce a fast version on the GPU so consider the code in the repository that accomplishes this as (particularly the .cu file) as well as the statistics in log.csv as the non-extra-credit portion. When we last spoke you mentioned that I could do a theoretical analysis using FLOPS. I have done this and it is included in Section V. These theoretical results are used in the results and experiments section too for aide my analysis. Please consider the theoretical stuff as extra credit, as the simplifying assumptions made during the derivations may not hold. In the repository I also included all of my previous attempts during this project to parallelize this algorithm along with a log (all in the file old_versions). You may disregard these. By the way, I also included a “more professional” looking version of this document in the repository that does not contain this section. I hope you enjoy this work.

Thanks!

Randy

II. INTRODUCTION

The greedy triangulation algorithm takes a set of points in the plane and greedily produces a triangulation. It builds the triangulation in three steps. In the first step all lines between each pair of points is generated, then these lines are sorted in increasing size. The third step is the triangulation step, during this step the smallest line is added to the triangulation and all lines that conflict with this line are removed, then the next smallest line is added and this step is repeated. The triangulation step is the most costly step of the algorithm. The first two steps combined take a very small fraction of the overall run time of the algorithm. Due to this fact we focused our efforts on parallelizing the triangulation step using CUDA.

The greedy triangulation algorithm is known to give a \sqrt{n} -approximation of the Minimum Weight Triangulation problem (MWT) [1]. This problem has been shown to be NP-hard [2]. Triangulation problems arise in many domains, such as computational geometry and geometry processing. Thus, quickly producing a triangulation is useful for many applications. We present a parallel version of the greedy triangulation algorithm which utilizes CUDA to run on many streaming multiprocessors. Previous work has produced a parallel version which runs on a cluster utilizing the Message Passing Interface (MPI) [1]. Our approach achieves significant speedup when compared to the cluster version. We present an empirical comparison to the previous work as well as a comparison to the serial version. We also present a theoretical analysis of our approach vs others. In the appendix we include images of some triangulations produced by our code.

In Section III we discuss the greedy triangulation algorithm and in Section IV we discuss our implementation. In Section V we present a theoretical analysis of our algorithm and in Section VI we present our empirical results. In Section VII we provide a discussion of our results before concluding in Section VIII. Lastly, we discuss potential future work.

III. BACKGROUND

Before discussing the greedy triangulation algorithm we first present some notation. Let P be a point set we wish to triangulate using the greedy algorithm; for convenience we let $n = |P|$. Let \mathcal{L} denote the set of all line segments formed by connecting every pair of points in P . Let P_i denote the i^{th} element of P . Given P the greedy triangulation algorithm produces the greedy triangulation \mathcal{T} . With this notation we present the greedy triangulation algorithm.

```
Input:  $P$  // The point set
Output:  $\mathcal{T}$  // The triangulation
 $\mathcal{L} \leftarrow \emptyset$ 
 $n \leftarrow |P|$ 
 $\text{idx} \leftarrow 0$ 
// Line generation step
for  $i \leftarrow 0$  to  $n$ 
    for  $j \leftarrow i+1$  to  $n$ 
         $\mathcal{L}_{\text{idx}} \leftarrow$  the line from  $P_i$  to  $P_j$ 
         $\text{idx} \leftarrow \text{idx} + 1$ 
// Sorting step
Sort  $\mathcal{L}$  in ascending order
```

¹The url for our repository: https://github.com/shoemarw/gpu_greedy_triangulation

```

// Triangulation step
 $\mathcal{T} \leftarrow \emptyset$ 
while  $\mathcal{L} \neq \emptyset$ 
   $I \leftarrow \mathcal{L}_0$ 
  remove  $I$  from  $\mathcal{L}$ 
   $\mathcal{T} \leftarrow \mathcal{T} \cup \{I\}$ 
  for  $J \in \mathcal{L}$ :
    if (not share_endpoint( $I, J$ )
        && intersects( $I, J$ )):
      remove  $J$  from  $\mathcal{L}$ 
return  $\mathcal{T}$ 

```

The function `share_endpoint` takes two lines and returns true iff the lines share an endpoint and `intersects` takes two lines and returns true iff they intersect. When implemented, the first two steps, the generation and sorting steps, take up a very small portion of the over all run time. Below, we have included some sample run times of a serial implementation of the algorithm. The run times are in seconds.

Points	Generation	Sort	Run time (seconds)
724	0.0149	0.0532	7.9712
1024	0.0236	0.1027	24.996
1448	0.0724	0.2175	98.2896
2048	0.0847	0.4431	385.6797
2896	0.1786	0.9174	1407.0545

Since the triangulation step is the most computationally intensive, our work focuses on parallelizing this step.

In the following section we present our parallel version of the algorithm.

IV. THE PARALLEL ALGORITHM

We now present our parallel version. First we must establish some further notation. Let T be an array of boolean values which are initially all true. During the execution of the algorithm, T_i becomes false exactly when \mathcal{L}_i is found to conflict with a line that is in \mathcal{T} . Once the algorithm has concluded T_i is true iff \mathcal{L}_i belongs to the triangulation. Below is the pseudo code for our implementation. Since the triangulation step is the only step parallelized we only include pseudo code for it.

```

Input:  $\mathcal{L}$  // The sorted lines
       $b$  // Number of blocks
       $t$  // Number of threads
Output:  $\mathcal{T}$  // The triangulation
// Triangulation step
 $T \leftarrow \{\text{True}_i\}_{i=1}^{|\mathcal{L}|}$  // Include every line
 $s \leftarrow 0$  // Index of smallest line
while  $s < |\mathcal{L}|$ 
  triangulate( $b, t$ ) ( $\mathcal{L}, s, T$ )
  find_new_small( $1, 1$ ) ( $s, T$ )
 $\mathcal{T} \leftarrow \{\mathcal{L}_i \in \mathcal{L} : T_i \text{ is True}\}$ 
return  $\mathcal{T}$ 

```

The implementation of this part of the algorithm runs on the host. The implementations of the functions `triangulate` and `find_new_small` run on the device and are launched

from the host. These kernel calls allow us to synchronize the blocks. `triangulate` takes a line and eliminates all lines which conflict with it by updating the corresponding elements of T to false. `find_new_small` updates the index s of the smallest line after a round of eliminations. These functions have no return values and manipulate T and s respectively on device memory to minimize data transfers between the host and the device. Below we have included pseudo code for `triangulate` and `find_new_small`. The number of blocks b and threads per block t are command line arguments for the implementation for experimental purposes.

```

triangulate( $\mathcal{L}, s, T$ )
   $j \leftarrow \text{blockID.x} * \text{blockDim.x} + \text{threadID.x}$ 
  while  $j < |\mathcal{L}|$ 
    if  $T_j$  is True and conflicts( $\mathcal{L}_s, \mathcal{L}_j$ )
       $T_j \leftarrow \text{False}$ 
     $j \leftarrow j + \text{blockDim.x} * \text{gridDim.x}$ 

find_new_small( $s, T$ )
   $s \leftarrow s + 1$ 
  while  $s < |\mathcal{L}|$  and  $T_s$  is false
     $s \leftarrow s + 1$ 

```

`find_new_small` simply increments the index of the smallest line so that the eliminated lines are skipped over. `triangulate` eliminates lines that conflict with the smallest line during a round in parallel. This task occurs over b blocks, each containing t threads. After a round of eliminations we must increment s on one block while the others wait. In order to accomplish this kind of synchronization we use kernel calls from the host.

The device function `triangulate` is written so that the cache is used as efficiently as possible, we used Nvidia's blogs² as a guide. The data that each thread operates on was tuned to achieve maximum memory coalesce. `triangulate` calls the device function `conflicts` which tests two lines to see if they conflict with one another, it in turn calls other functions. We include pseudo code for these functions in the appendix.

In the following section we present a theoretical analysis of our implementation.

V. THEORETICAL ANALYSIS

In this section we present an analysis of the number of floating point operations performed by the triangulation step of the greedy triangulation algorithm. For a given number of points n we compute an approximation of the number of floating point operations required by using an average case analysis.

Let N_n be the total number of floating point operations needed to produce a triangulation during the triangulation step. Notice that N_n can be computed if we know the number of conflict tests are executed and the number of floating point operations per conflict test. So let C_n be the number of

²Please see <https://devblogs.nvidia.com/how-access-global-memory-efficiently-cuda-c-kernels/>

conflict tests needed for a point set with n elements, and let F_n be the number of floating point operations required for one conflict test if there are n points in the point set. C_n clearly depends on n , to see why F_n depends on n notice that `share_endpoint` is computed for every conflict test and this function makes several calls `is_equal`. C_n depends on n because the probability that `is_equal` is true for two given points depends on the number of lines which is a function of n . Since we are doing an average case analysis this probability is used to estimate the floating point operations needed. In the appendix we provide the pseudo code for every function analyzed here.

We now compute an approximation of C_n . Let A_n be the average number of lines in the greedy triangulation of a point set with n elements. A_n is important because, on average, it is the number of iterations of the while loop in `triangulate` averaged across all point sets with n points. We ran experiments using different point sets to arrive at a suitable approximation of A_n . For several values of n we generated five unique point sets. The values of n used were 4, 16, 64, 256, 512, 1024, and 2048. For each n we found that the number of lines in the triangulation were always close. For each n we took the average of the number of lines in the triangulations³. The averages are included in the below table

Points	approximation of A_n
4	5
16	37.2
64	182.6
256	858.4
512	1963.4
1024	4874.8
2048	14056

We then fit a quadratic to the approximations of A_n , we use a quadratic because A_n appears to grow on the order of n^2 . The resulting function is

$$A_n = 2.01931 * 10^{-3} n^2 + 2.7221n + 5.27957$$

Now that we have the number of iterations of the while loop A_n , to compute C_n we need to estimate the number of conflict tests during a typical iteration of the while loop. For the sake of analysis we assume that an equal number of lines are eliminated during each iteration of the while loop. We believe that this assumption is reasonable when averaged across all point sets for a given n . Under this assumption $\frac{|\mathcal{L}_n| - A_n}{A_n}$ lines are eliminated during an iteration of the while loop. It follows that at iteration i there are $|\mathcal{L}_n| - (i - 1) \frac{|\mathcal{L}_n| - A_n}{A_n} - 1$ left that could potentially be part of the triangulation. That means that at iteration i each of these lines must be tested using `conflicts`. Therefore at iteration i there must be

³This data is included in the repository in the folder `find_An`. This folder also contains code to generate this experimental data.

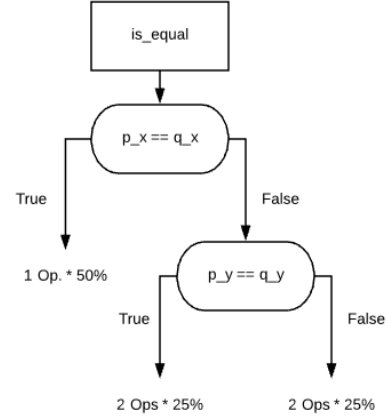


Fig. 1. The execution tree for `is_equal`.

$|\mathcal{L}_n| - (i - 1) \frac{|\mathcal{L}_n| - A_n}{A_n} - 1$ conflict tests and across all iterations and there must be a total of

$$C_n = \sum_{i=1}^{A_n} |\mathcal{L}_n| - (i - 1) \frac{|\mathcal{L}_n| - A_n}{A_n} - 1$$

conflict tests. This expression reduces to $C_n = \frac{1}{2}(A_n^2 + A_n(\mathcal{L}_n - 3) + \mathcal{L}_n)$. By plugging in the experimental approximation of A_n this becomes

$$C_n = 5.06866 * 10^{-4} n^4 + 0.685517 n^3 + 4.60191 n^2 + 8.71847 n + 6.01757$$

We now derive an approximation of the number of floating point operations per call to `conflicts`, F_n . We do an average case analysis to find F_n . To compute F_n we note that when evaluating a compound boolean expression in CUDA C short-circuiting is used. So, for example, if `!share_endpoint(I, J)` evaluates to false `conflicts` returns false immediately and no further computation occurs. To compute F_n we first compute the cost of all functions called by `conflicts`. We start with `is_equal`.

Consider the execution tree of `is_equal` in Figure 1. If we assume that each branch is equally likely then the average cost in terms of operations is $1 * .5 + 2 * .25 + 2 * .25 = 1.5$ operations. With this we can compute the average cost of `!share_endpoint(I, J)`, which is 1 plus the average cost of `share_endpoint(I, J)`.

Before we consider the execution tree of `share_endpoint(I, J)` we must consider the probability that two points are equal. Notice that each point is paired with $n - 1$ points and in total there are $\frac{n(n-1)}{2}$ lines. It follows that the probability that two randomly chosen points are equal is

$$\frac{n - 1}{\frac{n(n-1)}{2}} = \frac{2}{n}$$

Thus the probability that `is_equal` returns true is $\frac{2}{n}$. Consider the execution tree of `share_endpoint` in Figure 2.

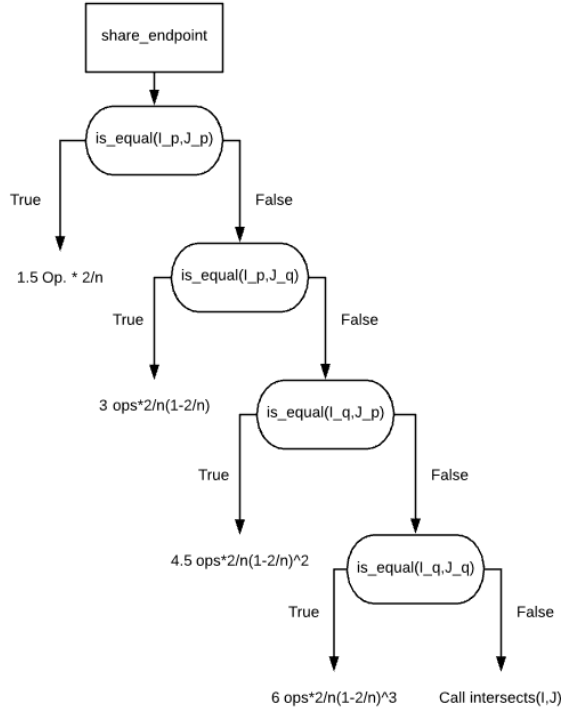


Fig. 2. The execution tree for share_endpoint.

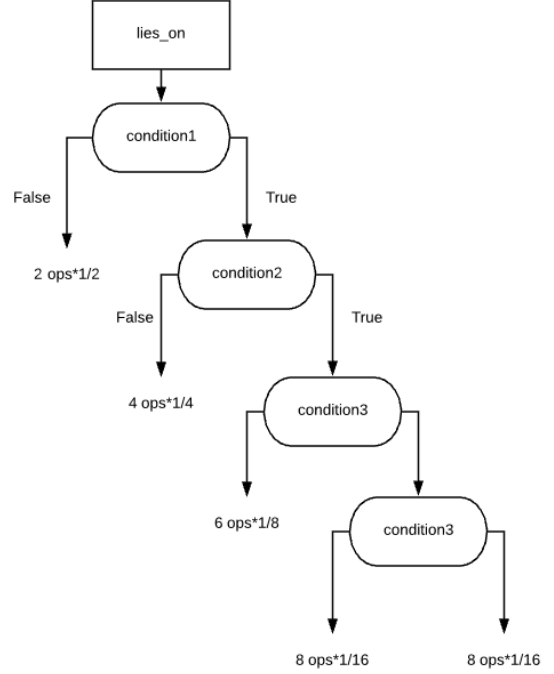


Fig. 3. The execution tree for lies_on.

To compute the average cost of share_endpoint we take the the sum of each branch in Figure 2, this sum reduces to $6 - \frac{18}{n} + \frac{24}{n^2} + \frac{12}{n^3}$.

Before analyzing intersects we analyze orient and lies_on. Each call to orient costs 9 floating point operations. For the cost of lies_on consider its execution tree in Figure 3. For simplicity we label each condition as condition 1, condition 2, etc. and assume that either truth value of each condition is equally likely. From the tree we can see that the average cost of lies_on is $\frac{15}{4}$.

We can now compute the average cost of intersects and conflicts. Since conflicts calls share_endpoint we must include the cost of it and the probability of taking its right most branch into account for each branch of the execution tree of intersects. We note that the execution tree of intersects is a sub tree of the execution tree of conflicts so we include them both in the same figure. Consider the execution tree in Figure 4. From the tree we can see that the total cost of conflicts is

$$\begin{aligned}
 F_n &= 7 - \frac{18}{n} + \frac{24}{n^2} - \frac{12}{n^3} + (1 - \frac{2}{n})^4 (\frac{46}{2} + \frac{51.75}{4} + \frac{57.5}{8} \\
 &\quad + \frac{63.25}{16} + \frac{68}{16}) \\
 &= 58.3281 + \frac{821.25}{n^4} - \frac{1654.5}{n^3} + \frac{1255.88}{n^2} - \frac{428.625}{n}
 \end{aligned}$$

We emphasize that F_n is the cost of an average call to conflicts, so summing F_n over many calls to conflicts yields more accurate estimates.

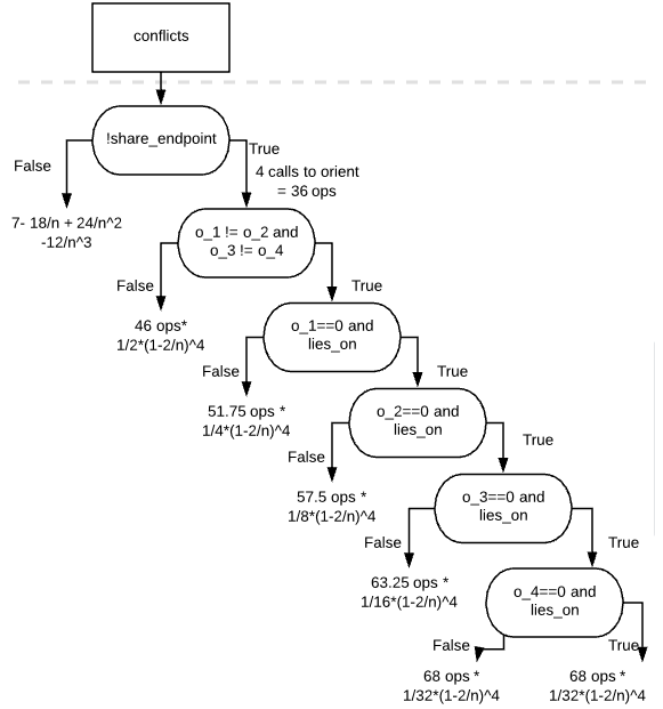


Fig. 4. The execution tree for conflicts.

Since we have approximations for F_n and C_n we can now compute N_n , the total number of operations needed to complete the triangulation step for n points. Notice that N_n is the product of C_n and F_n plus the total number of operations needed to keep track of the index of the smallest line. The latter cost is the total cost of all calls to `find_new_small`, which is $3|\mathcal{L}_n| - 2A_n$. Therefore the total cost of the triangulation step is

$$\begin{aligned} N_n &= C_n F_n + 3|\mathcal{L}_n| - 2A_n \\ &= 0.0295645n^4 + 39.7676n^3 - 24.7725n^2 - 603.873n \\ &\quad + 1259.71 + \frac{1319.19}{n} - \frac{3088.04}{n^2} - \frac{2796.03}{n^3} + \frac{4941.93}{n^4} \\ &\quad + 3|\mathcal{L}_n| - 2A_n \\ &= 0.0295645n^4 + 39.7676n^3 - 23.2765n^2 - 610.817n \\ &\quad + 1249.15 + \frac{1319.19}{n} - \frac{3088.04}{n^2} - \frac{2796.03}{n^3} + \frac{4941.93}{n^4} \end{aligned}$$

With this result, given the time taken to complete the triangulation step, we can estimate the FLOPS of our parallel version. We note that the same computation takes place in the serial version so we may also estimate the FLOPS of the serial version using N_n as a baseline for comparison. In the following section we present experimental results of our parallel version and in Section VII we give a discussion of our results.

VI. EXPERIMENTS

We now present our experimental results. The following data was gathered by executing the serial version of the algorithm on the bg9 computer in the Computer Science department of William and Mary. The values for N^n are reported in giga-operations and the FLOPS are reported as GFLOPS. The function N^n is the approximation of floating point operations required by the algorithm. It was derived in the previous section. This data is used as a baseline for our parallel implementation.

Points	N^n	Run time (seconds)	GFLOPS
724	23.2	7.9712	2.91
1024	75.1816	24.996	3.01
1448	250.656	98.2896	2.55
2048	861.606	385.6797	2.23
2896	3045.22	1407.0545	2.16

We now present the statistics for our parallel CUDA-C implementation. The following data was gathered by executing our parallel version on the GPU in the bg9 computer in the Computer Science department of William and Mary. We only present a subset of the data using the best combination of the number of blocks (B) and number of threads per block (T) for each point set. To see all of the data generated by our tests please refer to the Github repository⁴. As with the values for the serial version N^n is reported in giga-operations.

⁴Specifically the file `log.csv` contains all experiments performed using our best parallel version. `log.csv` was generated using `experiments.sh` which is also included in the repository.

n	B	T	N^n	Time (s)	GFLOPS
724	256	512	23.2	0.2794	83.0
1024	256	1024	75.18	0.4629	162.4
1448	256	1024	250.7	0.9115	275.0
2048	256	1024	861.61	2.2405	384.6
2896	256	1024	3045.2	6.2929	484.0
4096	512	1024	11054	19.5484	565.5
5792	512	1024	40999	66.6196	615.4
8192	512	1024	155008	240.101	645.6

We were able to triangulate point sets with much larger values due to the massive speedup. In the serial version we did not triangulate point sets this large due to the large amount of time it would take.

In the following section we give a discussion of our results.

VII. RESULTS

We now discuss our results. In the previous section we presented data from our experiments. We can see that the serial version scales poorly in terms of GFLOPS, in fact as the number of points in the point set grows, the theoretical GFLOPS actually decreases. We believe that this is due to the effects of moving the processed lines in and out of the cache. While this occurs the CPU is idle, thus lowering the amount of FLOPS. As the number of lines grows thrashing in the cache grows too. We believe that this explains the diminishing performance of the serial version as the number of points grows.

Now consider our parallel version. As the number of lines in the point set increases we can see that the theoretical GFLOPS actually increases, though it appears increase asymptotically. We believe that this increase in GFLOPS is due to the effects of memory coalescing. The way the lines are accessed by threads across blocks in the parallel version is designed so that cache evictions are minimized. Each block of threads only accesses lines that are near each other in the line array, so the lines stay in the cache of the streaming multiprocessors longer, thus minimizing evictions. We believe that this explains the increasing performance of the parallel version.

The parallel version enables larger point sets to be triangulated in a practical amount of time. Consider the serial version, as the number of points increases by a factor of $\sqrt{2}$ the run time roughly increases by a factor of 4. This means that triangulating a point set with $8192 = 2896 * \sqrt{2} * \sqrt{2} * \sqrt{2}$ points should take roughly $1407 * 4 * 4 * 4 = 90,048$ seconds, which is about 25 hours. The parallel version with 512 blocks and 1024 blocks per thread only takes about 240 seconds, or 4 minutes.

We now compare our results to the previous MPI version in [1]. We do not perform a comparison using our theoretical analysis because the MPI implementation is very different than this version, so N^n is not an accurate approximation of the number of floating point operations needed to perform the triangulation step. Here we only compare the run times of the triangulation step of the two versions. Below is a table containing the amount of time needed to accomplish the triangulation step using the MPI version. The stats are

from [1] and were run on the cluster at the Computer Science department at James Madison University in April 2019.

Points	Time on 64 P in (s)
707	1.2657
1000	2.3979
1414	5.3011
2000	12.742
2828	46.203
4000	177.73

We can see that with point sets of comparable size, the GPU version is faster, thus enabling larger point sets to be triangulated. For example, our version is able to compute a point set with 4096 elements in about 20 seconds, where as the MPI version takes about 3 minutes. This is a speedup by a factor of 9.

We can conclude that this implementation allows larger point sets to be triangulated more quickly and efficiently than before. In the following section we provide our conclusions and afterwards we discuss future work.

VIII. CONCLUSION

The CUDA-C implementation presented here scales nicely in terms of both run time and theoretical GFLOPS. Our results and theoretical analysis show that our implementation can triangulate large point sets quickly and efficiently. When compared to the serial version our results show that our implementation is significantly faster, allowing a point set of 8192 points to be triangulated in 4 minutes instead of 25 hours. When compared to the previous work in [1] our implementation is also quicker, being able to triangulate point sets up to nine times faster. We can conclude that them implementation presented here is robust and enables greedy triangulations of large point sets to be carried out in a reasonable amount of time.

IX. FUTURE WORK

We believe that this implementation can be further improved. One way it can be improved is by inlining all the functions discussed in the appendix so the code for them does not have to be moved in and out of the cache. We also believe that this algorithm can be improved by producing a hybrid of this version and the MPI version. This future hybrid implementation should run on multiple nodes, each containing several streaming multiprocessors. The inter-node communications should be handled as in the MPI version and the intra-node processing should occur as in this work.

We also believe that the theoretical approximation of floating point operations can be improved to be more accurate. We believe that it would be beneficial to do a study to test the assumption that an equal number of lines are eliminated during each round of eliminations. If this assumption is false, which is likely, a function $g_{n,r}$ should be derived using experimental data that approximates the number of lines eliminated at round r for a point set of size n . We also think that experimental data should be used to compute the probability that certain paths

are taken in the execution trees. This can be used to derive a more accurate N_n .

REFERENCES

- [1] E. Shoemaker and R. Shoemaker, "Parallelizing the greedy triangulation." submitted to James Madison University Research Journal (JMURJ).
- [2] W. Mulzer and G. Rote, "Minimum-weight triangulation is np-hard," *Journal of the ACM (JACM)*, vol. 55, no. 2, p. 11, 2008.

X. APPENDIX

Here we include the pseudo code for all device functions needed to check if two lines conflict with one another. We note that in these functions zero corresponds to false and a non-zero number corresponds to True. Please refer to the Github repository for implementations of these functions.

The function `conflicts` takes two lines and returns true iff the two lines conflict with one another.

```
conflicts(I, J)
    return !share_endpoint(I, J)
        and intersects(I, J)
```

The function `share_endpoint` takes two lines and returns true iff the two lines share an endpoint. This function is important because if two lines share an endpoint they intersect but they do not conflict with one another. We let I_p and I_q denote the endpoints of the line I .

```
share_endpoint(I, J)
    return is_equal(I_p, J_p)
        or is_equal(I_p, J_q)
        or is_equal(I_q, J_p)
        or is_equal(I_q, J_q)
```

The function `is_equal` takes two points and returns true if they are the same point. We let p_x and p_y denote the x and y coordinates respectively of the point p .

```
is_equal(p, q)
    return p_x == q_x and p_y == q_y
```

The function `intersects` takes two lines and tells us if they intersect.

```
intersects(I, J)
    o1 ← orient(I_p, I_q, J_p)
    o2 ← orient(I_p, I_q, J_q)
    o3 ← orient(J_p, J_q, I_p)
    o4 ← orient(J_p, J_q, I_q)
    return (o1 != o2 and o3 != o4)
        or (o1 == 0 and lies_on(I_p, J_p, I_q))
        or (o2 == 0 and lies_on(I_p, J_q, I_q))
        or (o3 == 0 and lies_on(J_p, I_p, I_q))
        or (o4 == 0 and lies_on(J_p, I_q, J_q))
```

The function `orient` takes three points p, q, r and the return value depends on the orientation of the traversal from p to q to r . If the traversal is clockwise 1 is returned, if the traversal is counter-clockwise -1 is returned, if the points are colinear then 0 is returned. The function `sign` simply returns the sign of its argument.

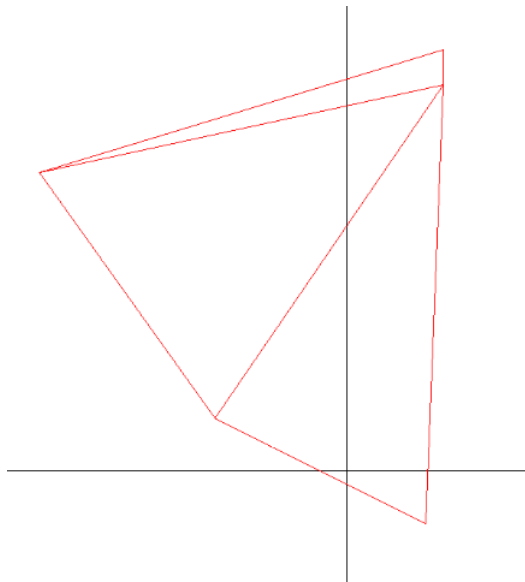


Fig. 5. A point set with 5 points

```
orient(p, q, r)
    return sign((qy - py)(rx - qx)
               - (qx - px)(ry - qy))
```

The function `lies_on` takes three points p, q, r and returns true iff q lies on the line segment formed by p and r . The functions `max` and `min` return the max and min respectively of their arguments.

```
lies_on(p, q, r)
    return qx ≤ max(px, rx) and qx ≥ min(px, rx)
           and qy ≤ max(py, ry) and qy ≥ min(py, ry)
```

Below are some images of various triangulations produced by our parallel version, we include them because they are the visually appealing product of this work.

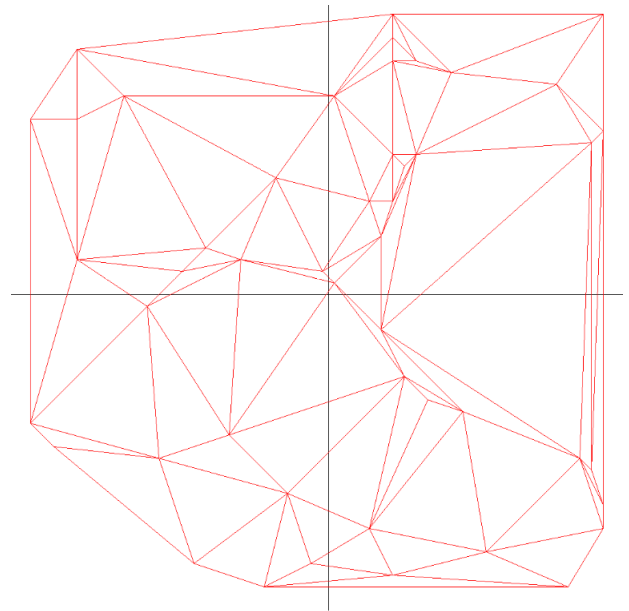


Fig. 6. A point set with 50 points

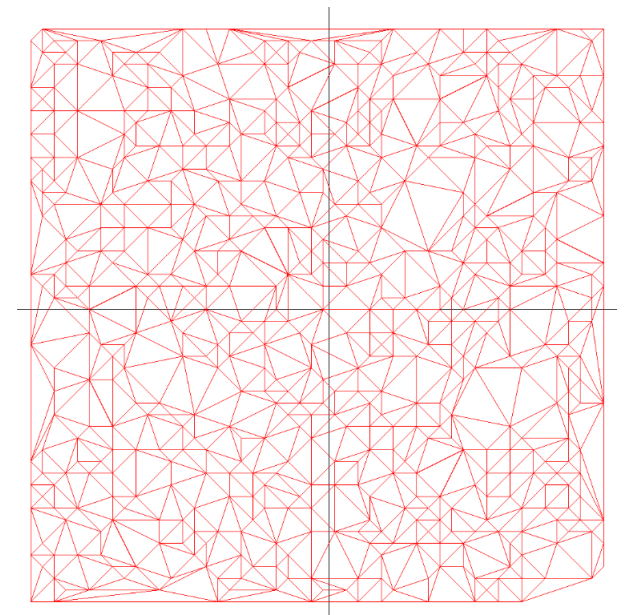


Fig. 7. A point set with 707 points

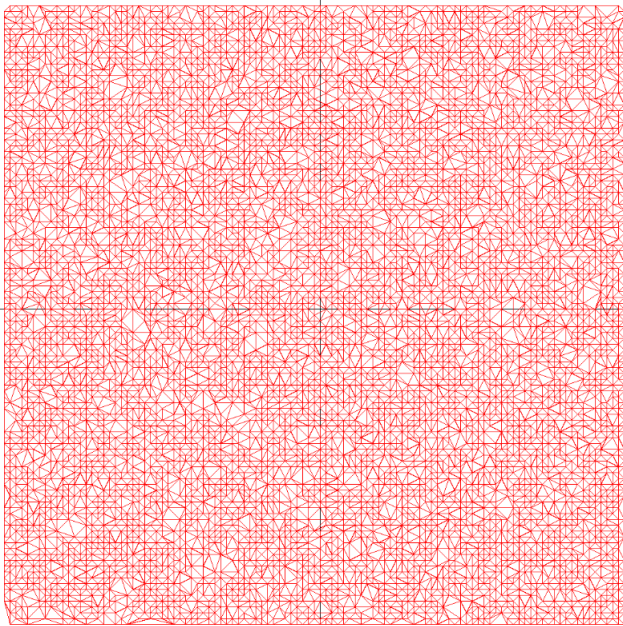


Fig. 8. A point set with 11585 points