Leveraging Multiple GPUs and CPUs for Graphlet Counting in Large Networks

Ryan A. Rossi Palo Alto Research Center rrossi@parc.com

ABSTRACT

Massively parallel architectures such as the GPU are becoming increasingly important due to the recent proliferation of data. In this paper, we propose a key class of hybrid parallel graphlet algorithms that leverages multiple CPUs and GPUs simultaneously for computing k-vertex induced subgraph statistics (called graphlets). In addition to the hybrid multi-core CPU-GPU framework, we also investigate single GPU methods (using multiple cores) and multi-GPU methods that leverage all available GPUs simultaneously for computing induced subgraph statistics. Both methods leverage GPU devices only, whereas the hybrid multi-core CPU-GPU framework leverages all available multi-core CPUs and multiple GPUs for computing graphlets in large networks. Compared to recent approaches, our methods are orders of magnitude faster, while also more cost effective enjoying superior performance per capita and per watt. In particular, the methods are up to 300+ times faster than a recent state-of-the-art method. To the best of our knowledge, this is the first work to leverage multiple CPUs and GPUs simultaneously for computing induced subgraph statistics.

Categories and Subject Descriptors

H.2.8 [Database Applications]: Data Mining; I.2.6 [Artificial Intelligence]: Learning; G.1.0 [Numerical Analysis]: Parallel Algorithms

Keywords

Graphlets; induced subgraphs; network motifs; graphlet decomposition; GPU computing; multi-GPU; heterogeneous computing; parallel algorithms; graph mining; roles; relational learning; graph classification; graph kernels; node and edge embedding

1. INTRODUCTION

Graphlets are small k-vertex induced subgraphs and are important for many predictive and descriptive modeling tasks [21,

CIKM '16 Indianapolis, Indiana USA

© 2016 ACM. ISBN 0-12345-67-8/90/01...\$15.00 DOI: 10.475/123_4 Rong Zhou Palo Alto Research Center rzhou@parc.com

18, 12] in a variety of disciplines including bioinformatics [29, 27], cheminformatics [25], and image processing and computer vision [31, 32]. Given a network G, our approach counts the frequency of each k-vertex induced subgraph patterns (See Table 1). These counts represent powerful features that succinctly characterize the fundamental network structure [27]. Indeed, it has been shown that such features accuratly capture the local network structure in a variety of domains [14, 8, 9]. As opposed to global graph parameters such as diameter for which two or more networks may have global graph parameters that are nearly identical, yet their local structural properties may be significantly different.

Despite the practical importance of graphlets, existing algorithms are slow and are limited to small networks (e.g., even modest graphs can take days to finish), require vast amounts of memory/space-inefficient, are inherently sequential (inefficient and difficult to parallelize), and have many other issues. Overcoming the slow performance and scalability limitations of existing methods is perhaps the most important and challenging problem remaining. This work proposes a fast hybrid parallel algorithm for computing both connected and disconnected subgraph statistics (of size k) including macro-level statistics for the global graph G as well as microlevel statistics for individual graph elements such as edges. Furthermore, a number of important machine learning tasks are likely to benefit from the proposed methods, including graph anomaly detection [20, 5], entity resolution [6], as well as features for improving community detection [26], role discovery [22], and relational classification [10].

The recent rise of Big Data has brought many challenges and opportunities [28, 1, 2]. Recent heterogeneous high performance computing (HPC) architectures offer viable platforms for addressing the computational challenges of mining and learning with big graph data. General-purpose graphics processing units (GPGPUs) [7] are becoming increasingly important with applications in scientific computing, machine learning, and many others [11]. Heterogeneous (hybrid) computing architectures consisting of *multi-core CPUs* and *multiple GPUs* are an attractive alternative to traditional homogeneous HPC clusters (HPCC). Despite the impressive theoretical performance achievable by such hybrid architectures, leveraging this computing power remains an extremely challenging problem.

Graphics processing units (GPUs) offer cost-effective highperformance solutions for computation-intensive data mining applications. As an example, the Nvidia Titan Black GPU has 2880 cores capable of performing over 5 trillion floating-point operations per second (TFLOPS). For comparison a Xeon

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.



Figure 1: The distribution of graphlet computation times for the edge neighborhoods obey a power-law. The time taken to count $k = \{2, 3, 4\}$ graphlets for each edge is shown above and clearly obeys a power-law distribution (tech-internet-as).

E5-2699v3 processor can perform about 0.75 TFLOPS, but can cost 4x as much. Besides TFLOPS, GPUs also enjoy a significant advantage over CPUs in terms of memory bandwidth, which is more important for data-intensive algorithms such as graphlet decomposition. For Titan Black, its maximum memory bandwidth is 336 GB/sec; whereas E5-2699v3's is only 68 GB/sec. Higher memory bandwidth means more graph edges can be traversed on the GPU than the CPU in the same amount of time, which is one of the main reasons why our approach can outperform the state-of-the-art.

However, unlike the single GPU-based approach proposed in [19], we adopt a hybrid parallelization strategy that exploits the complementary features of multiple GPUs and CPUs to achieve significantly better performance. We begin with the key observation that the amount of work required to compute graphlets for each edge in G obeys a power-law (See Figure 1). Strikingly, a handful of edges require a lot of work to determine the local graphlet counts (due to the density and structure of the local edge neighborhood), whereas the vast majority of other edges require only a small amount of work. Such heterogeneity can cause significant load balancing issues, especially for GPUs that are designed to solve problems with uniform workloads (e.g., dense matrix-matrix multiplications). This motivates our hybrid approach that dynamically divides up the work between GPU and CPU, to reduce inter-processor communication, synchronization and data transfer overheads.

This work demonstrates that parallel graphlet methods designed for heterogeneous computing architectures consisting of many multi-core CPUs and multiple GPUs can significantly improve performance over GPU-only and CPU-only approaches. In particular, our hybrid CPU-GPU framework is designed to leverage the advantages and key features offered by each type of processing unit (multi-core CPUs and GPUs). As such, our approach leverages the fact that graphlets can be computed via M independent edge-centric neighborhood computations. Therefore, the method dynamically distributes the edge-centric graphlet computations to either a CPU or a GPU. In particular, the edge-centric graphlet computations that are fundamentally unbalanced and highly-skewed are given to the CPUs whereas the GPUs work on the more well-balanced and regular edge neighborhoods (See Figure 4). These approaches capitalize on the fact that GPUs are generally best for computations that are well-balanced and regular, whereas CPUs are designed for a wide variety of applications and thus more flexible [15, 16]. Our approach also leverages *dynamic load balancing* and *work stealing strategies* to ensure all GPU and CPU workers (cores) remain fully utilized.

2. RELATED WORK

Recently, Milinković *et al.* [19] proposed a GPU algorithm for counting graphlets based on a recent sequential graphlet algorithm called ORCA [13]. However, this paper is fundamentally different. First and foremost, that approach is not hybrid and is only able to use a single GPU for computing graphlets. In addition, that work focuses on computing connected graphlets only, whereas we compute both connected and disconnected induced subgraphs. Moreover, that approach computes graphlets for each vertex in parallel (vertex-centric), whereas our methods are naturally edge-centric and search edge neighborhoods in parallel. Furthermore, that work does not provide any comparison to understand the utility and speedup (if any) offered by their approach.

In this work, we propose a heterogeneous graphlet framework for hybrid multi-GPU and CPU systems that leverages all available GPUs and CPUs for efficient graphlet counting. Our single-GPU, multi-GPU, and hybrid CPU-GPU algorithms are largely inspired by the recent *state-of-the-art* parallel (CPUbased) exact graphlet decomposition algorithm called PGD [4, 3], which is known to be significantly faster and more efficient than other methods including RAGE [17], FANMOD [30], and ORCA [13]. Moreover, PGD has been parallelized for multi-core CPUs and is publicly available¹. Our approach is evaluated against PGD and a recent ORCA-GPU approach in Section 5.

Table 1: Summary of the graphlet notation and properties

Graphlets are grouped by number of vertices (k-graphlets) and categorized into connected and disconnected graphlets. Connected graphlets of each size are then ordered by density. The complement of each connected graphlet is shown on the right and represent the disconnected graphlets. Note 4-path is a self-complementary^{*}. Graphlets of size k=2 are included for completeness.

	k	Со	NNEC	CTED	DIS	CONN	ECTED
	2	I	H_1	edge	•	H_2	2-node-independent
-	3	\square	H_3	triangle	•••	H_6	3-node-independent
LETS		Γ	H_4	2-star	•	H_5	3-node-1-edge
APH		\bowtie	H_7	4-clique	•••	H_{17}	4-node-independent
GI		\square	H_8	chordal-cycle	•••	H_{16}	4-node-1-edge
	4	Ν	H_9	tailed-triangle		H_{15}	4-node-2-star
		\Box	H_{10}	4-cycle		H_{14}	4-node-2-edge
		\square	H_{11}	3-star	arsigma	H_{13}	4-node-1-triangle
		\square	H_{12}	4-path^{\star}			

3. GRAPHLET DECOMPOSITION

Graphlets are at the heart and foundation of many network analysis tasks (*e.g.*, relational classification, network alignment) [21, 18, 12]. Given the practical importance of graphlets,

¹www.github.com/nkahmed/pgd

this paper proposes a hybrid CPU-GPU algorithm for computing the number of embeddings for both connected and disconnected k-vertex induced subgraphs (See Table 1).

3.1 Preliminaries

Let G = (V, E) be an undirected graph where V is the set of vertices and E is its edge set. The number of vertices is n = |V| and number of edges is m = |E|. We assume all vertex and edge sets are *ordered*, i.e., $V = \{v_1, v_2, ..., v_i, ..., v_N\}$ such that v_{i-1} appears before v_i and so forth. Similarly, the ordered edges are denoted $E = \{e_1, e_2, ..., e_i, ..., e_M\}$. Given a vertex $v \in V$, let $\Gamma(v) = \{w|(v, w) \in E\}$ be the set of vertices adjacent to v in G. The degree d_v of $v \in V$ is the size of the neighborhood $|\Gamma(v)|$ of v. We also define Δ to be the largest degree in G.

DEFINITION 1 (**Graphlet**). A graphlet $H_i = (V_k, E_k)$ is a subgraph consisting of a subset $V_k \subset V$ of k vertices from G = (V, E) together with all edges whose endpoints are both in this subset $E_k = \{ \forall e \in E \mid e = (u, v) \land u, v \in V_k \}.$

Let $\mathcal{H}^{(k)}$ denote the set of k-vertex induced subgraphs and $\mathcal{H} = \mathcal{H}^{(2)} \cup \cdots \cup \mathcal{H}^{(k)}$. A k-graphlet is simply an *induced* subgraph with k-vertices.

3.2 Problem Formulation

It is important to distinguish between the two fundamental classes of graphlets, namely, *connected* and *disconnected* graphlets (see Table 1). A graphlet is connected if there is a path from any node to any other node in the graphlet (see Definition 2). Table 1 provides a summary of the connected and disconnected k-graphlets of size $k = \{2, 3, 4\}$.

DEFINITION 2 (Connected graphlet). A k-graphlet $H_i = (V_k, E_k)$ is connected if there exists a path from any vertex to any other vertex in the graphlet H_i , $\forall u, v \in V_k, \exists P_{u-v} : u, \ldots, w, \ldots, v$, such that $d(u, v) \geq 0 \land d(u, v) \neq \infty$. By definition, a connected graphlet H_i has only one connected component (i.e., |C| = 1).

DEFINITION 3 (**Disconnected graphlet**). A k-graphlet $H_i = (V_k, E_k)$ is disconnected if there is not a path from any vertex $v \in H_i$ to any other vertex $w \in H_i$.

Unlike most existing work that is only able to compute connected graphlets of a certain size (such as k = 4), the goal of this work is to compute the frequency of both connected and disconnected graphlets of size $k \in \{2, 3, 4\}$. More formally,

PROBLEM 1 (Global graphlet counting). Given the graph G, find the number of embeddings (appearances) of each graphlet $H_i \in \mathcal{H}$ in the input graph G. We refer to this problem as the global graphlet counting problem. A graphlet $H_i \in \mathcal{H}$ is embedded in G, iff there is an injective mapping $\sigma : V_{H_i} \to V$, with $e = (u, v) \in E_{H_i}$ if and only if $e' = (\sigma(u), \sigma(v)) \in E$.

4. HYBRID CPU-GPU FRAMEWORK

This work proposes parallel graphlet decomposition methods that are designed to leverage (i) parallelism (multiple cores on a CPU or GPU) as well as (ii) heterogeneity that leverages simultaneous use of a CPU and GPU (as well as multiple CPUs and GPUs). To the best of our knowledge, this is the first work to use multiple GPUs (and of course multiple GPUs and CPUs) for computing induced subgraph statistics.



Figure 2: Let T be the set of nodes completing a triangle with the edge $(v, u) \in E$, and let S_v and S_u be the set of nodes that form a 2-star with v and u, respectively. Note that $S_u \cap S_v = \emptyset$ by construction and $|S_u \cup S_v| = |S_u| + |S_v|$. Further, let E_c be the set of edges that complete a cycle (of size 4) w.r.t. the edge (v, u) where for each edge $(p, q) \in E_c$ such that $p \in S_v$ and $q \in S_u$ and both $(p \cap S_u) \cup (q \cap S_v) = \emptyset$, that is, p is not adjacent to u $(p \notin \Gamma(u))$ and q is not adjacent to v $(q \notin \Gamma(v))$.

- Single GPU Methods (using multiple cores)
- Multi-GPU Methods.
- Hybrid Multi-core CPU-GPU Methods.

Methods from all three classes are shown to be effective on a large collection of graphs from a variety of domains (e.g., biological, social, and information networks [24]). In particular, methods from these classes have three important benefits. First, the performance is orders of magnitude faster than the state-of-the-art. Second, the GPU methods are cost effective enjoying superior performance per capita. Third, the performance per watt is significantly better than existing traditional CPU methods.

4.1 Parallel Graphlet Computations

Our algorithm searches over the set of edges E of the input graph G = (V, E). Given an edge $e = (u, v) \in E$, let $\Gamma(e)$ denote the edge neighborhood of e defined as:

$$\Gamma(e) = \Gamma(u, v) = \{ \Gamma(u) \cup \Gamma(v) \setminus \{u, v\} \},$$
(1)

where $\Gamma(u)$ and $\Gamma(v)$ are the neighbors of u and v, respectively. For convenience, let $\Gamma_e = G(\{\Gamma(v) - u\} \cup \{\Gamma(u) - v\})$ be the (explicit) edge-induced neighborhood subgraph. Given an edge $e = (u, v) \in E$, we explore the subgraph surrounding e(called the egonet of e), *i.e.*, the subgraph induced by both its endpoints and the vertices in its neighborhood. Our approach uses a specialized graph encoding based on the edge-CSC representation [23].

The parallel scheme leverages the fact that the induced subgraph (graphlet) problem can be solved independently for each edge-centric neighborhood $\Gamma(e_i) \in \{\Gamma(e_1), \ldots, \Gamma(e_M)\}$ in *G*, and therefore may be computed simultaneously in parallel. A processing unit denoted by ω refers to a single CPU/GPU worker (core). In the context of message-passing and distributed memory parallel computing, a node refers to another machine on the network with its own set of multicore CPUs, GPUs, and memory. Other important properties include the search order II in which edges are solved in parallel, the batch size *b* (number of jobs/tasks assigned to a worker by a dynamic scheduling routine), and the dynamic assignment of jobs (for load balancing).

While there are only a few such parallel graphlet algorithms, with the exception of PGD all of these methods are based on searching over the vertices (as opposed to the edges). However, as we shall see, the parallel performance of these approaches are guaranteed to suffer more from load balancing issues, communication costs, and other issues such as curse of the last reducer, etc.

Improved Load Balancing: Let $\mathbf{z}_j \in \mathbb{R}^N$ and $\mathbf{x}_j \in \mathbb{R}^M$ be counts of an arbitrary graphlet $H_j \in \mathcal{H}$ for vertices N = |V|and edges M = |E|, respectively. Given a vertex $v_i \in V$ and an edge $e_k \in E$, let z_{ij} and x_{kj} denote the number of vertex and edge incident counts of a graphlet H_j for vertex v_i and edge e_k , respectively. Furthermore, let

$$Z_j = \sum_{v_i \in V} z_{ij} \text{ and } X_j = \sum_{e_k \in E} x_{kj}$$
⁽²⁾

where Z_j and X_j are the global frequency of graphlet H_j in the graph G. Thus, it is straightforward to verify that $Z_j = X_j$. Further, let $\overline{Z}_j = Z_j/N$ and $\overline{X}_j = X_j/M$ be the mean vertex and edge count for graphlet $H_j \in \mathcal{H}$, respectively. Now, assuming $N \ll M$,² then $\overline{X}_j < \overline{Z}_j$. Clearly, more work is required to compute graphlets for each vertex v_i on average (compared to the number of graphlets counted per edge). This implies that edge-centric parallel algorithms are guaranteed to have better load balancing (among other important advantages) than existing vertex-centric algorithms.

4.2 Preprocessing Steps

Our approach benefits from the preprocessing steps below and the useful computational properties that arise.

- **P1** The vertices $V = \{v_1, \ldots, v_N\}$ are sorted from smallest to largest degree and relabeled such that $d(v_1) \leq d(v_2) \leq d(v_i) \leq d(v_N)$.
- **P2** For each $\Gamma(v_i) \in {\Gamma(v_1), \ldots, \Gamma(v_N)}$, order the neighbors $\Gamma(v_i) = {\ldots, w_j, \ldots, w_k, \ldots}$ s.t. j < k if $f(w_j) \ge f(w_k)$. Thus, the set of neighbors $\Gamma(v_i)$ are ordered from largest to smallest degree.
- **P3** Given an edge $(v, u) \in E$, we ensure that v is always the vertex with largest degree d_v , that is, $d_v \geq d_u$. This gives rise to many useful properties and as we shall see can lead to a significant reduction in runtime. For instance, our approach avoids searching both S_v and S_u for computing 4-cycles, and instead, allows us to compute 4-cycles by simply searching one of those sets. Thus, our approach always computes 4-cycles using S_u since (by the property above) is guaranteed to be less work/faster than if S_v is used (and the runtime difference can be quite significant).

Each step above is computed in $\mathcal{O}(N)$ or $\mathcal{O}(M)$ time and is easily parallelized.

4.3 Hybrid CPU-GPU

The algorithm begins by computing an edge ordering $\Pi = \{e_1, \ldots, e_M\}$ where the edges that are most difficult (with highly skewed, irregular, unbalanced degrees) are placed upfront, followed by edges that are more evenly balanced. In other words, Π is a permutation of the edges by some function or graph property $f(\cdot)$ such that k < j for e_k and e_j if

 $f(e_k) > f(e_j)$, and ties are broken arbitrarily (*e.g.*, using ids). For instance, edges can be ordered from largest to smallest degree (a proxy for the difficulty and unbalanced nature of the edge). For implementation purposes, Π is essentially a double-ended queue (dequeue) where elements can be added or removed from either end (*i.e.*, push and pop operations at either end). Afterwards, the previous edge ordering Π is split into three *initial* sets:

$$\Pi = \left\{ \underbrace{e_1, \dots, e_k}_{\Pi_{\rm cpu}}, \underbrace{\underbrace{e_{k+1}, \dots, e_{j-1}}_{\Pi_{\rm unproc}}}_{\Pi_{\rm unproc}}, \underbrace{e_j, e_{j+1}, \dots, e_M}_{\Pi_{\rm gpu}} \right\} \quad (3)$$

Now, the edges $\Pi_{\text{gpu}} = \{e_j, e_{j+1}, e_{j+2}, \ldots, e_M\}^3$ are split into p disjoint sets $\{I_1, I_2, \ldots, I_p\}$ with approximately equal work among each GPU device. This is accomplished by partitioning the edges in a round-robin fashion. Each GPU computes the induced subgraphs centered at each of the edges in I_i , which can be thought of as a local job queue for a particular GPU. Similarly, the CPU workers compute the induced subgraphs centered at each edge is assigned, it is removed from the corresponding local queue (for either CPUs or GPUs).

Once a CPU worker finishes all edges (or more generally tasks) in its local queue, it takes (dequeues) the next *b* unprocessed edges from the **front** of Π_{unproc} (the global queue of remaining/unassigned work) and pushes them to its local queue. On the other hand, once a GPU's local queue becomes empty (and thus the GPU becomes idle), it is assigned the next chunk of unprocessed edges from the *back* of the queue (*i.e.*, dequeued and pushed onto that GPU's local queue). Unlike the CPU, we must transfer the assigned edges to the corresponding GPU.

Finally, the graphlet counts from the multiple CPUs and GPU devices are combined to arrive at the final global and local graphlet counts. Similarly, if the local graphlet counts for each edge are warranted (also known as micro graphlet counts), then one would simply combine the per edge results to arrive at the final counts for each edge (and each graphlet). Recall that to compute all $k \in \{2,3,4\}$ -vertex graphlets, our algorithm only requires us to store counts for triangles, cliques, and cycles, and from these, we can easily derive the other counts for both connected and disconnected graphlets in constant time. This not only avoids communication costs (and thus significantly reduces the amount of communications that would otherwise be needed by other algorithms), but also reduces space and time.

Other important aspects include:

- Dynamic load balancing is performed for both CPUs and GPUs. For the CPUs, once a worker completes the tasks in its local queue, it immediately takes the next b tasks from the front of Π. We typically use a very small chunk size b for the CPU. The intuition is that these tasks are likely skewed and thus may take a significant amount of time to complete. Moreover, the overhead associated with this dynamic load balancing on the CPU is quite small (relative to the GPU of course, where the communication costs are significantly larger).
- To avoid communication costs and other performance degrading behavior, it is important that the GPUs are

²This holds in practice for nearly all real-world graphs

³The edges can be thought of as edge neighborhoods, since the induced subgraphs are counted over each edge neighborhood.

initially assigned a large fraction of the edges to process (which are then split among the GPUs). For majority of large real-world networks (with power-law), GPUs are initially assigned about 80% of the edges, and this seemed to work well, as it avoids both extremes (that is, significantly under- or oversubscribing the GPUs). In particular, significantly oversubscribing the GPUs causes a lot of work to be stolen by the GPUs (and the overhead associated with it, such as additional communication costs for the edges that will be stolen by the CPU (or even another GPU), whereas significantly undersubscribing the GPUs increase the load balancing overhead (causing additional communication costs, etc.). Ideally, we would want to assign the largest fraction of edges to the GPUs such that both the GPUs and CPUs finish at exactly the same time, and thus, avoid any communication or other costs associated with load balancing, etc.

- Note b is a chunk size, and is different for CPU and GPU devices. In particular, for CPU, we typically set b = 1, since these tasks are the most difficult to compute, and the runtime for each is likely to be skewed and irregular. This also helps avoid costs associated with work-stealing, which occurs when all such edge-centric tasks have been assigned, but not yet finished. Hence, to avoid the case where all but a single CPU workers have finished, and the remaining CPU worker has many edge-centric tasks in its local queue. In this case, of course, the tasks remaining in that CPU workers local queue would be stolen and distributed among the idle CPU workers. Note that as discussed later, one can also divide tasks at a much finer level of granularity, as many of the core computations carried out for a single edge-centric task can be computed independently. For instance, cliques and cycles are completely independent once the sets T and S_u are computed.
- Shared job queue with *work stealing* so that every multicore CPU and GPU remains fully utilized.
- Data is never moved once it is partitioned and distributed to the workers.
- It is straightforward to see that $|T_u| \leq |\Gamma(u)|$, that is, the number of triangles centered at u is bounded above by the number of neighbors (degree of u). Similarly, $|S_u| \leq |N(u)|$, that is, the number of 2-stars centered at u is bounded above by $|\Gamma(u)|$. Hence, $|S_u| + |T_u| = |\Gamma(u)|$. See Figure 2 for further intuition. We use the above fact to reduce the space requirements as well as improve locality.
- Given two vertices $w_i, w_j \in T$ that form a triangle with $e_k = (v, u)^4$, what is the most efficient search strategy (Alg. 3)? As a result of the past ordering and relabeling, searching the neighbors of w_j for w_i always results in less work. Observe that since the vertices in the set T are ordered from largest to smallest degree, then $|\Gamma(w_j)| < |\Gamma(w_i)|$, and thus, we search $|\Gamma(w_j)|$ for vertex w_i . Hence, if $(w_i, w_j) \in E$, then the subgraph induced by $\{v, u, w_i, w_j\}$ denoted $G(\{v, u, w_i, w_j\})$ is a 4-clique H_7 (See Table 1)

More importantly, our hybrid multi-core algorithm ensures all CPUs and GPU devices are always working simultaneously. The only exception arises when all but a single edge remains Algorithm 1 This algorithm computes all 3-vertex graphlets via neighbor iteration and also leverages a hash table for fast lookups. Used mainly for CPU-based enumeration of 3-vertex induced subgraphs.

Output:

A set T of vertices that complete triangles with e_k A set S_u of vertices that form 2-stars centered at u with e_k

1 procedure 3-GRAPHLETS (G, e_k, Ψ)

- 2 Set $T \leftarrow \emptyset$ and $S_u \leftarrow \emptyset$ 3 **parallel for** $w \in \Gamma(v)$ where $w \neq u$ do
- 4 Set $\Psi(w) = \lambda_1$
- 5 parallel for $w \in \Gamma(u)$ where $w \neq v$ do
- 6 if $\Psi(w) = \lambda_1$ then
- 7 $T \leftarrow T \cup \{w\}$ and set $\Psi(w) = \lambda_3$ \triangleright triangle
- 8 else $S_u \leftarrow S_u \cup \{w\}$ and set $\Psi(w) = \lambda_2$ \triangleright 2-star

9 Set $\mathbf{X}_{k,3} \leftarrow |T|$

Algorithm 2 This algorithm computes all 3-vertex graphlets using *binary search*. For multi-core architectures with limited memory (such as GPUs).

1 p	rocedure 3-GRAPHLETS-BINSEARCH(G , ϵ	$e_k)$
2	Let $e_k = (v, u) \in E$ be an edge in G s.t. d	$d_v \ge d_u$
3	Set $T \leftarrow \emptyset$ and $S_u \leftarrow \emptyset$	
4	parallel for each $w \in \Gamma(u)$ where $w \neq w$	v do
5	if $v \in \Gamma(w)$ via binary search then	$\triangleright \mathcal{O}(\log d_w)$
6	$T \leftarrow T \cup \{w\}$	⊳ triangle
7	else $S_u \leftarrow S_u \cup \{w\}$	\triangleright 2-star
8	end parallel	
9	Set $\mathbf{X}_{k,3} \leftarrow T $	

(is currently being processed by one or more other workers). In this case, there is no further work remaining (as the edge task has already been split and assigned to other workers). Nevertheless, we stress that this case is unlikely due to the fact that edges with large degrees are always processed first by both GPU and CPU workers. Thus, the remaining edge is likely to be easily computed. As mentioned, even a single edge neighborhood has many independent components that can be computed in parallel by other workers. Work-stealing is used to ensure that all processing units are fully utilized.

4.4 Finer-Granularity and Work Stealing

Each job in Π may represent the computations required for a single edge, or they may represent an even smaller unit of work. For a single edge $e_k \in E$, we compute (i) the sets T and S_u , then (ii) we find the total k-cliques for a given edge using T, and (iii) the total cycles of size k using S_u . Note that (ii) and (iii) are independent and thus can be computed in parallel. A job may also represent these smaller units of work. For instance, if T (or S_u) are large and computationally expensive to compute, then the worker in the simplest case can push T (or S_u) on the job queue for others to work on. Further, one can split T (or S_u) into even more finer grained independent jobs such that:

$$T = \{\underbrace{w_1, \cdots, w_i}_{T_{1:i}}, \underbrace{w_{i+1}, \cdots, w_t}_{T_{i+1:t}}\}$$
(4)

⁴Thus, there are two triangles centered at $e_k = (v, u)$, namely, $\{v, u, w_i\}$ and $\{v, u, w_j\}$.

Algorithm 3 Clique counts restricted to searching T

 $\begin{array}{cccc} 1 & \textbf{procedure } \text{CLIQUERES}(\ T,\ e_k \) \\ 2 & \text{Set } \mathbf{X}_{k,7} \leftarrow 0 \\ 3 & \textbf{parallel for each } w_i \in T \text{ in order } w_1, w_2, \cdots \text{ of } T \text{ do} \\ 4 & \textbf{for all each } w_j \in \{w_{i+1}, ..., w_{|T|}\} \text{ in order do} \\ 5 & \text{ if } w_i \in \Gamma(w_j) \text{ via binary search then} \\ 6 & \mathbf{X}_{k,7} \leftarrow \mathbf{X}_{k,7} + 1 \\ 7 & \textbf{end parallel} \\ 8 & \textbf{return } \mathbf{X}_{k,7} \end{array}$

Algorithm 4 Cycle counts restricted to S_u a	and \mathcal{S}_v	
--	---------------------	--

1 p	rocedure CycleRes(S_u, S_v, e_k)
2	Set $\mathbf{X}_{k,10} \leftarrow 0$
3	parallel for each $w \in S_u$ do
4	for all $r \in S_v$ do
5	if $r \in \Gamma(w)$ via binary search then
6	$\mathbf{X}_{k,10} \leftarrow \mathbf{X}_{k,10} + 1$
7	end parallel
8	return $\mathbf{X}_{k,10}$

Whenever possible work is stolen locally to minimize communication. Notice that if work is stolen locally from a GPU, *i.e.*, a GPU worker computes 4-cliques for a given edge e_k using the already computed set T, then expensive communication is avoided due to the set T being stored in global memory. In fact, one could store T and S_u contiguously in a single global array, and since $|T| + |S_u| = d_u$, then we can allocate contiguous memory for storing such sets for each edge. Alternatively, one could simply allocate an array of size $\Delta \cdot P_i$ where P_i is the number of cores (unique workers) for the i^{th} GPU (or CPU). Hence, Δ serves as an upper bound and in most cases requires significantly less space than the previous approach, since each GPU worker simply indexes into their own subarray which stores T and S_{u} . However, one may need to communicate T so that the GPU worker to compute 4-cliques from T has a copy (stored in their subarray). The only exception is if that is the last edge to be assigned to that GPU, and thus, the state of the subarray can persist until termination. However, suppose work is stolen from another GPU (not locally), then we would need to communicate not only the edge id, but the portion of T (or S_u) assigned to the GPU worker would also be needed.

To avoid bulk synchronization, workers store and aggregate statistics locally, and thus avoid unnecessary communications. In the case of global graphlet statistics computed for G, each worker maintains local aggregates and communicates them only once upon completion and thus has a cost of $O(\kappa)$ where κ is the number of induced subgraph statistics.

4.5 Unrestricted counts

The GPU workers use binary search to derive the sets T and S_u from Alg. 2, whereas the CPU workers compute T and S_u from Alg. 1 (or Alg. 2 if memory is limited and/or dynamically selected). The key difference is that the CPU workers create a fast hash table on the neighbors of v for $e_k = (v, u)$ allowing for o(1) time checks. Moreover, the hash table is also exploited for encoding nodes with particular types and enabling us to check an arbitrary type of vertex in o(1) time. Note that the λ_i 's in Alg. 1 represent distinct vertex types, and are used

later for finding cliques (Alg. 5) and cycles (Alg. 6) extremely fast. We also avoid the $O(|\Gamma(v)|)$ time it costs to reset Ψ for each edge by defining λ_i to be unique for each edge. Observe that each CPU worker in Alg. 1 maintains Ψ taking O(N)space, whereas Alg. 2 only requires $O(|T| + |S_u|) = O(d_u)$ space to store T and S_u (which Alg. 1 also uses). Note that S_v is easily computed (if needed) from Alg. 1 and Alg. 2 by simply setting $S_v \leftarrow \Gamma(v)$, and then removing each vertex $w \in T_e$ (on-the-fly), that is, $S_v \leftarrow S_v \setminus \{w\}$.

The proposed approach derives all k-graphlets for $k \in \{2, 3, 4\}$ using only the local edge-based counts of triangles, cliques, and cycles, along with a few other constant time graph and vertex parameters such as number of vertices N = |V|, edges M = |E|, as well as vertex degree $d_v = |\Gamma(v)|$. Given an edge $e_k \in E$ from G, let $\mathbf{X}_{k,3}, \mathbf{X}_{k,7}$, and $\mathbf{X}_{k,10}$ be the frequency of triangles, cliques, and cycles centered at the edge $e_k \in E$ in the graph G, respectively. Observe that $\mathbf{X}_{k,i}$ (or simply x_i) is the count of the induced subgraph H_i for an arbitrary edge e_k (See Table 1). The local (micro-level) 3-graphlets for edge e_k are as follows:

$$x_3 = |T| \tag{5}$$

$$x_4 = (d_u + d_v - 2) - 2|T| \tag{6}$$

$$x_5 = N - x_4 + |T| - 2 \tag{7}$$

$$x_6 = \binom{N}{3} - (x_3 + x_4 + x_5) \tag{8}$$

Further, notice that given $x_3 = |T|$ for $e_k = (v, u) \in E$, we can derive $|S_u|$ and $|S_v|$ (that is, the number of 2-star patterns centered at u and v of e_k , respectively) as:

$$|S_u| = d_u - |T| - 1 \tag{9}$$

$$|S_v| = d_v - |T| - 1 \tag{10}$$

Therefore, the number of two-stars centered at e_k denoted x_4 can be rewritten simply as $x_4 = |S_u| + |S_v|$. These 3-vertex induced subgraph statistics are then used as a basis to derive the induced subgraphs of size k + 1.

Notice that GPU workers compute $\mathbf{X}_{k,3}$ (as well as T and S_u) for edge e_k using Alg. 2, whereas CPU workers compute $\mathbf{X}_{k,3}$ from Alg. 1. Afterwards, $\mathbf{X}_{k,7}$ and $\mathbf{X}_{k,10}$ can be computed

Al	gorithm 5 Clique counts via neighbor iteration
1	procedure CLIQUE(Ψ , T, e_k)
2	Set $\mathbf{X}_{k,7} \leftarrow 0$
3	parallel for each $w \in T$ do
4	for each $r \in \Gamma(w)$ do
5	if $\Psi(r) = \lambda_3$ then Set $\mathbf{X}_{k,7} \leftarrow \mathbf{X}_{k,7} + 1$
6	Reset $\Psi(w)$ to 0
$\overline{7}$	end parallel
8	$\mathbf{return} \ \mathbf{X}_{k,7}$

Algorithm 6 Cycle counts via neighbor i	iteration	n
--	-----------	---

1 procedure CYCLE(Ψ , S_u , e_k)

2 Set $\mathbf{X}_{k,10} \leftarrow 0$

3 parallel for each $w \in S_u$ do

4 for each $r \in \Gamma(w)$ do

- 5 **if** $\Psi(r) = \lambda_2$ then set $\mathbf{X}_{k,10} \leftarrow \mathbf{X}_{k,10} + 1$
- 6 Reset $\Psi(w)$ to 0
- 7 end parallel
- 8 return $\mathbf{X}_{k,10}$

in any order as they are completely independent, and can even be stolen by another parallel worker (CPU and/or GPU worker that requires more work). A GPU worker computes the number of 4-cliques $\mathbf{X}_{k,7}$ centered at edge e_k via Alg. 3, whereas a CPU worker mainly leverages Alg. 5 for computing $\mathbf{X}_{k,7}$, but may also exploit Alg. 3 if determined (dynamically) that it requires less work than the other approach. Similarly, a GPU worker computes the number of 4-cycles $\mathbf{X}_{k,10}$ centered at e_k via Alg. 4, whereas a CPU worker computes $\mathbf{X}_{k,10}$ via Alg. 6.

Given only the triangles $\mathbf{X}_{k,3}$, cliques $\mathbf{X}_{k,7}$, and cycles $\mathbf{X}_{k,10}$ for each edge $e_k = (v, u) \in E$, we derive the unrestricted counts for connected and disconnected graphlets of size $k \in \{3, 4\}$. We first derive the unrestricted counts for connected and disconnected 3-graphlets as follows:

$$C_3 = \sum_{e_k = (v,u) \in E} \mathbf{X}_{k,3} = \sum_{e_k = (v,u) \in E} |T|$$
(11)

$$C_4 = \sum_{e_k = (v,u) \in E} |S_v| + |S_u| \tag{12}$$

$$C_5 = \sum_{e_k = (v,u) \in E} N - (|S_v| + |S_u| + |T|) - 2 \qquad (13)$$

Note that C_6 is not needed since X_6 can be computed directly from X_3 , X_4 , and X_5 . The unrestricted counts for the connected 4-graphlets (4-vertex connected induced subgraphs) are derived as follows:

$$C_7 = \sum_{e_k = (v,u) \in E} \mathbf{X}_{k,7}$$
(14)

$$C_8 = \sum_{e_k = (v,u) \in E} \binom{T}{2}$$
(15)

$$C_{9} = \sum_{e_{k} = (v,u) \in E} |T| \cdot |S_{v}| \cdot |S_{u}|$$
(16)

$$C_{10} = \sum_{e_k = (v,u) \in E} \mathbf{X}_{k,10}$$
(17)

$$C_{11} = \sum_{e_k = (v,u) \in E} \left(\begin{smallmatrix} |S_v| \\ 2 \end{smallmatrix} \right) + \left(\begin{smallmatrix} |S_u| \\ 2 \end{smallmatrix} \right)$$
(18)

$$C_{12} = \sum_{e_k = (v,u) \in E} |S_v| \cdot |S_u|$$
(19)

Given an arbitrary edge $e_k \in E$, let us define $D_e = N - (|S_v| + |S_u| + |T|) - 2$ for convenience. The unrestricted counts for the disconnected 4-graphlets (4-vertex disconnected induced subgraphs) are derived as follows:

$$C_{13} = \sum_{e_k = (v, u) \in E} |T| \cdot D_e$$
(20)

$$C_{14} = \sum_{e_k = (v,u) \in E} M - d_v - d_u + 1$$
(21)

$$C_{15} = \sum_{e_k = (v,u) \in E} (|S_v| + |S_u|) \cdot D_e$$
(22)

$$C_{16} = \sum_{e_k = (v,u) \in E} \begin{pmatrix} D_e \\ 2 \end{pmatrix}$$
(23)

Recall that all unrestricted counts $C_3, ..., C_{16}$ are computed in $\mathcal{O}(M)$ time and easily parallelized.

4.6 Global Graphlet Frequencies

Now, using the above unrestricted counts, we can derive

the connected and disconnected global (macro-level) graphlet counts of size $k \in \{2, 3, 4\}$ for the graph G as:

$$X_1 = M \tag{24}$$

$$X_2 = \left(\begin{smallmatrix} N\\2 \end{smallmatrix}\right) - M \tag{25}$$

$$X_3 = \frac{1}{3} \cdot C_3 \tag{26}$$

$$X_4 = \frac{1}{2} \cdot C_4 \tag{27}$$

$$X_5 = C_5 \tag{28}$$

$$X_{6} = \binom{N}{3} - (X_{3} + X_{4} + X_{5})$$
(29)
$$X_{--} V_{--} C$$
(20)

$$X_7 = \frac{1}{6} \cdot C_7 \tag{30}$$

$$X_8 = C_8 - C_7 \tag{31}$$

$$X_8 = C_8 - C_7 \tag{31}$$

$$X_9 = \frac{1}{2}(C_0 - 4X_8) \tag{32}$$

$$X_{10} = \frac{1}{4} \cdot C_{10}$$
(52)
$$X_{10} = \frac{1}{4} \cdot C_{10}$$
(33)

$$X_{10} = \frac{1}{4} \cdot C_{10}$$
(33)
$$X_{11} = \frac{1}{3}(C_9 - X_9)$$
(34)

$$X_{12} = C_{12} - C_{10} \tag{35}$$

$$X_{13} = \frac{1}{3} \cdot \left(C_{13} - X_9\right) \tag{36}$$

$$X_{14} = \frac{1}{2} \cdot \left(C_{14} - 6X_7 - 4X_8 - 2X_9 - 4X_{10} - 2X_{12} \right)$$
(37)

$$X_{15} = \frac{1}{2} \cdot (C_{15} - 2X_{12}) \tag{38}$$

$$X_{16} = C_{16} - 2X_{14} \tag{39}$$

$$X_{17} = \binom{N}{4} - \sum X_i \quad \text{for } i = 7, \dots, 16$$
 (40)

4.7 Complexity

This section gives the space and time complexity. Let T_{max} and S_{max} denote the maximum number of triangles and stars incident to a selected edge $e \in E$. Our algorithm solves the graphlet decomposition problem for k-vertex induced subgraphs in:

$$\mathcal{O}\left(M\Delta(T_{\max}+S_{\max})\right)$$

Using *M* processors (cores, workers), this reduces to $\mathcal{O}(\Delta(T_{\max} + S_{\max}))$. For the local graphlet problem, finding all graphlets centered at an edge e = (v, u) in *G* is solved in $\mathcal{O}(d_u(|T| + |S_u|))$ time.

Given an arbitrary edge $e_k \in E$ in G, Alg 2 computes Tand S_u in $\mathcal{O}(\sum_{w \in \Gamma(u)} \log d_w)$ time. This is due to the fact that each vertex $w \in \Gamma(u)$ takes $\mathcal{O}(\log d_w)$ time to check if $(v,w) \in E$. However, Alg 1 finds T and S_u in $\mathcal{O}(d_v + d_u)$ time. In particular, Alg 1 first marks the neighbors of v in $\mathcal{O}(d_v)$ time. Now, for each $w \in \Gamma(u)$, we check in o(1) time if $(v, w) \in E$ (using the fast lookup table), as this implies that w completes a triangle with (v, u). Thus, taking a total of $\mathcal{O}(d_v + d_u)$ time. Note that in terms of space, Alg 2 is more efficient, since Alg 1 requires an additional $\mathcal{O}(N)$ space for Ψ . Note that each parallel worker maintains a local hash table Ψ , and thus too expensive for GPUs that have thousands of workers/cores. Thus, Alg. 2 is used for GPUs since they have limited memory while also having many more cores (workers) than CPUs. Now, given T and S_u for edge e_k , we compute 4cliques in $\mathcal{O}(\Delta|T|)$ time. More precisely, the 4-cliques centered at $e_k \in E$ are computed in exactly $\mathcal{O}(\sum_{w_i \in T} \Gamma(w_i))$ time. In a similar fashion, Similarly, we compute 4-cycles in $\mathcal{O}(\Delta | S_u |)$ time, and more precisely, $\mathcal{O}(\sum_{w_i \in T} \Gamma(w_i))$.

Space: Each GPU has a copy of the graph taking $\mathcal{O}(|E| + |V|+1)$ space and a set of edges I_i . In addition, each GPU has an array of length $P_i \cdot \Delta_i$ where Δ_i is the maximum degree of any vertex u in the set of edges I_i and P_i is the total cores

Table 2: Results demonstrate the effectiveness of the hybrid parallel graphlet decomposition algorithms. In particular, we find that by leveraging the unique features and advantages of CPUs and GPUs, one can obtain significant speedups over existing methods that leverage only CPUs or GPUs, but not both.

					Speed	up (time	es faster)
						Multi-	
graph	\mathbb{K}	Δ	$\Delta_{\rm gpu}$	α	GPU	GPU	Hybrid
socfb-Texas84	81	6312	450	0.031	$4.65 \times$	21.91 x	263.26x
socfb-UF	83	8246	370	0.05	1.6x	55.65x	165.63x
socfb-MIT	72	708	266	0.7	$11.98 \times$	28.47x	106.14x
socfb-Stanford3	91	1172	365	0.05	$21.07 \times$	63x	133.15x
socfb-Wisc87	60	3484	300	0.04	$17.88 \times$	142.41 x	189.08x
socfb-Indiana	76	1358	329	0.04	$22.25 \times$	96.89x	207.11x
soc-flickr	309	4369	4196	0.04	$7.32 \times$	31.85x	102.24x
soc-google-plus	135	1790	328	0.07	4.95x	$11.98 \times$	56.03x
soc-youtube	49	25409	1079	0.07	3.87x	26.82x	180.64x
soc-brightkite	52	1134	132	0.12	2.51 x	8.09×	17.67x
soc-livejournal	213	2651	157	0.05	8.92x	70.01x	98.83x
soc-twitter	125	51386	13533	0.05	2.68x	21.76x	372.72x
soc-orkut	230	27466	646	0.05	6.12x	57.71x	129.26x
ia-enron-large	43	1383	243	0.176	2.94 x	10.79x	28.30x
ia-wiki-Talk	58	1220	1034	0.02	23.35x	37.50x	85.46x
ca-HepPh	238	491	169	0.35	1.42 x	$6.62 \times$	17.14x
brain-mouse-ret1	121	744	712	0.26	3.21 x	5.14 x	32.71×
web-baidu-baike	78	97848	11919	0.03	4.83x	39.55x	156.45x
web-arabic05	101	1102	49	0.14	5.19x	29.51 x	60.02x
tech-internet-as	23	3370	208	0.35	$1.26 \times$	3.55x	12.78x
tech-as-skitter	111	35455	4768	0.08	$0.62 \times$	1.89x	58.62x
C500-9	432	468	450	0.42	$3.13 \times$	21.99 x	33.23x
p-hat500-1	86	204	199	0.1	13.8x	$28.02 \times$	46.23×
p-hat1000-1	163	408	323	0.1	$14.95 \times$	67.87x	117.3x

(workers/processing units) of the i^{th} GPU. Thus, Δ_i is an upper bound on the maximum size any edge neighborhood task would require, see Figure 2 for intuition. Thus, the total space for the i^{th} GPU is:

$$\mathcal{O}\left(|E| + |V| + |I_i| + (P_i \cdot \Delta_i)\right)$$

Note the above is for global (macro) graphlet counts. Now, suppose we want to compute the local graphlet counts centered at each edge. This would require three additional arrays all of size $|I_i|$. Notice that for global macro counts, each GPU simply communicates the total number of triangles, cliques, and cycles, which can be easily aggregated locally (as opposed to the count of triangles, cliques, and cycles for each edge in I_i). This avoids expensive and unnecessary communications.

Each CPU worker has a local hash table Ψ taking $\mathcal{O}(N)$ space, as well as two arrays for storing T and S_u of length Δ , thus the total space per CPU worker is $\mathcal{O}(N + 2\Delta)$. Assuming P workers, the total space is: $\mathcal{O}(P(N + 2\Delta)) = \mathcal{O}(PN)$. The graph is $\mathcal{O}(|V| + |E| + 1)$ and shared among the CPU workers.

4.8 Representative Methods from Framework

Observe that our approach succinctly generalizes across the spectrum of GPU graphlet methods, including the following three classes:

- Single GPU algorithm uses only a single GPU device for computing induced subgraphs of size $k \in \{2, 3, 4\}$.
- Multi-GPU algorithm that leverages multiple GPU devices for graphlets.

• Hybrid Multi-core CPU-GPU algorithm that leverages the unique computing capabilities of each type of processing unit.

For instance, if we manually set the upper bound on the local CPU queue to $\alpha = 0$, then this gives rise to the Multi-GPU algorithm that leverages *only* GPU devices. Similarly, if we manually set $\alpha = 0$ and set the maximum number of GPU devices to 1 (*i.e.*, -GPUs 1), then we have the single GPU algorithm. We investigate each of these methods in Section 5 and evaluate them against the state-of-the-art CPU parallel graphlet decomposition (PGD) framework [4, 3].

5. EXPERIMENTS

This experiments in this section are designed to answer the following two questions that lie at the heart of this work. First, does the GPU and multi-GPU only algorithms improve performance over the state-of-the-art CPU method? Second, can we leverage the unique features and advantages of both CPUs and GPUs to further improve performance?

For these experiments, two Intel Xeon CPU E5-2687 @ 3.10GHz were used with 8 cores each. Further, we used 8 GeForce GTX TITAN Black GPUs and each GPU has 2880 cores (889 MHz base) and 6144 MB memory. We demonstrates the effectiveness of our approach on a variety of real-world networks from a range of domains with different properties. We evaluate three parallel graphlet methods from the proposed framework including: Single GPU only, Multi-GPUs, and Hybrid approach that uses multiple CPUs and GPUs. Table 2 shows the speedup relative to the state-of-the-art method called PGD [4, 3]. These results demonstrate the effectiveness of the proposed methods. In particular, the multi-GPU only and hybrid (CPU+GPU) are orders of magnitude faster than PGD. For certain types of real-world networks, we find that the methods are over 100 times faster, with the largest improvement being 372x for soc – twitter. In Table 2, Δ_{gpu}

Table 3: *Connected 4-graphlet* frequencies for a variety of the real-world networks investigated from different domains.

		Co	nnected	Graph	lets	
graph	\bowtie	\square	\square	\square	\square	\square
socfb-Texas84	$70.7 \mathrm{M}$	376M	1.2B	215M	664M	3.9B
socfb-UF	98M	433M	708M	186M	778M	874M
socfb-MIT	$13.7 \mathrm{M}$	88.5M	909M	50.9 M	498M	3.8B
socfb-Stanford3	37.1M	226M	659M	$151 \mathrm{M}$	600M	1.8B
socfb-Wisc87	23M	121M	1.9B	59.3M	1.3B	3.8B
socfb-Indiana	60.2M	269M	1.6B	141M	495M	3.9B
soc-flickr	311M	1B	208M	252M	1.2B	3.7B
soc-google-plus	186M	994M	204M	463M	668M	3.7B
soc-youtube	3.8M	156M	1.2B	162M	1B	2.3B
soc-livejournal	307M	1.9B	1.8B	465M	778M	3.5B
soc-twitter	430M	2.3B	1.7B	990M	314M	1.9B
soc-orkut	280M	3.2B	953M	595M	520M	2B
ia-enron-large	2.3M	22.5M	376M	$6.8 \mathrm{M}$	185M	1.4B
ia-wiki-Talk	2.2M	32.3M	668M	33.8M	766M	1.5B
ca-HepPh	150M	$35.2 \mathrm{M}$	462M	821k	143M	204M
prain-mouse-ret1	71.4M	303M	1.1B	47.4M	1.1B	1.1B
web-baidu-baike	$27.8 \mathrm{M}$	248M	476M	653M	1.3B	1.2B
web-arabic05	232M	3.4M	$26.5 \mathrm{M}$	79.2k	490M	27.3M
tech-as-skitter	149M	2.4B	571M	817M	808M	2.8B
C500-9	656M	909M	201M	50.2M	7.3M	22.3M
p-hat1000-1	20.3M	265M	1.3B	282M	1.2B	3B



Figure 3: Runtimes of the existing single-GPU method is normalized w.r.t. the runtime of our approach. For the BA graph, our approach is $\approx 55x$ faster than Orca-GPU, and $\approx 71x$ faster than Orca-GPU-S3. Furthermore, for the ER graph, our approach is $\approx 124x$ faster than Orca-GPU, and $\approx 181x$ faster than Orca-GPU-S3.

represents the maximum degree of an edge assigned to the GPUs. Strikingly, we observe that Δ_{gpu} is usually much less than Δ . Graphlet statistics for a few of the graphs are shown in Table 3.

In addition, we also compare to the approach proposed by [19], see Figure 3. Recall that approach is not hybrid (as the one proposed in this work) nor does it use multiple GPUs. Furthermore, we compute all connected and disconnected graphlets up to size k = 4 (thus including $k \in \{2, 3, 4\}$), whereas [19] is only able to compute connected graphlets of size k = 4 — a much smaller subset. Despite this difference, our approach is still orders of magnitude faster as shown in Figure 3. In particular, Figure 3 demonstrates the effectiveness of our approach. For comparison, we report results using the same benchmark graphs (see [19]). In particular, that work used random graphs generated from Barabási-Albert (BA), Erdős-Rénvi (ER) and geometric algorithms (GEO). From each, we selected the largest graph for comparison (1K vertices and 150K edges). Strikingly, our approach is orders of magnitude faster with up to 181x improvement. It is important to note that our approach is significantly faster for real-world networks where the work associated with each edge is fundamentally unbalanced and highly skewed. The random graphs in Figure 3 represent prime candidates for that approach. Observe that unlike many real-world networks, the degree distribution of these synthetic graphs is not skewed (does not obey a power-law) and are relatively dense (30%). Nevertheless, our approach is still orders of magnitude faster.

Table 4: Varying edge ordering can significantly impact performance. Results demonstrate the effectiveness of the initial ordering technique.

$\begin{array}{c c} & Descending & Reverse \\ \hline graph & d & vol & d^{-1} & v \\ \hline \hline & & & & & & \\ \hline & & & & & & & \\ \hline & & & &$	
$\frac{\text{graph}}{d} \text{vol} d^{-1} d^$	ORDER
	ol ⁻¹
soctb-Texas84 263.3x 284.1x 23.5x	0.8x

Table 4 demonstrates the impact of various edge orderings (using the Hybrid GPU-CPU approach). In particular, we investigate ordering edges from largest to smallest degree **d** and volume vol $(i.e., vol(e_k) = vol(u, v) = \sum_{w \in \Gamma(u,v)} d_w$, which is the sum of degrees of vertices in $\Gamma(u, v)$). The impact of the reverse ordering is also investigated, *i.e.*, ordering edges from smallest to largest degree and degree volume denoted by \mathbf{d}^{-1} and vol^{-1} , respectively. Notably, we observe that the ordering can significantly influence performance and in some instances may even lead to slower performance than a GPU-only approach.



Figure 4: CPU and GPU processing time for each edge neighborhood. The x-axis represents edge neighborhoods which are computed by either a CPU or GPU, whereas the y-axis is the time (ms) for computing each edge neighborhood. See text for discussion.

Figure 4 validates the proposed hybrid (GPU+CPU) approach. Recall that our framework dynamically partitions the edges among the CPUs and GPUs based on some notion of difficulty. In particular, we see that edge neighborhoods assigned to the CPU are indeed difficult and require significantly more time to compute than the edge neighborhoods processed by the GPUs. In addition, our approach is more space-efficient for the GPU, which has significantly less RAM than the CPU. Recall that GPUs are assigned neighborhoods that are significantly more sparse than those given to the CPUs, and therefore requires less space as well as avoiding expensive communications. Results in Figure 4 also demonstrate the effectiveness of the dynamic load balancing approach used in the hybrid graphlet algorithm, as it assigns edge neighborhoods to the corresponding "best" processor type. Moreover, the above also demonstrates the effectiveness and importance of the initial edge ordering Π . Notably, we found degree to be a useful approximation of the actual work required to compute the local graphlet counts for each edge neighborhood. However, ordering edges by vol can lead to improvements in performance over degree (Table 4). Moreover, we also exploit different graphlet algorithms (that essentially trade-off space for time) by dynamically selecting the appropriate one at runtime based on simple heuristics that can be derived in constant time. The GPUs are then used to compute graphlets for edge neighborhoods that are more well-balanced and regular, which is exactly the type of problems for which they are most effective.

Finally, memory use for GPUs is shown in Figure 5. We notice that the graph usually exceeds the others. For as-skitter, the memory used to store T and S_u for each GPU worker is



Figure 5: Average memory (MB) per GPU for three real-world networks.

slightly larger due to the large degrees. Finally, as expected the set of edges I_i is always less than the others. We also note that the communication overhead is insignificant compared to the time to compute graphlets.

6. CONCLUSION

This work proposed a parallel graphlet decomposition method that effectively leverages multiple CPUs and GPUs simultaneously. The algorithm is designed to exploit the unique features and strengths of each type of processing unit and is shown to be orders of magnitude faster than existing work that is based on only a single type of computing device. In particular, the proposed methods were shown to be up to 300+ times faster than the recent state-of-the-art. Besides being orders of magnitude faster, our approach is also more energy efficient. The proposed methods are also well-suited for unbiased graphlet estimation, and we plan to investigate this problem in future work.

7. REFERENCES

- N. K. Ahmed, N. Duffield, J. Neville, and R. Kompella. Graph sample and hold: A framework for big graph analytics. In *SIGKDD*, pages 1–10, 2014.
- [2] N. K. Ahmed, J. Neville, and R. Kompella. Space-efficient sampling from social activity streams. In SIGKDD BigMine, pages 1–8, 2012.
- [3] N. K. Ahmed, J. Neville, R. A. Rossi, and N. Duffield. Efficient motif counting for large-scale network analysis and mining. In *ICDM*, pages 1–10, 2015.
- [4] N. K. Ahmed, J. Neville, R. A. Rossi, N. Duffield, and T. L. Willke. Graphlet decomposition: Framework, algorithms, and applications. In *KAIS*, pages 1–32, 2016.
- [5] L. Akoglu, H. Tong, and D. Koutra. Graph based anomaly detection and description: a survey. *Data Mining and Knowledge Discovery*, pages 1–63, 2014.
- [6] I. Bhattacharya and L. Getoor. Entity resolution in graphs. *Mining graph data*, page 311, 2006.
- [7] A. R. Brodtkorb, T. R. Hagen, and M. L. Sætra. Graphics processing unit (gpu) programming strategies and trends in gpu computing. *Journal of Parallel and Distributed Computing*, 73(1):4–13, 2013.
- [8] K. Faust. A puzzle concerning triads in social networks: Graph constraints and the triad census. *Social Networks*, 32(3):221–233, 2010.
- [9] O. Frank. Triad count statistics. Annals of Discrete Mathematics, 38:141–149, 1988.

- [10] L. Getoor and B. Taskar. Introduction to statistical relational learning. MIT press, 2007.
- [11] A. Gharaibeh, E. Santos-Neto, L. B. Costa, and M. Ripeanu. The energy case for graph processing on hybrid cpu and gpu systems. In *IA3 Workshop*, 2013.
- [12] W. Hayes, K. Sun, and N. Pržulj. Graphlet-based measures are suitable for biological network comparison. *Bioinformatics*, 29(4):483–491, 2013.
- [13] T. Hočevar and J. Demšar. A combinatorial approach to graphlet counting. *Bioinformatics*, 30(4):559–565, 2014.
- [14] P. W. Holland and S. Leinhardt. Local structure in social networks. *Sociological Meth.*, 7:pp. 1–45, 1976.
- [15] V. W. Lee, C. Kim, J. Chhugani, M. Deisher, et al. Debunking the 100X GPU vs. CPU myth: an evaluation of throughput computing on CPU and GPU. In ACM SIGARCH Comp. Arch. News, pages 451–460, 2010.
- [16] A. D. Malony, S. Biersdorff, S. Shende, H. Jagode, S. Tomov, G. Juckeland, R. Dietrich, D. Poole, and C. Lamb. Parallel performance measurement of heterogeneous parallel systems with gpus. In *ICPP*, pages 176–185, 2011.
- [17] D. Marcus and Y. Shavitt. Rage–a rapid graphlet enumerator for large networks. *Computer Networks*, 56(2):810–819, 2012.
- [18] T. Milenkoviæ and N. Pržulj. Uncovering biological network function via graphlet degree signatures. *Cancer* informatics, 6:257, 2008.
- [19] A. Milinković, S. Milinković, and L. Lazicć. A contribution to acceleration of graphlet counting. In *Infoteh Jahorina Symposium*, volume 14, pages 741–745.
- [20] C. Noble and D. Cook. Graph-based anomaly detection. In SIGKDD, pages 631–636, 2003.
- [21] N. Pržulj, D. G. Corneil, and I. Jurisica. Modeling interactome: scale-free or geometric? *Bioinformatics*, 20(18):3508–3515, 2004.
- [22] R. Rossi and N. Ahmed. Role discovery in networks. *TKDE*, 27(4):1112–1131, 2015.
- [23] R. A. Rossi. Fast Triangle Core Decomposition for Mining Large Graphs. In *PAKDD*, pages 310–322, 2014.
- [24] R. A. Rossi and N. K. Ahmed. The network data repository with interactive graph analytics and visualization. In AAAI, pages 4292–4293, 2015.
- [25] M. Rupp and G. Schneider. Graph kernels for molecular similarity. *Molecular Informatics*, 29(4):266–273, 2010.
- [26] S. Schaeffer. Graph clustering. Computer Science Review, 1(1):27–64, 2007.
- [27] N. Shervashidze, T. Petri, K. Mehlhorn, K. M. Borgwardt, and S. Vishwanathan. Efficient graphlet kernels for large graph comparison. In *AISTATS*, 2009.
- [28] K. Taylor-Sakyi. Big data: Understanding big data. arXiv preprint arXiv:1601.04602, 2016.
- [29] S. Vishwanathan, N. Schraudolph, R. Kondor, and K. Borgwardt. Graph kernels. *JMLR*, 11, 2010.
- [30] S. Wernicke and F. Rasche. Fanmod: a tool for fast network motif detection. *Bioinformatics*, 22(9):1152–1153, 2006.
- [31] L. Zhang, R. Hong, Y. Gao, R. Ji, Q. Dai, and X. Li. Image categorization by learning a propagated graphlet path. *TNNLS*, 27(3):674–685, March 2016.
- [32] L. Zhang, M. Song, Z. Liu, X. Liu, J. Bu, and C. Chen. Probabilistic graphlet cut: Exploiting spatial structure cue for weakly supervised image segmentation. In *CVPR*, pages 1908–1915, 2013.