

KTRUSSEXPLORER: Exploring the Design Space of K-truss Decomposition Optimizations on GPUs

Safaa Diab*, Mhd Ghaith Olabi*, Izzat El Hajj
American University of Beirut
{syd04, moo02}@mail.aub.edu, izzat.elhajj@aub.edu.lb

*Equal contribution

Abstract—K-truss decomposition is an important method in graph analytics for finding cohesive subgraphs in a graph. Various works have accelerated k-truss decomposition on GPUs and have proposed different optimizations while doing so. The combinations of these optimizations form a large design space. However, most GPU implementations focus on a specific combination or set of combinations in this space.

This paper surveys the optimizations applied to k-truss decomposition on GPUs, and presents KTRUSSEXPLORER, a framework for exploring the design space formed by the combinations of these optimizations. Our evaluation shows that the best combination highly depends on the graph of choice, and analyses the conditions that make each optimization attractive. Some of the best combinations we find outperform previous Graph Challenge champions on many large graphs.

I. INTRODUCTION

A k-truss is a subgraph of a graph such that every edge in the subgraph participates in at least $k - 2$ triangles in the subgraph. K-truss decomposition discovers the maximal k-trusses in a graph for $k \geq 2$. Samsi et al. [11] describe the problem in detail. A popular algorithm for finding a k-truss is to count the triangles that each edge participates in (i.e., an edge’s *support*), delete the edges with support less than $k - 2$ (i.e., *weak* edges), and iteratively repeat this process until no edges are deleted. Triangles are typically found by intersecting the adjacency lists of an edge’s endpoints.

Various works have accelerated k-truss decomposition on GPUs and have proposed different optimizations while doing so. These optimizations include using a directed graph [4], [9], directing edges by degree [4], [9], tiling the adjacency matrix [9], [13], parallelizing list intersection operations [4], [9], removing deleted edges from the graph data structure in between iterations [4], [5], [6], only recounting the support of edges affected by deletions [1], [4], [6], [8], [10], and others. The combinations of these optimizations form a large design space. However, most GPU implementations focus on a specific combination or set of combinations in this space.

To address this issue, we present KTRUSSEXPLORER, a framework for exploring the design space of k-truss decomposition optimizations on GPUs. KTRUSSEXPLORER consists of multiple configurable kernel implementations. The framework is highly parameterized, allowing users to specify any combination of many of the optimizations mentioned previously. Our evaluation shows that the best optimization combination highly

depends on the graph of choice. We also perform a quantitative analysis of the conditions that make each optimization attractive. Some of the best combinations we find outperform previous Graph Challenge champions on many large graphs. KTRUSSEXPLORER has been open-sourced to help further advance research on k-truss decomposition optimizations.

II. KTRUSSEXPLORER

This section surveys the literature on optimizations for k-truss decomposition on GPUs while describing how these optimizations are supported in KTRUSSEXPLORER.

A. Edge-Centric and Vertex-Centric Parallelization

One key distinction between parallel implementations is whether they are edge-centric [1], [6], [10], [12] or vertex-centric [3], [4]. Edge-centric implementations assign a thread (or group of threads) to each edge to find the support of that edge. Vertex-centric implementations assign a thread (or group of threads) to each vertex to find the support of that vertex’s outgoing edges.

The edge-centric approach makes load balancing easier than the vertex-centric approach where different vertices may have a substantially different number of outgoing edges whose support needs to be computed. On the other hand, the vertex-centric approach works with a CSR representation of the graph whereas the edge-centric approach typically uses more memory to store a COO+CSR representation to look up the endpoints of each edge. Moreover, the vertex-centric approach enables optimizations that target edges that share the same source vertex, such as laying out the source vertex’s adjacency list as a bitmap for all its outgoing edges to share [3], [4].

Linear algebraic formulations [5] can also be classified as edge-centric or vertex-centric depending on how the sparse matrix multiplication is implemented. Assigning threads to nonzeros corresponds to an edge-centric implementation, whereas assigning threads to rows corresponds to a vertex-centric implementation.

KTRUSSEXPLORER currently only implements the edge-centric parallelization approach.

B. Graph Directedness

Another key distinction between implementations is the use of undirected graphs [1], [8] or directed graphs [3], [4], [9]. In

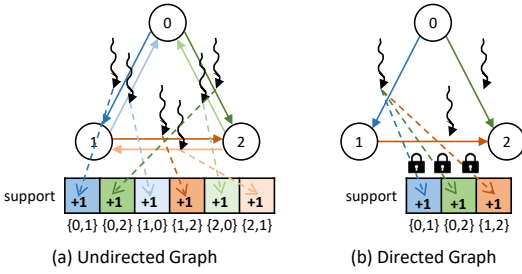


Fig. 1: Discovering Triangles Based on Graph Directedness

linear algebraic terms, using directed graphs is equivalent to operating on triangular matrices [5]. The distinction between using undirected and directed graphs is illustrated in Fig. 1. With an undirected graph, each edge’s thread discovers the triangles that the edge participates in independently. However, each triangle is redundantly discovered six times as shown in Fig. 1(a). With a directed graph, each triangle is discovered only once. For example, in Fig. 1(b), only edge $\{0, 1\}$ ’s endpoints have a common neighbor so only that edge’s thread will discover a triangle. However, although redundancy is reduced, the thread that discovers the triangle must update the support of all edges involved which requires atomic operations.

Using a directed graph has many advantages. First, fewer edges are stored which reduces memory capacity and bandwidth requirements. Second, triangles are not discovered redundantly which reduces the amount of work. Third, adjacency lists are shorter which makes intersecting them faster.

On the other hand, using undirected graphs also has advantages owing to the fact that each edge’s thread finds the edge’s support independently. First, no atomic operations are needed to update support values. Second, a thread may stop an intersection operation earlier if certain conditions are met [12]. One such condition is when one of the two adjacency lists is smaller than $k - 2$. In this case, it is impossible for the thread’s edge to participate in $k - 2$ triangles so there is no need to attempt the intersection. Another such condition when $k - 2$ triangles are found but the intersection operation is not yet complete. In this case, the thread does not need to continue the operation because it has already established that the edge is not weak. These early stopping conditions cannot be applied with directed graphs because a thread is also responsible for updating the support of other edges, not just the one it owns.

Still, undirected graphs suffer from the redundancy of discovering triangles six times. As an optimization, some works use undirected graphs but only count triangles for one of the edge directions [6], [10], [12] which reduces the redundancy from six to three. In linear algebraic terms, this optimization is equivalent to using the full adjacency matrix as an input, but only computing a triangular matrix as an output, taking advantage of the fact that the output is symmetric. The drawback of this approach is that if an edge has weak support and needs to be deleted, finding the reverse edge to delete it as well can be expensive.

KTRUSSEXPLORER provides an option for selecting between directed and undirected graphs. For undirected graphs, both conditions for stopping the intersection operation early are implemented, and triangles are counted for both directions of each edge (i.e., each triangle is discovered six times).

C. Directing Edges by Degree

A simple approach for converting an undirected graph to a directed one is to keep only the edges from the vertex with the lower index to the vertex with the higher index, i.e., direct edges by index. Another approach is to keep the edges from the vertex with lower degree to the vertex with higher degree, i.e., direct edges by degree. The latter approach is good for keeping adjacency lists short. With undirected graphs, high-degree vertices have long adjacency lists which are expensive to intersect. Directing edges by degree significantly shrinks the adjacency lists of high-degree vertices.

Directing edges by degree has been done in two ways. One way is to simply remove the edges from higher-degree to lower-degree vertices from the graph [6], [9]. Another way is to sort vertices by increasing degree, relabel vertices according to their sort index, and then direct edges by index [4].

KTRUSSEXPLORER provides an option for directing graphs by index or by degree. Directing by degree is currently implemented by sorting, relabelling, then directing by index.

D. Tiling

Tiling refers to partitioning of the adjacency matrix into tiles. Tiling has been evaluated in the literature for triangle counting on GPUs [9], [13]. KTRUSSEXPLORER applies tiling to k -truss decomposition on GPUs. We discuss tiling at more length than other optimizations because it has received less attention in the literature and its implementation varies across different works.

An example of how tiling is implemented in KTRUSSEXPLORER is shown in Fig. 2. Fig. 2(b) shows the logical adjacency matrix of the graph in Fig. 2(a) and Fig. 2(c) shows how it is physically stored as a CSR data structure (we use a hybrid COO+CSR but omit the COO part in the figure). Fig. 2(d) shows the same adjacency matrix with tiling applied and Fig. 2(e) shows how it is physically stored as a tiled CSR data structure. The matrix is logically divided into 2D square tiles, with each tile’s edges stored contiguously. We refer to the number of tiles in each dimension as the *tiling factor*.

Fig. 3 is used to illustrate the advantages of tiling. In each subfigure, the three matrices represent different logical views of the same physical data structure. The view to the bottom right depicts how threads are assigned to edges. The view to the left depicts how a thread traverses the adjacency list of its edge’s source vertex. The view on top depicts how a thread traverses the adjacency list of its edge’s destination vertex.

The first advantage of tiling is that it improves data locality. In Fig. 3(a), the first four threads in the grid are assigned to the first four edges of the CSR data structure which all have the same source vertex but different destination vertices. Hence, there is high spatial and temporal locality when accessing the

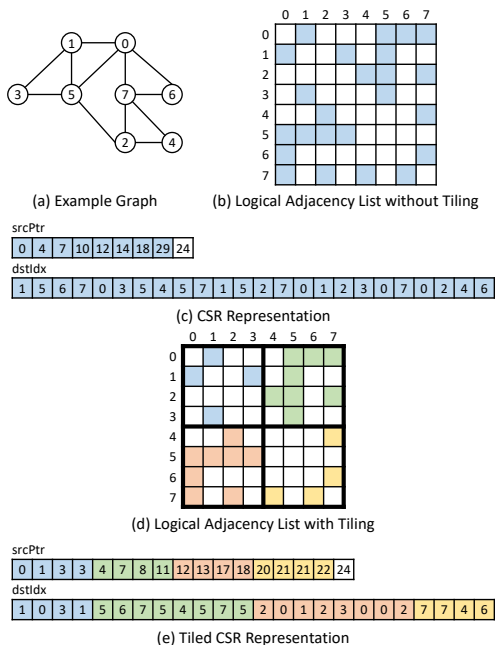


Fig. 2: Example of Tiling an Adjacency Matrix

source vertex adjacency lists, but poor spatial and temporal locality when accessing the destination vertex adjacency lists. In contrast, in Fig. 3(b), the four threads are assigned to the first four edges of a tile. In this case, accesses to both source and destination vertex adjacency lists exhibit good spatial and temporal locality because of tiling.

We also tried loading tiles to shared memory for faster access, however, we found that doing so significantly hurt performance. The reason is that some data in a tile may not be needed so loading that data to shared memory is wasteful. For example, in Fig. 3(b), vertex 2’s adjacency list is not needed by any thread. Loading the entire tile to shared memory would result in vertex 2’s adjacency list being loaded even though it is not needed. For this reason, we do not use shared memory and rely on the L1 cache for fast access to reused data.

The second advantage of tiling is that it partitions long intersection operations into multiple shorter *sub-intersections*. This partitioning makes intersection operations faster because it allows skipping edges in a sublist if the other sublist is empty or has reached the end. For example, in Fig. 3(c), the thread assigned to edge {2, 5} needs multiple steps to intersect the adjacency lists of vertices 2 and 5. In contrast, in Fig. 3(d) with tiling, the intersection is partitioned into two sub-intersections. Sub-intersection 1 takes zero steps because the sublist of vertex 2’s adjacency list is empty. Sub-intersection 2 also takes zero steps because the sublist of vertex 5’s adjacency list is empty. Therefore, the intersection takes fewer steps overall and loads less data from memory. This affect is similar to that achieved by Intersect-Path [7], except that Intersect-Path finds sublists dynamically rather than partitioning the data structure.

KTRUSSEXPLOER provides an option for specifying whether or not to apply tiling and the tiling factor of choice.

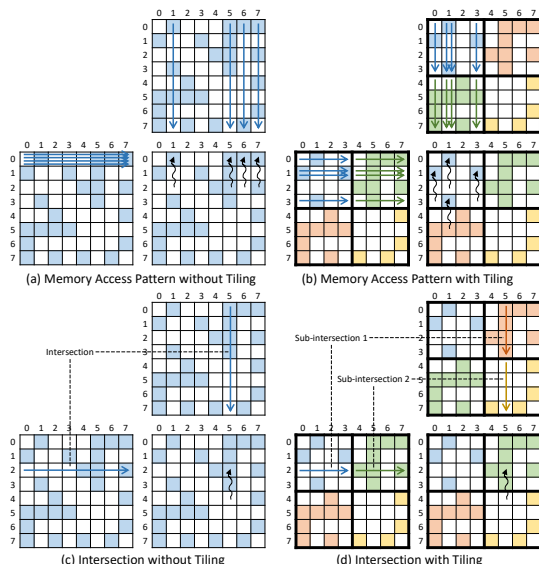


Fig. 3: Impact of Tiling

E. Parallelizing Intersections

Various works parallelize individual intersection operations to further extract parallelism from the computation. This parallelization typically divides one of the lists across multiple threads and has each thread find the corresponding edge(s) in the other list via binary search [9] or a bitmap lookup [3], [4].

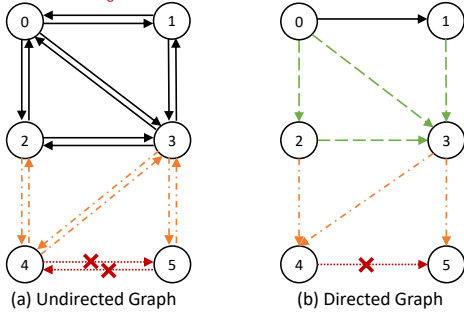
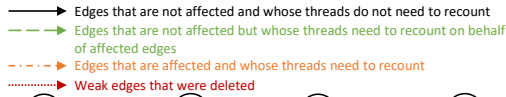
KTRUSSEXPLOER parallelizes intersections in the context of tiling. Recall that with tiling, each intersection operation is partitioned into multiple sub-intersections. These sub-intersections can be performed in parallel. KTRUSSEXPLOER assigns multiple threads to each intersection operation (up to the tiling factor) and divides the sub-intersections across these threads. Each sub-intersection is performed sequentially by one thread. In the absence of tiling, the whole intersection operation is performed sequentially by one thread.

We refer to the number of threads assigned per intersection as the *parallelization factor*. KTRUSSEXPLOER provides an option for experimenting with different parallelization factors.

F. Removing Deleted Edges Intermediately

When a weak edge is deleted during an iteration, it may be marked as deleted and kept in the graph data structure or it may be removed from the data structure entirely before the next iteration. The advantage of removing edges intermediately is that it shrinks adjacency lists making intersections faster, and reduces the graph’s memory footprint. The disadvantage is that removing edges is more expensive than simply marking them as deleted. For this reason, some works do not remove edges in between iterations [9], [10], [12], some remove them every iteration [3], [4], [5], [6], and some remove them only if enough edges have been deleted [1].

Removing deleted edges from the graph data structure can be done via a stream compaction operation that filters out edges marked as deleted. Stream compaction can be done on the entire edge list [1] or on each vertex’s adjacency



01: parallel for $e = \{u, v\} \in E$ do
 02: if e is deleted then
 03: mark u as affected, mark v as affected
 04: parallel for $e = \{u, v\} \in E$ do
 05: if e is not deleted and (u is affected or v is affected) then
 06: mark e as affected
 07: if u is not affected then mark u as needs to recount
 08: else if v is not affected then mark v as needs to recount
 09: parallel for $e = \{u, v\} \in E$ do
 10: if e is not deleted and e is not affected then
 11: if u needs to recount or v needs to recount then
 12: mark e as needs to recount

(c) Pseudocode for Marking Affected Edges

Fig. 4: Recomputing Support Values for Affected Edges

list separately [3], [4], [5]. Compacting the entire edge list is more expensive because the array to be compacted is larger, and because the CSR pointers need to be recomputed after the compaction. On the other hand, compacting each vertex’s adjacency list separately requires more metadata to track where each vertex’s adjacency list ends.

KTRUSSEXPLORER provides an option for specifying for how many initial iterations (if any) edges should be removed. This approach is based on the observation that the largest number of edges is removed in the first few iterations and decreases as iterations proceed. Edges are removed via stream compaction of the entire edge list using the in-place stream compaction primitive in Thrust [2]. Compacting each vertex’s adjacency list separately is not yet supported.

G. Recomputing Support Values

Deleting a weak edge breaks the triangles that the edge participates in which may affect the support of the edges that share a vertex with the deleted edge. These edges are referred to as *affected* edges. When recomputing support values across iterations, it is sufficient to recompute the support of affected edges and unnecessary to recompute for all edges.

The advantage of recomputing the support of only affected edges is that it eliminates unnecessary work. However, the disadvantage is that it incurs additional overhead for identifying which edges are affected. For this reason, some works recompute the support of all undeleted edges [5], [9], [12] while others recompute the support of only the affected edges [1], [3], [4], [6], [8], [10].

Deciding which edges’ threads should recount triangles differs between undirected graphs and directed graphs. This

TABLE I: Design Space Explored

Optimization	Options
Graph Directedness	undirected, directed by index, directed by degree
Tiling	no tiling, tiling with a factor of {2, 4, 8, 16}
Parallelizing Intersections	no parallelization, parallelization with a factor of {2, 4, 8, 16}
Removing Edges Intermediately	no iterations, first {1, 2, 4, 8} iterations, all iterations (edges are also always removed at the end for all combinations)
Recomputing Support Values	all edges, affected edges

distinction is shown in Fig. 4. In the undirected graph in Fig. 4(a), each thread is solely responsible for computing its edge’s support. Therefore, only the affected edges’ threads need to recount triangles. However, for directed graphs, an edge may not be affected, but its thread may be responsible for finding a triangle on behalf of an affected edge and updating its support. Therefore, it is not sufficient for only the affected edges’ threads to recount. Threads assigned to any edge that shares a vertex with an affected edge must also recount. For example, in Fig. 4(b), edge {2,3} is not affected by the deletion of edge {4,5}. However, its thread must recount because it must notify edge {3,4} which is affected that it is part of the triangle {2,3,4}. The pseudocode for marking affected edges is shown in Fig. 4(c).

Another optimization related to recomputing support values is specific to the case when $k = 3$. When $k = 3$, any edge deleted in the first iteration does not belong to a triangle. Hence, no triangles are broken in the first iteration, so no further iterations are needed to recompute support values [3].

KTRUSSEXPLORER provides an option for specifying whether to recompute support values of all undeleted edges or only affected edges for both directed and undirected graphs. If $k = 3$, only one iteration is performed.

III. METHODOLOGY

We evaluate KTRUSSEXPLORER using a Volta V100 GPU with 16GB of device memory coupled with an AMD EPYC 7551P CPU with 15GB of main memory. We evaluate with all the data sets in the graph challenge collection [11] except for Friendster, graph500-scale24-ef16, and graph500-scale25-ef16 due to limited device memory capacity. We report results for $k = 3$, and for $k = k_{max}$ when k_{max} is not 3.

The design space explored is summarized in Table I. For $k = 3$, the space is searched exhaustively for graphs with less than 17 million edges. For larger graphs, only a subset of the combinations are searched based on which combinations did best with the other graphs. For graphs where the best parallelization factor was 16, a parallelization factor of 32 is also attempted. Since $k = 3$ only needs one iteration to converge, edges are removed once at the end and recomputing support values is not relevant.

For each combination, we report the mean of 10 runs after discarding 5 warm up runs. For large graphs, we take the mean of 5 runs with no warm up runs. The time reported includes: counting triangles, marking edges as deleted, data transfer from device to host to check for convergence, removing deleted edges intermediately for relevant combinations, removing deleted edges at the end for all combinations, and marking affected edges for relevant combinations. The time reported

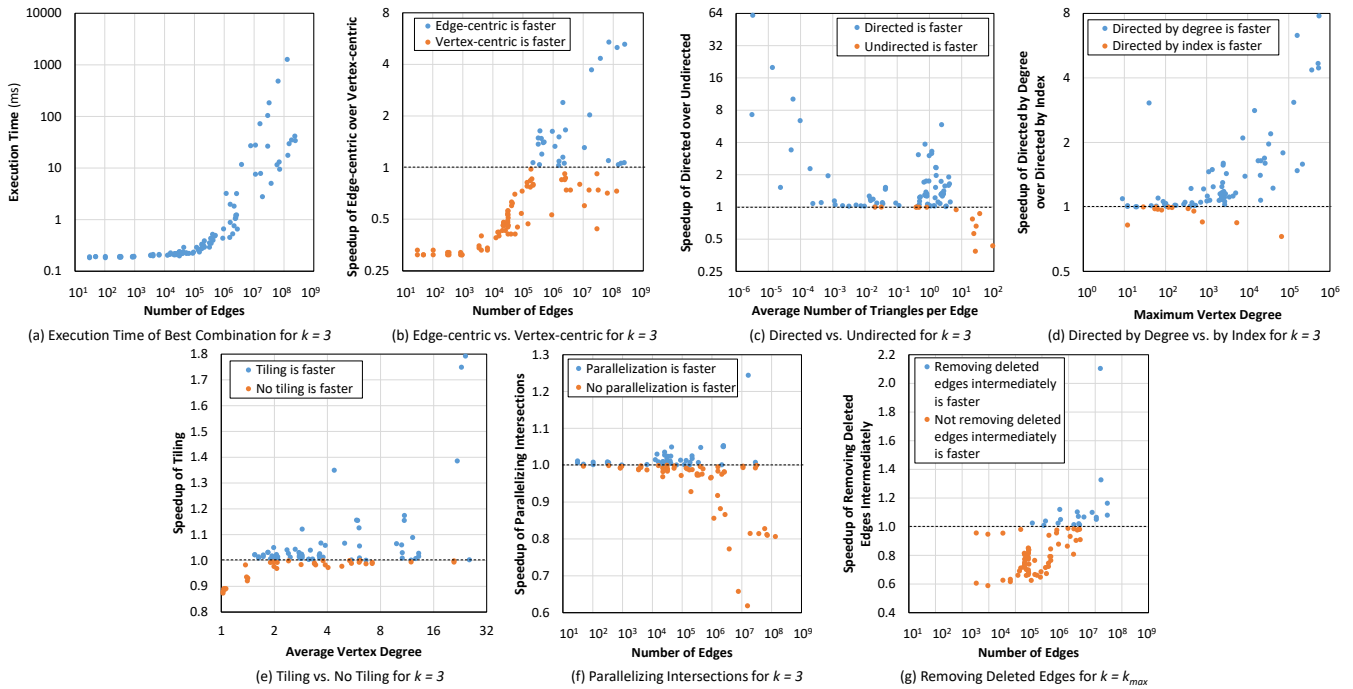


Fig. 5: Evaluation

does not include: allocation and deallocation time, initial and final copy time, and initial and final graph conversion time (COO to CSR, relabelling vertices, undirected to directed).

IV. EVALUATION

Fig. 5(a) shows how the execution time (in *milliseconds*) for the best optimization combination scales with the number of edges for $k = 3$. Details on the best combination found for each graph and the breakdown of the execution time can be found in Table II for $k = 3$ and Table III for $k = k_{max}$. It is clear that there is no best combination for all graphs and that the best combination depends on the graph of choice. In the rest of this section, we quantitatively analyze the conditions that make each design decision or optimization attractive.

A. Edge-centric and Vertex-centric Parallelization

Since we do not currently support a vertex-centric implementation, we compare the edge-centric implementation in KTRUSSEXPLOER to the vertex-centric implementation in Bisson & Fatica [4] using the kernel times reported in that work. Bisson & Fatica [4] are also the 2018 Graph Challenge champions. Fig. 5(b) shows the speedup of the edge-centric implementation over the vertex-centric implementation for $k = 3$. It is clear that the vertex-centric implementation is superior for smaller graphs, but as the graphs get larger, the edge-centric implementation becomes competitive on many graphs (up to $5.4\times$ faster). Recall from Section II-A that the advantage of edge-centric parallelization over vertex-centric parallelization is that edge-centric parallelization makes load balancing easier. The importance of load balance becomes more pronounced as the graphs grow in size.

Note however that this comparison has some limitations because different optimizations are applied to each implementation so there may be other factors impacting the performance difference. For this reason, we plan to support a vertex-centric implementation in KTRUSSEXPLOER as future work in order to have a more controlled comparison.

B. Graph Directedness

Table II shows that most graphs have better performance when a directed graph is used. However, a significant number of graphs perform better when an undirected graph is used. The preference for undirected graphs is independent of graph size as there are very small as well as very large graphs that do better when the graph is undirected. We observe that the preference for undirected graphs is correlated with the average number of triangles per edge in the graph.

Fig. 5(c) shows how the average number of triangles per edge impacts the speedup of the best combination that uses a directed graph over the best combination that uses an undirected graph. It is clear that when the number of triangles per edge becomes large, undirected graphs perform better. Recall from Section II-B that the advantages of undirected graphs are that they do not need atomic operations to update support values, and that intersection operations can be stopped early as soon as the threshold for an edge not being weak is met. When the number of triangles per edge is high, support values will be updated frequently, intersection operations will be long, and the chances of meeting the threshold early will be high. Hence, avoiding atomic operations and stopping intersections early are features that make undirected graphs attractive.

C. Directing Edges by Degree

Table II shows that the overwhelming majority of cases where a directed graph is preferred, directing edges by degree is the faster option. We observe that the extent to which directing edges by degree is better than by index correlates with the maximum vertex degree in the graph.

Fig. 5(d) shows how the maximum vertex degree impacts the speedup of the best combination that directs edges by degree over the best combination that directs edges by index. It is clear that as the maximum vertex degree increases, the benefit of directing edges by degree increases. Recall from Section II-C that the advantage of directing edges by degree is reducing the size of the adjacency lists for high-degree vertices. Hence, directing the graph by degree is more attractive as the degree of high-degree vertices increases.

D. Tiling

Table II shows that some graphs have better performance when tiling is applied while others are better off without tiling. We observe that the preference for tiling is correlated with the average vertex degree in the graph. Fig. 5(e) shows how the average vertex degree impacts the speedup of the best combination that applies tiling over the best combination that uses does not apply tiling. It is clear that as the average vertex degree increases, the benefit of tiling increases. Recall from Section II-D that one of the advantages of tiling is that it partitions long intersection operations into shorter sub-intersections. As the average vertex degree increases, the intersection operations become longer and partitioning them becomes more important. Hence, tiling is more attractive when the average vertex degree is high.

E. Parallelizing Intersections

Table II shows that some graphs have better performance when individual intersection operations are parallelized while others do not. Fig. 5(f) shows how the size of the graph (number of edges) impacts the speedup of the best combination that parallelizes intersections over the best combination that does not. It is clear that the advantage of parallelizing intersections diminishes for large graphs. Recall from Section II-E that the advantage of parallelizing intersection operations is to extract more parallelism from the computation. For large graphs, there is a sufficient amount of parallelism to fully utilize the device because there are many intersections to perform. Hence, extracting more parallelism becomes less attractive as the graph gets larger.

F. Removing Deleted Edges Intermediately

For $k = 3$, only one iteration is needed for convergence so edges are always removed once at the end. Table II shows that a significant fraction of the computation is spent removing deleted edges at the end. This overhead is particularly high for small graphs, where removing deleted edges at the end accounts for up to 80% of the execution time. We plan to reduce this overhead by further optimizing the stream compaction operation, as well as supporting stream compaction on

each vertex's adjacency list separately [3], [4], [5] as opposed to the entire edge list.

For $k = k_{max}$, Table III shows that most graphs have better performance when deleted edges are not removed intermediately, but we expect that reducing the overhead of stream compaction should allow more graphs to show benefit. We observe that the preference for removing deleted edges intermediately is correlated with the number of edges in the graph. Fig. 5(g) shows how the number of edges impacts the speedup of the best combination that removes deleted edges intermediately over the best combination that does not. It is clear that the advantage of removing deleted edges intermediately increases with the number of edges. Recall from Section II-F that the benefit of removing deleted edges is shrinking intersections and reducing the memory footprint, which is more critical as graphs get larger.

G. Recomputing Support Values

For $k = 3$, only one iteration is needed for convergence so recomputing edge support values is not relevant. For $k = k_{max}$, Table III shows that all except a few very large graphs have better performance when the support of all edges is recomputed, not just affected edges. We plan to further optimize the feature of recomputing the support of affected edges. First, we observe that in the initial iterations where many edges are deleted, the number of affected edges is very large. Hence, the overhead of tracking affected edges is not worth the effort it saves. For this reason, we plan to provide the option to track affected edges for only later iterations where the affected edges are few. Second, our current implementation launches threads for all edges and each thread checks if it needs to recount or not. Hence, computational resources are still allocated for threads that do not need to recount and there is high control divergence. Instead, we plan to reduce this divergence by creating a frontier of edges whose threads need to recount and only launching threads for those edges.

V. CONCLUSION

This paper surveys the optimizations applied to k-truss decomposition on GPUs, and presents KTRUSSEXPLORER, a framework for exploring the design space formed by the combinations of these optimizations. The optimizations supported include using a directed graph, directing edges by degree, tiling the adjacency matrix, parallelizing list intersection operations, removing deleted edges from the graph in between iterations, and recomputing support values for only affected edges. Future work includes supporting vertex-centric parallelization, expanding and enhancing the selection of supported optimizations, and leveraging properties of the graph to prune the search space or infer the best combination rather than search the space exhaustively.

ACKNOWLEDGMENTS

We would like to thank Amer Mouawad and Jad Ismail for the valuable discussions we had with them. This work is supported by the University Research Board of the American University of Beirut (URB-AUB-103782-25509).

