

# VERTEXWISE TRIANGLE COUNTING

The topic of the project is triangle counting of large unweighted and undirected graphs, by calculating a vector  $c_3$  that express the number of triangles that every node takes part in. In graph theory, the triangle graph is the complete graph  $K_3$ , consisting of three vertices and three edges. Triangle counting has gained increased popularity in the fields of network and graph analysis. The application of triangle detection, location, and counting are multi-fold: detection of minimal cycles, graph-theoretic clustering techniques, recognition of median, claw-free, and line graphs, and test of automorphism. [Source code](#) consist of two different algorithms ( $v_3$  and  $v_4$ ) of calculating  $c_3$  vector. For every version, there is a sequential implementation and some parallel ones, that use *Cilk*, *OpenMP* and *PThreads*. The parallel implementations were tested in the *AUTH High Performance Computing (HPC) infrastructure* and data of operation run time were exported. Afterwards, we used these data to analyze the behavior of our code.

## WHAT IS PARALLEL COMPUTING?

Parallel computing is a type of computation where many calculations or the execution of processes are carried out simultaneously. Large problems can often be divided into smaller ones, which can then be solved at the same time. Parallel computing can do great work in reducing of operation time of large problems.

## CRITICAL PARTS OF OUR CODE

At this point, we are going to explain some critical functions and points of our code.

- **find\_triangles (v3)**: This function is used to calculate  $c_3$  vector. We have already read the symmetric COO format of the graph and converted it to CSC. So, through this function we take every node  $i$  (from 0 to  $n-1$ ) and find all the pairs  $(j, k)$  of its neighbors (=nodes that are connected with) and then, by calling *check\_edge* function, we find out if these two are connected with an edge to each other too. If this happens, we have found a triangle, so we increase  $c_3[i]$  by 1. We do not increase  $c_3[j]$  and  $c_3[k]$  in this iteration to avoid data race in parallel implementations. Note that we will find the same triangle three times with all possible nodes as first, so we do not loose any triangle.
- **check\_edge (v3)**: As already mentioned, it is used inside *find\_triangles* to check if there is an edge between a pair of nodes. It uses binary search to make problems with high amount of data more efficient ( $O(\log n)$  steps).
- **readmtxvalues, openmtxfile (v3 & v4)**: Used to read *mtx* file and convert it to two COO arrays.
- **find\_triangles (v4)**: In this version, the function is different. Here, we take every pair of nodes that are connected to each other (masking with adjacent  $A$ ) and we are counting their common neighbors (array  $A^*A$ ) with *common\_neighbors* function. The result is the number of triangles that they both take part in. We only increase  $c_3[i]$  though, for the same reason we explained earlier.
- **common\_neighbors (v4)**: Counts the common neighbors of two given nodes in a very efficient way, taking advantage of the fact that the two lists are sorted. We implement a method similar to “merging of two sorted lists” and we replace a  $(length_1 * length_2)$  complexity algorithm (checking all combinations) with an  $O(length_1 + length_2)$  algorithm.
- **Cilk implementation (v3 & v4)**: We chose to use only one “*cilk\_for*”, as we suppose that  $n \gg num\_of\_threads$ . In this way, each thread takes over one iteration (of  $n$  total) at a time and when it finishes with that, it takes the next available from the queue, working dynamically. So we utilize all of our threads every moment and we maximize the speedup.

- **OpenMP implementation (v3 & v4):** Same implementation with *Cilk*.
- **PThreads implementation (v4):** Here, we implemented dynamic scheduling with *PThreads*. Each thread, locks the global variable “iteration” to avoid data race that would cause multiple computations of the same iteration. The thread keeps the value of iteration using a local variable “*t\_it*” and increase the global “iteration” by one. So, this thread computes a specific iteration, while the next thread that will be available will take over the next iteration of the virtual queue in the same way.

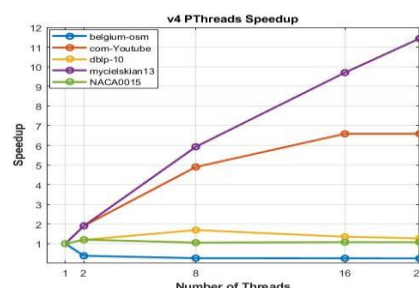
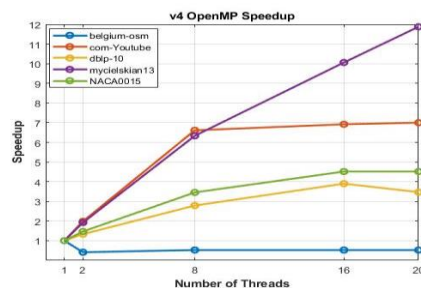
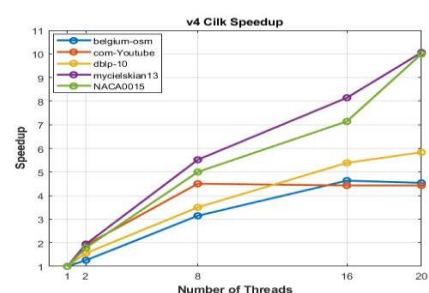
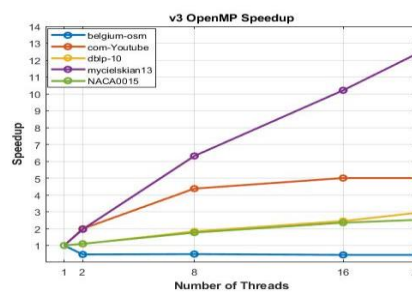
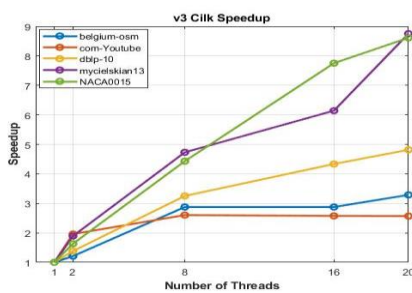
## ANALYZING THE BEHAVIOR OF OUR CODE

As we have mentioned, we also used the *AUTH High Performance Computing (HPC) infrastructure* to test our code behavior for specific graphs given as input. We present our graphs through the following table:

graph	n (nodes)	m (edges)	m/n (density)	m×n (complexity)	triangles
belgium_osm	1,441,295	1,549,970	1.07	$2.23 \times 10^{12}$	2420
com-Youtube	1,134,890	2,987,624	2.63	$3.39 \times 10^{12}$	3,056,386
dblp-2010	326,186	807,700	2.47	$2.63 \times 10^{11}$	1,676,652
mycielskian13	6,143	613,871	99.9	$3.77 \times 10^9$	0
NACA0015	1,039,183	3,114,818	3	$3.23 \times 10^{12}$	2,075,635

We notice that the size of the problem is determined of  $n$  and  $m$ , but we cannot sort the graphs by size, as the exact dependence is not known, we could, though, compare the graphs according to their density and complexity, as we defined them. So, we could say that *com-Youtube* is the most complex one, followed by *NACA0015* and *belgium\_osm*, while *mycielskian13* is extremely dense and *belgium\_osm* extremely sparse. It is really remarkable the fact that *mycielskian13* does not have any triangles despite its density. *Dblp-2010* is a balanced graph in both criteria.

In our experiment in *HPC*, we used every single version of our code, for every single graph and with different amount of threads each time (1, 2, 8, 16, 20). So, we recorded the time of every operation and created the speedup diagrams that are presented below. We calculated speedup through the formula:  $speedup(n\_threads) = operation\_time(1\_thread) / operation\_time(n\_threads)$  and we are going to analyze the scalability of our code.



**Speedup definition:**

In computer architecture, speedup is a number that measures the relative performance of two systems processing the same problem. More technically, it is the improvement in speed of execution of a task executed on two similar architectures with different resources.

**Scalability definition:**

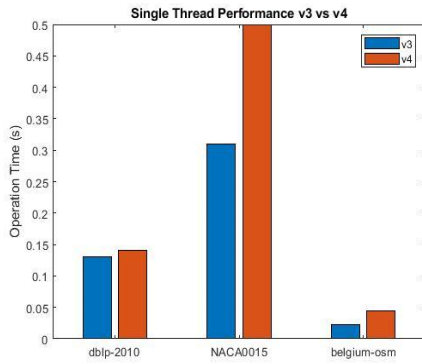
The scalability of a parallel algorithm on a parallel architecture is a measure of its capacity to effectively utilize an increasing number of processors. ... For a fixed problem size, it may be used to determine the optimal number of processors to be used and the maximum possible speedup that can be obtained.

**Amdahl's law:**

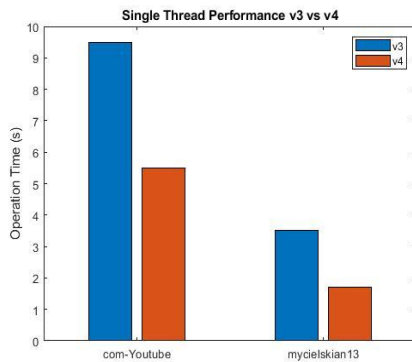
In computer architecture, Amdahl's law is a formula which gives the theoretical speedup in latency of the execution of a task at fixed workload that can be expected of a system whose resources are improved. Amdahl's law gives an *upper bound* of speedup of every parallel implementation, as there is a part of every problem that has to be computed sequentially.

Amdahl's law is confirmed through our diagrams. We can notice that almost every process meets an upper bound on its speedup in the range of 20 cpu threads. Also, according to Amdahl's law, the processes that have not met the upper bound in the 20 thread experiment, are going to do for some wider range of threads.

Moreover, it is clear that the parallel implementation is really efficient in most of the cases. *Mycielskian13* shows the greatest speedup in all of the methods. Its high density could be factor of that, as every thread has a huge amount of work to do in each iteration and the work is almost equally distributed. On the other hand, *belgium\_osm*, who is a complex and sparse graph, shows very bad parallel behavior in *OpenMP* and *PThreads*, as the distribution of work is not achieved properly. However, *Cilk* looks to make much better scheduling with *belgium\_osm*. *Db1p-10*, who is an average graph regarding to our criteria, density and complexity, displays an average behavior regarding to its speedup too, in comparison to the other graphs. Overall, v4 looks to present better speedup performance in comparison to v3. Last but not least, our attempt to implement dynamic scheduling in *PThreads* looks to perform great in the 'dense' *mycielskian13* and the 'complex' *com-Youtube*, but not that good on the other graphs (almost equally to single thread). We have to keep in mind that our *PThreads* implementation is the only one that uses mutex, so it is clear that sometimes the part that is serialized could harm our scalability, according to *Amdahl's law*. The process that is computing during the *lock* is done sequentially, as just one thread can access that memory part at a time, so the part of the work that has to be done sequentially is getting bigger, and the upper bound lower.



At last, we are about to compare  $v3$  and  $v4$  performance in single thread processes. We can notice that for the graphs with low computation time,  $v3$  is quite faster, while for those who demand more time,  $v4$  clearly does better job. Also, in combination with the previous diagrams, we can export that the more the single thread operation time increases, the more the scalability of speedup using more threads increases. *Com-Youtube* and *mycielskian13* are clearly the most time-demanding graphs and they also have the best overall scalability on the first diagrams, in contrast with *belgium\_osm* that has the lowest computation time and the worst overall parallel behavior.



## SUMMING UP

Parallel computing is a great way to improve behavior of our code and solve very complex problems much faster than we would do sequentially. However, not every problem is proper to be computed in parallel and we have to be very careful to avoid inappropriate results.

GITHUB URL OF OUR PROJECT REPOSITORY:

<https://github.com/georgegito/vertexwise-triangle-counting>