

GraphMineSuite: Enabling High-Performance and Programmable Graph Mining Algorithms with Set Algebra

Maciej Besta^{1*}, Zur Vonarburg-Shmaria¹, Yannick Schaffner¹, Leonardo Schwarz¹, Grzegorz Kwasniewski¹, Lukas Gianinazzi¹, Jakub Beranek², Kacper Janda³, Tobias Holenstein¹, Sebastian Leisinger¹, Peter Tatkowski¹, Esref Ozdemir¹, Adrian Balla¹, Marcin Copik¹, Philipp Lindenberger¹, Pavel Kalvoda¹, Marek Konieczny³, Onur Mutlu¹, Torsten Hoefler^{1*}

¹Department of Computer Science, ETH Zurich; ²Faculty of Electrical Engineering and Computer Science, VSB; ³Department of Computer Science, AGH-UST Krakow; *Corresponding authors

ABSTRACT

We propose GraphMineSuite (GMS): the first benchmarking suite for graph mining that facilitates evaluating and constructing high-performance graph mining algorithms. First, GMS comes with a benchmark specification based on extensive literature review, prescribing representative problems, algorithms, and datasets. Second, GMS offers a carefully designed software platform for seamless testing of different fine-grained elements of graph mining algorithms, such as graph representations or algorithm subroutines. The platform includes parallel implementations of more than 40 considered baselines, and it facilitates developing complex and fast mining algorithms. High modularity is possible by harnessing set algebra operations such as set intersection and difference, which enables breaking complex graph mining algorithms into simple building blocks that can be separately experimented with. GMS is supported with a broad concurrency analysis for portability in performance insights, and a novel performance metric to assess the throughput of graph mining algorithms, enabling more insightful evaluation. As use cases, we harness GMS to rapidly redesign and accelerate state-of-the-art baselines of core graph mining problems: degeneracy reordering (by up to >2×), maximal clique listing (by up to >9×), *k*-clique listing (by 1.1×), and subgraph isomorphism (by up to 2.5×), also obtaining better theoretical performance bounds.

Website: http://spcl.inf.ethz.ch/Research/Parallel_Programming/GMS

1 INTRODUCTION AND MOTIVATION

Graph mining is used in many compute-related domains, such as social sciences (e.g., studying human interactions), bioinformatics (e.g., analyzing protein structures), chemistry (e.g., designing chemical compounds), medicine (e.g., drug discovery), cybersecurity (e.g., identifying intruder machines), healthcare (e.g., exposing groups of people who submit fraudulent claims), web graph analysis (e.g., providing accurate search services), entertainment services (e.g., predicting movie popularity), and many others [63, 73, 113, 120]. Yet, graphs can reach one trillion edges (the Facebook graph (2015) [71]) or even 12 trillion edges (the Sogou webgraph (2018) [143]), requiring unprecedented amounts of compute power to solve even simple graph problems such as BFS [143]. For example, running PageRank on the Sogou webgraph using 38,656 compute nodes (10,050,560 cores) on the Sunway TaihuLight supercomputer [96] (nearly the full scale of TaihuLight) takes 8 minutes [143]. Harder problems, such as mining cliques, face even larger challenges.

At the same time, massive parallelism has become prevalent in modern compute devices, from smartphones to high-end servers [18],

bringing a promise of high-performance parallel graph mining algorithms. Yet, several issues hinder achieving this. First, a large number of graph mining algorithms and their variants make it hard to identify the most relevant baselines as either promising candidates for further improvement, or as appropriate comparison targets. Similarly, a plethora of available networks hinder selecting relevant input datasets for evaluation. Second, even when experimenting with a single specific algorithm, one often faces numerous design choices, for example which graph representation to use, whether to apply graph compression, how to represent auxiliary data structures, etc.. Such choices may significantly impact performance, often in a non-obvious way, and they may require a large coding effort when trying different options [78]. This is further aggravated by the fact that developing efficient parallel algorithms is usually challenging [14] because one must tackle issues such as deadlocks, data conflicts, and many others [14].

To address these issues, we introduce **GraphMineSuite (GMS)**, a benchmarking suite for high-performance graph mining algorithms. GMS provides an exhaustive benchmark specification **S**. Moreover, GMS offers a novel performance metric **M** and a broad theoretical concurrency analysis **C** for deeper performance insights beyond simple empirical run-times. To maximize GMS' usability, we arm

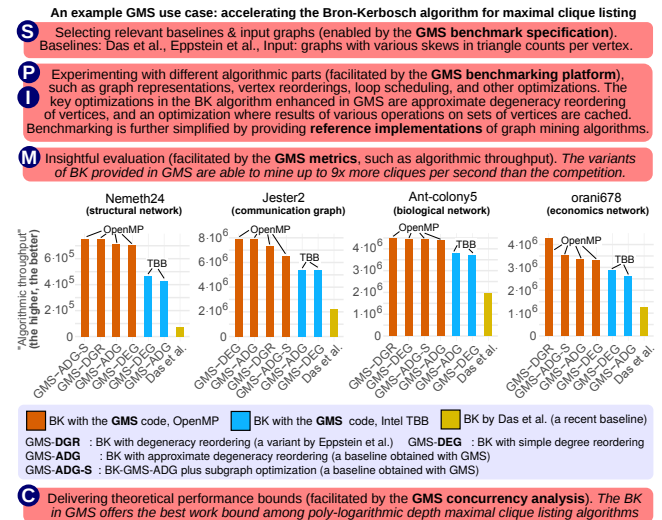


Figure 1: Performance advantages of the parallel Bron-Kerbosch (BK) algorithm implemented in GMS over a state-of-the-art implementation by Das et al. [79] and a recent algorithm by Eppstein et al. [91] (GMS-DGR) using a novel performance metric "algorithmic throughput" that shows a number of maximal cliques found per second. Details of experimental setup: Section 8.

Reference / Infrastructure	Focus on what problems?	Pattern Matching					Learning			Opt	Vr	Remarks	
		mC?	kC?	dS?	sI?	fS?	vS?	IP?	cl?				cD?
[B] Cyclone [201]	Graph database queries	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐*	☐**	*Only shortest paths. **Only degree centrality.
[B] GBBS [84] + Ligma [192]	More than 10 “low-complexity” algorithms	✗	☐	☐	✗	✗	✗	✗	✗	✗	☐*	☐*	*Support for degeneracy, but no explicit rank derivation.
[B] GraphBIG [165]	Mostly vertex-centric schemes	✗	☐*	✗	✗	✗	✗	✗	✗	✗	☐**	☐	*Only $k = 3$. **Only shortest paths and one coloring scheme.
[B] GAPBS [20]	Seven “low-complexity” algorithms	✗	☐*	✗	✗	✗	✗	✗	✗	✗	☐**	✗	*Only $k = 3$. **Only shortest paths.
[B] LDBC [51]	Graph database queries	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐**	✗	*Only one clustering coefficient. **Only shortest paths.
[B] WGB [12]	Mostly online queries	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐**	✗	*Only one clustering scheme. **Only shortest paths.
[B] PBBS [44]	General parallel problems	✗	✗	✗	✗	✗	✗	✗	✗	☐	☐	✗	Only graph optimization problems are considered
[B] Graph500 [162]	Graph traversals	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐*	✗	*Support for shortest paths only.
[B] HPCS [15]	Two “low-complexity” algorithms	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐*	✗	*Just one clustering scheme is considered
[B] Han et al. [106]	Evaluation of various graph processing systems	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐*	✗	*Support for Shortest Paths and Minimum ST
[B] CRONO [6]	Focus on futuristic multicores	✗	✗	✗	✗	✗	✗	✗	✗	☐	☐*	☐**	*Only shortest paths. **Only triangle counting.
[B] GARDENIA [218]	Focus on future accelerators	✗	✗	✗	✗	✗	✗	✗	✗	✗	☐*	☐**	*Only shortest paths. **Triangle counting and vertex coloring.
[F] A framework, e.g., Peregrine [118] or Fractal [86] (more at the end of Section 1)		☐*	☐*	☐*	☐*	☐*	✗	✗	✗	✗	✗	✗	*No good performance bounds (focus on expressiveness), not competitive to specific parallel mining algorithms
[B] GMS [This paper]	General graph mining	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	Details in Table 4 and Section 4

Table 1: Related work analysis, part 1: a comparison of GMS to selected existing graph-related benchmarks (“[B]”) and graph mining frameworks (“[F]”), focusing on supported graph mining problems. We exclude benchmarks only partially related to graph processing, with no focus on mining algorithms (Lonestar [57], Rodinia [65], Parboil [197], BigDataBench [212], BDGS [158], LinkBench [13], and SeBS [74]). **mC**: maximal clique listing, **kC**: k -clique listing, **dS**: densest subgraph, **sI**: subgraph isomorphism, **fS**: frequent subgraph mining, **vS**: vertex similarity, **IP**: link prediction, **cl**: clustering, **cD**: community detection, **Opt**: optimization, **Vr**: vertex rankings, ☐: Supported, ☐: Partial support, ✗: no support.

it with an accompanying software platform **P** with reference implementations of algorithms **I**. We motivate the GMS platform in Figure 1, which illustrates example performance advantages (even more than 9×) of the GMS code over a state-of-the-art variant of the Bron-Kerbosch (BK) algorithm. This shows the key benefit of the platform: it facilitates developing, redesigning, and enhancing algorithms considered in the benchmark, and thus it enabled us to rapidly obtain large speedups over fast existing BK baselines. GMS aims to propel research into different aspects of high-performance graph mining algorithms, including design, implementation, analysis, and evaluation.

To construct GMS, we first identify representative graph mining problems, algorithms, and datasets. We conduct an extensive literature review [5, 11, 63, 97, 120, 136, 137, 142, 146, 176–178, 200, 213], and obtain a benchmark specification **S** that can be used as a reference point when selecting relevant comparison targets.

Second, GMS comes with a benchmarking platform **P**: a highly modular infrastructure for easy experimenting with different design choices in a given graph mining algorithm. A key idea for high modularity is exploiting *set algebra*. Here, we observe that data structures and subroutines in many mining algorithms are “set-centric”: they can be expressed with sets and set operations, and the user can seamlessly use different implementations of the same specific “set-centric” part. This enables the user to seamlessly use new graph representations, data layouts, architectural features such as vectorization, and even use numerous graph compression schemes. We deliver ready-to-go parallel implementations of the above-mentioned elements, including parallel *reference implementations* **I** of graph mining algorithms, as well as representations, data layouts, and compression schemes. Our code is *public* and can be reused by anyone willing to use it as a basis for trying new algorithmic ideas, or simply as comparison baselines.

For more insightful performance analyses, we propose a novel *performance metric* **M** that assesses “algorithmic efficiency”, i.e., “how efficiently a given algorithm mines selected graph motifs”.

To ensure performance insights that are portable across different machines and independent of various implementation details, GMS also provides *the first extensive concurrency analysis* **C** of a

wide selection of graph mining algorithms. We use *work-depth*, an established theoretical framework from parallel computing [42, 45], to show which algorithms come with more potential for high performance on today’s massively parallel systems. Our analysis enables developers to *reduce time* spent on implementation: instead of spending days or weeks to implement an algorithm that would turn out not scalable, one can use our theoretical insights and guidelines for deciding against mounting an implementation effort.

To show the potential of GMS, we *enhance state-of-the-art algorithms* that target some of the most researched graph mining problems. This includes maximal clique listing [79], k -clique listing [78], degeneracy reordering (core decomposition) [152], and subgraph isomorphism [59, 60]. By being able to rapidly experiment with different design choices, we get *speedups of* $>9\times$, *up to* $1.1\times$, $>2\times$, and $2.5\times$, respectively. We also *improve theoretical bounds*: for example, for maximal clique listing, we obtain $O(dm3^{(2+e)d/3})$ work and $O(\log^2 n + d \log n)$ depth (d, m, n are the graph degeneracy, #edges, and #vertices, respectively). This is the best work bound among poly-logarithmic depth maximal clique listing algorithms, improving upon recent schemes [79, 91, 92].

To summarize, we provide the specific contributions:

- We propose GMS, the first benchmark for graph mining, with a specification based on more than 300 associated research papers.
- We deliver a GMS *benchmarking platform* that facilitates developing and tuning high-performance graph mining algorithms, with reference implementation of more than 40 algorithms, and high modularity obtained with set algebra, enabling experimenting with different fine- and coarse-grained algorithmic elements.
- We propose a novel *performance metric* for assessing the algorithmic throughput of graph mining algorithms.
- We support GMS with the first extensive *concurrency analysis* of graph mining for performance insights that are portable and independent of various implementation details.
- As an example of using GMS, we *enhance state-of-the-art baselines* for core graph mining problems (degeneracy, maximal clique listing, k -clique listing, and subgraph isomorphism), obtaining respective speedups of $>9\times$, up to 10%, $>2\times$, and $2.5\times$. We also *enhance their theoretical bounds*.

Reference / Infrastructure	Summary of focus (functionalities)	New Alg			Gen. APIs				Metrics				Storage				Compres.				Th.			
		∃	na	sp	N	G	S	P	rt	me	fg	mf	af	ag	bg	aa	ba	ad	of	fg	en	re	∃	nb
[B] Cyclone [201]	Graph databases	✗	✗	✗	✗	☐	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
[B] GBBS [84] + Ligra [192]	General graph processing	✗	✗	✗	☐	☐	☐	☐	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] GraphBIG [165]	General graph processing	✗	✗	✗	☐	☐	✗	✗	☐	☐	✗	☐	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
[B] GAPBS [20]	General graph processing	✗	✗	✗	☐	✗	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] Graphalytics LDBC [51]	Graph databases	✗	✗	✗	✗	✗	✗	✗	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*	☐*
[B] WGB [12]	General graph processing	✗	✗	✗	☐	✗	✗	✗	☐	☐	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
[B] PBBS [44]	General graph processing	✗	✗	✗	✗	✗	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] Graph500 [162]	Graph traversals	☐	☐	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] HPCS [15]	General graph processing	✗	✗	✗	☐	✗	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] Han et al. [106]	Evaluation of graph processing systems	✗	✗	✗	☐	✗	✗	✗	☐	☐	✗	☐	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] CRONO [6]	Multicore systems	✗	✗	✗	✗	✗	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] GARDENIA [218]	Accelerators	✗	✗	✗	☐	✗	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Arabesque [204]	Graph pattern matching	☐	☐	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] NScale [174]	Ego-network analysis	☐	☐	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] G-Thinker [219]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] G-Miner [66]	Graph pattern matching	☐	☐	☐	☐	✗	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Nuri [124]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] RStream [210]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] ASAP [116]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Fractal [86]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Kaleido [224]	Graph pattern matching	☐	☐	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] AutoMine+GraphZero [153, 154]	Graph pattern matching	☐	☐	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Pangolin [67]	Graph pattern matching	☐	✗	☐	✗	☐	✗	✗	☐	☐*	☐*	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] PrefixFPM [220]	Graph Pattern Mining	☐	✗	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[F] Peregrine [118]	Graph Pattern Mining	☐	✗	☐	✗	☐	✗	✗	☐	☐	✗	✗	✗	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐
[B] GMS [This paper]	Graph mining algorithms	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐	☐

Table 2: Related work analysis, part 2: a comparison of GMS to graph benchmarks (“[B]”) and graph pattern matching frameworks (“[F]”), focusing on supported functionalities important for developing fast and simple graph mining algorithms. We exclude benchmarks only partially related to graph processing, with no focus on mining, such as Lonestar [57], Rodinia [65], Parboil [197], BigDataBench [212], BDGS [158], LinkBench [13], and SeBS [74]. **New alg?** (∃): Are there any new/enhanced algorithms offered? **na**: do the new algorithms have provable performance properties? **sp**: are there any speedups over tuned existing baselines? **Modularity**: Is a given infrastructure modular, facilitating adding new features? The numbers (1)–(5) indicate aspects of modularity, details in Sections 3–4. In general: **Gen. APIs**: Dedicated generic APIs for a seamless integration of an arbitrary graph mining algorithm with: **N** (an arbitrary vertex neighborhood), **G** (an arbitrary graph representation), **S** (arbitrary processing stages, such as preprocessing routines), **P** (PAPI infrastructure). **Metrics**: Supported performance metrics. **rt**: (plain) run-times. **me**: (plain) memory consumption. **fg**: support for fine-grained analysis (e.g., providing run-time fraction due to preprocessing). **mf**: metrics for machine efficiency (details in § 4.3). **af**: metrics for algorithmic efficiency (details in § 4.3). **Storage**: Supported graph representations and auxiliary data structures. **ag**: graph representations based on (sparse) integer arrays (e.g., CSR). **bg**: graph representations based on (sparse or dense) bitvectors [1, 107]. **aa**: auxiliary structures based on (sparse) integer arrays. **ba**: auxiliary structures based on (sparse or dense) bitvectors. **Compression**: Supported forms of compression or space-efficient data structures. **ad**: compression of adjacency data. **of**: compression of offsets into the adjacency data. **fg**: compression of fine-grained elements (e.g., single vertex IDs). **en**: various forms of the encoding of the adjacency data (e.g., Varint [40]). **re**: support for relabeling adjacency data (e.g., degree minimizing [40]). **Th.**: Theoretical analysis. ∃: Any theoretical analysis is provided. **Nb**: Whether any new bounds (or other new theoretical results) are derived. ☐: Support. ☐: Partial support. ☐*: A given metric is supported via an external profiler. ✗: No support.

1.1 GMS vs. Graph-Related Benchmarks

We motivate GMS as the **first benchmark for graph mining**. There exist graph processing benchmarks, but they do not focus on graph mining; we illustrate this in Table 1 (“[B]”). They focus on graph *database workloads* (LDBC [51], Cyclone [201], LinkBench [13]), extreme-scale graph *traversals* (Graph500 and GreenGraph500 [162]), and different “*low-complexity*” (i.e., with run-times being low-degree polynomials in numbers of vertices or edges) parallel graph algorithms such as PageRank, triangle counting, and others, researched intensely in the parallel programming community (GAPBS [20], GBBS & Ligra [83], WGB [12], PBBS [44], HPCS [15], GraphBIG [165], Lonestar [57], Rodinia [65], Parboil [197], BigDataBench [212], BDGS [158]). Despite some similarities (e.g., GBBS provides implementations of k -clique listing), none of these benchmarks targets general graph mining, and they do not offer novel performance metrics or detailed control over graph representations, data layouts, and others. We broadly analyze this in Table 2, where we compare GMS to other benchmarks in terms of the modularity of their software infrastructures, offered metrics, control over storage schemes, support for graph compression, provided theoretical analyses, and whether they improve state-of-the-art algorithms. Finally, GMS is the only benchmark that is used

to directly enhance core state-of-the-art graph mining algorithms, achieving both better bounds and speedups in empirical evaluation.

Unlike other benchmarks, GMS proposes to exploit *set algebra* as a driving enabler for *modularity*, *simplicity*, but also *high-performance*. This design decision comes from our key observation that established formulations of many relevant graph mining problems and algorithms heavily rely on set algebra.

1.2 GMS vs. Pattern Matching Frameworks

Many graph mining *frameworks* have recently been proposed, for example Peregrine [118] and others [66, 67, 86, 116, 124, 153, 154, 204, 219, 220, 224]. GMS does *not* compete with such frameworks. First, as Table 1 shows, such frameworks do not target broad graph mining. Second, key offered functionalities also differ, see Table 2. These frameworks focus on *programming models* and *abstractions*, and on the underlying *runtime systems*¹. Contrarily, GMS focuses on benchmarking and tuning *specific parallel algorithms*, with provable performance properties, to accelerate the most competitive existing baselines.

¹We do not include these aspects in Table 2 due to space constraints – these aspects are *not* in the focus of GMS and any associated columns would have “✗” for GMS

2 NOTATION AND BASIC CONCEPTS

We first present the most basic used concepts. However, GMS touches many different areas, and – for clarity – *we will present any other background information later, when required*. Table 3 lists the most important symbols used in this work.

$G = (V, E)$	An graph G ; V, E are sets of vertices and edges.
n, m	Numbers of vertices and edges in G ; $ V = n, E = m$.
$\Delta(v), N(v)$	The degree and neighbors of $v \in V$.
Δ, \bar{d}	The maximum and the average degree in G ($\bar{d} = m/n$).

Table 3: The most important symbols used in the paper.

2.1 Graph Model

We model an undirected graph G as a tuple (V, E) ; V is a set of vertices and $E \subseteq V \times V$ is a set of edges; $|V| = n$ and $|E| = m$. The maximum degree of a graph is Δ . The neighbors and the degree of a given vertex v are denoted with $N(v)$ and $\Delta(v)$, respectively. The vertices are identified by integer IDs: $V = \{1, \dots, n\}$.

2.2 Set Algebra Concepts

GMS uses basic *set operations*: $A \cap B, A \cup B, A \setminus B, |A|$, and $\in A$. Any set operation can be implemented with various *set algorithms*. Sets usually contain vertices and at times edges. A set can be *represented* differently, for example with a bitvector or an integer array.

2.3 Graph Representation

By default, we use a standard sorted *Compressed Sparse Row (CSR) graph representation*. For an unweighted graph, CSR consists of a contiguous array with IDs of neighbors of each vertex ($2m$ words) and offsets to the neighbor data of each vertex (n words). We also use more complex representations such as compressed bitvectors.

3 OVERVIEW OF GMS

We start with an overview; see Figure 2.

The GMS **benchmark specification** **S** (details in Section 4) motivates representative graph mining problems and state-of-the-art algorithms solving these problems, relevant datasets, performance metrics **M**, and a taxonomy that structures this information. The specification, in its entirety or in a selected subpart, enables choosing relevant comparison baselines and important datasets that stress different classes of algorithms.

The specification is implemented in the **benchmarking platform** **P** (details in Section 5). The platform facilitates developing and evaluating high-performance graph mining algorithms. The former is enabled by incorporating set algebra as the key driver for modularity *and* high performance. For the latter, the platform forms a processing pipeline with well-separated parts (see the bottom of Figure 2): loading the graph from I/O, constructing a graph representation **1** – **2**, optional preprocessing **3**, running selected graph algorithms **4** – **5**, **5***, and gathering data.

The **reference implementation of algorithms** **1** (details in Section 6) offers publicly available, fast, and scalable baselines that effectively use massive parallelism in today’s architectures. We implement algorithms to make their design *modular*, i.e., different

building blocks of a given algorithm, such as a preprocessing optimization, can be replaced with user-specified codes. As data movement is dominating runtimes in irregular graph computations, we also provide a large number of *storage schemes*: graph representations, data layout schemes, and graph compression. We describe selected implementations, focusing on how they achieve high performance and modularity, in Section 6.

The **concurrency analysis** **C** (details in Section 7) offers a theoretical framework to analyze performance, storage, and the associated tradeoffs. We use work and depth [42, 45] that respectively describe the total work done by all executing processors, and the length of the associated longest execution path.

In the next sections, we detail the respective parts of GMS. We will also describe in more detail example use cases, in which we show how using GMS ensures speedups over state-of-the-art baselines for k -clique listing [78] and maximal clique listing [91].

4 BENCHMARK SPECIFICATION

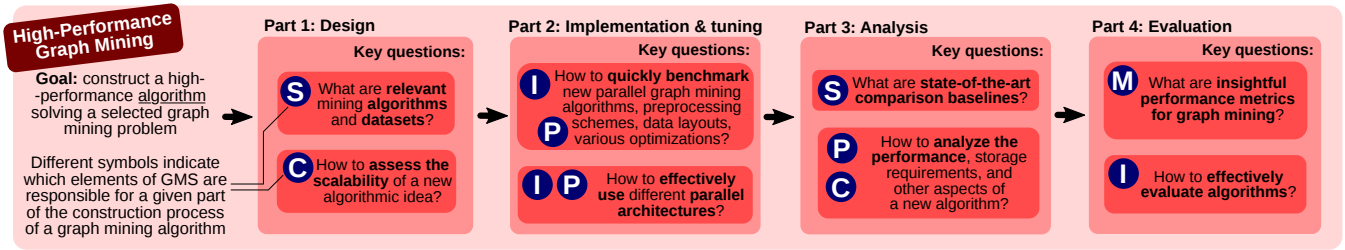
To construct a specification of graph mining algorithms, we extensively reviewed related work [5, 11, 63, 97, 120, 136, 137, 142, 146, 176–178, 200, 213]. The GMS specification has four parts: graph mining *problems, algorithms, datasets, and metrics*².

4.1 Graph Problems and Algorithms

We identify **four** major classes of graph mining problems and the corresponding algorithms: **pattern matching, learning, reordering**, and (partially) **optimization**. For each given class of problems, we aimed to cover a wide range of problems and algorithms that differ in their design and performance characteristics, for example P and NP problems, heuristics and exact schemes, algorithms with time complexities described by low-degree and high-degree polynomials, etc.. The specification is summarized in Table 4. Additional details are provided in the appendix, in Section A.

4.1.1 Graph Pattern Matching. One large class is graph pattern matching [120], which focuses on finding specific subgraphs (also called *motifs* or *graphlets*) that are often (but not always) *dense*. Most algorithms solving such problems consist of the *searching part* (finding candidate subgraphs) and the *matching part* (deciding whether a given candidate subgraph satisfies the search criteria). The search criteria (the details of the searched subgraphs) influence the time complexity of both searching and matching. First, we pick **listing all cliques** in a graph, as this problem has a long and rich history in the graph mining domain, and numerous applications. We consider both **maximal cliques** (an NP-hard problem) and **k -cliques** (a problem with time complexity in $O(n^k)$), and the established associated algorithms, most importantly Bron-Kerbosch [56], Chiba-Nishizeki [69], and their various enhancements [61, 78, 91, 151, 207]. Next, we cover a more general problem of listing **dense subgraphs** [117, 136] such as k -cores, k -star-cliques, and others. GMS also includes the Frequent Subgraph Mining (FSM) problem [120], in which one finds **all subgraphs** (not just dense) that occur *more often than a specified threshold*. Finally, we include the established NP-complete **subgraph isomorphism** (SI) problem, because of its prominence in both the

²We encourage participation in the GMS effort. If the reader would like to include some problem or algorithm in the specification and the platform, the authors would welcome the input.



Challenges & questions
 Solutions & answers

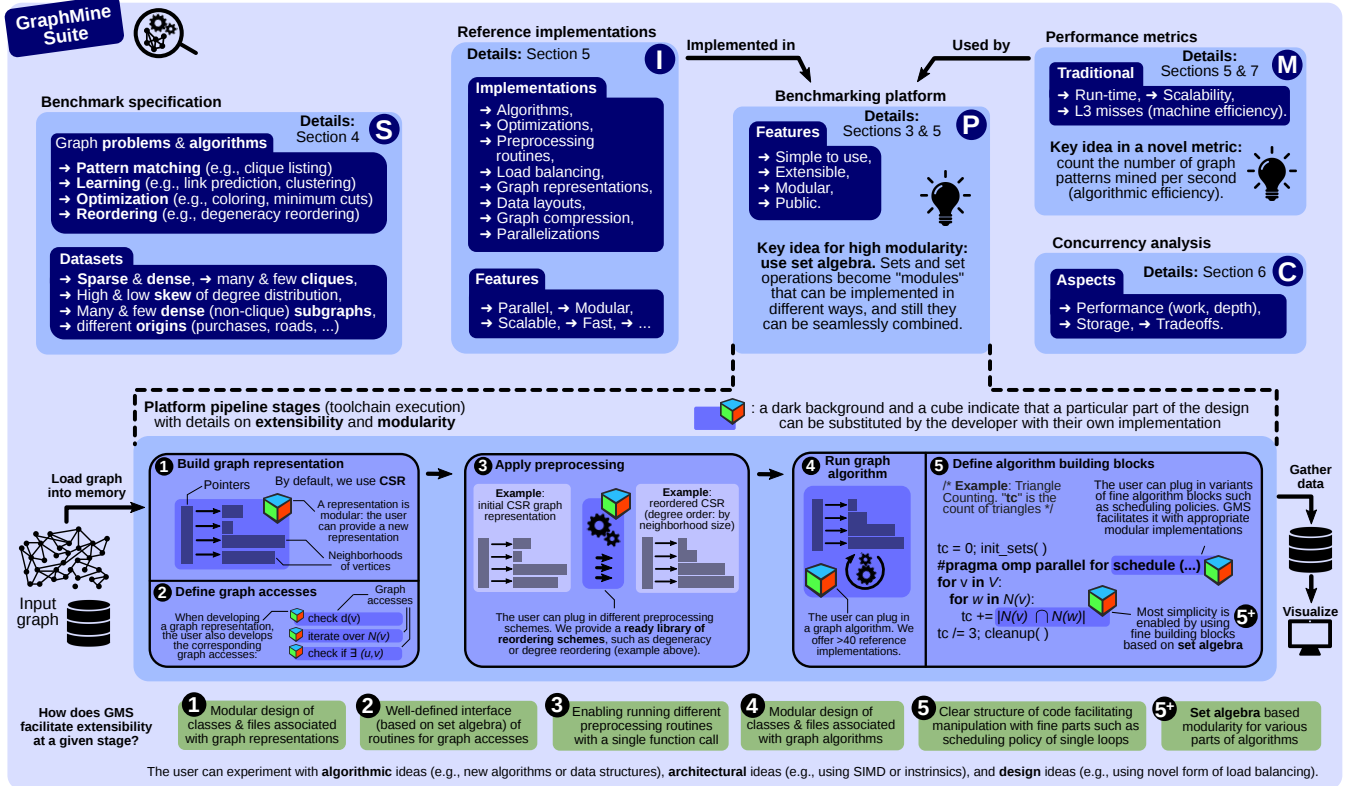


Figure 2: The overview of GMS and how it facilitates constructing, tuning, and benchmarking graph mining algorithms. The upper red part shows a process of constructing a graph mining algorithm, and the associated research questions. The middle blue part shows the corresponding different elements of the GMS suite (S – M). The bottom blue part illustrates the details of the GMS design benchmarking, with the stages of the GMS pipeline (execution toolchain) for running a given graph mining algorithm (1 – 5, 5+).

theory and practice of pattern matching, and because of a large number of variants that often have different performance characteristics [59, 75, 108, 155, 209]; SI is also used as a subroutine in the matching part of FSM.

4.1.2 *Graph Learning.* We also consider various problems that can be loosely categorized as graph learning. These problems are mostly related to clustering, and they include **vertex similarity** [137, 179, 179] (verifying how similar two vertices are), **link prediction** [10, 142, 146, 202, 211] (predicting whether two non-adjacent vertices can become connected in the future, often based on vertex similarity scores), and **Clustering and Community Detection** [46, 119, 175] (finding various densely connected groups of vertices, also often incorporating vertex similarity as a subroutine).

4.1.3 *Vertex Reordering.* We also consider reordering of vertices. Intuitively, the order in which vertices are processed in some algorithm may impact the performance of this algorithm. For example, when counting triangles, ordering vertices by degrees (prior to counting) minimizes the number of times one triangle is (unnecessarily) counted more than once. In GMS, we first consider the above-mentioned **degree ordering**. We also provide two algorithms for the **degeneracy ordering** [94] (**exact** and **approximate**), which was shown to improve the performance of maximal clique listing or graph coloring [24, 61, 91, 207].

4.1.4 *Optimization.* While GMS focuses less on optimization problems, we also include a representative problem of graph coloring and selected other problems.

Graph problem	Corresponding algorithms	E.?	P.?	Why included, what represents? (selected remarks)	
Graph Pattern Matching	• Maximal Clique Listing [87]	Bron-Kerbosch [56] + optimizations (e.g., pivoting) [61, 91, 207]	👍👉	👎	Widely used, NP-complete, example of backtracking
	• k -Clique Listing [78]	Edge-Parallel and Vertex-Parallel general algorithms [78], different variants of Triangle Counting [184, 193]	👍👉	👎	P (high-degree polynomial), example of backtracking
	• Dense Subgraph Discovery [5]	Listing k -clique-stars [117] and k -cores [94] (exact & approximate)	👍👉	👎	Different relaxations of clique mining
	• Subgraph isomorphism [87]	VF2 [75], TurboISO [108], Glasgow [155], VF3 [58, 60], VF3-Light [59]	👍	👎	Induced vs. non-induced, and backtracking vs. indexing schemes
	• Frequent Subgraph Mining [5]	BFS and DFS exploration strategies, different isomorphism kernels	👍	👎	Useful when one is interested in many different motifs
Graph Learning	• Vertex similarity [137]	Jaccard, Overlap, Adamic Adar, Resource Allocation, Common Neighbors, Preferential Attachment, Total Neighbors [179]	👍👉	👎	A building block of many more complex schemes, different methods have different performance properties
	• Link Prediction [202]	Variants based on vertex similarity (see above) [10, 142, 146, 202], a scheme for assessing link prediction accuracy [211]	👍👉	👎	A very common problem in social network analysis
	• Clustering [183]	Jarvis-Patrick clustering [119] based on different vertex similarity measures (see above) [10, 142, 146, 202]	👍👉	👎	A very common problem in general data mining; the selected scheme is an example of overlapping and single-level clustering
	• Community detection	Label Propagation and Louvain Method [195]	👍	👎	Examples of convergence-based on non-overlapping clustering
Optimization problems	• Minimum Graph Coloring [168]	Jones and Plassmann's (JP) [123], Hasenplaugh et al.'s (HS) [110], Johansson's (J) [121], Barenboim's (B) [17], Elkin et al.'s (E) [90], sparse-dense decomposition (SD) [109]	👍	👎	NP-complete; uses vertex prioritization (JP, HS), random palettes (J, B), and adapted distributed schemes (E, SD)
	• Minimum Spanning Tree [76]	Boruvka [53]	👍	👎	P (low complexity problem)
	• Minimum Cut [76]	A recent augmentation of Karger-Stein Algorithm [125]	👍	👎	P (superlinear problem)
Vertex Ordering	• Degree reordering	A straightforward integer parallel sort	👍	👉	A simple scheme that was shown to bring speedups
	• Triangle count ranking	Computing triangle counts per vertex	👍👉	👉	Ranking vertices based on their clustering coefficient
	• Degeneracy reordering	Exact and approximate [94] [127]	👍👉	👉	Often used to accelerate Bron-Kerbosch and others

Table 4: Graph problems and algorithms considered in GMS. “E.?” (Extensibility) indicates how extensible given implementations are in the GMS benchmarking platform: “👉” indicates full extensibility, including the possibility to provide new building blocks based on set algebra (👉 – 📍). “👍”: an algorithm that does not straightforwardly (or extensively) use set algebra, offering modularity levels (👍 – 📍). “P.?” (Preprocessing) indicates whether a given algorithm can be seamlessly used as a preprocessing routine; in the current GMS version, this feature is reserved for the vertex reordering algorithms.

4.1.5 Taxonomy and Discussion. Graph pattern matching, clustering, and optimization are related in that the problems from these classes focus on *finding certain subgraphs*. In the two former classes, such subgraphs are usually “local” groups of vertices, most often dense (e.g., cliques, clusters) [2–4, 23, 115, 169, 206], but sometimes can also be sparse (e.g., in FSM or SI). In optimization, a subgraph to be found can be “global”, scattered over the whole graph (e.g., vertices with the same color). Moreover, *clustering and community detection (central problems in graph learning) are similar to dense subgraph discovery (a central problem in graph pattern matching)*. Yet, the latter use the notion of *absolute density*: a dense subgraph S is some relaxation of a clique (i.e., one does not consider what is “outside S ”). Contrarily, the former use a concept of *relative density*: one compares different subgraphs to decide which one is dense [5].

4.2 Graph Datasets

We aim at a dataset selection that is computationally challenging for all considered problems and algorithms, cf. Table 4. We list both *large* and *small* graphs, to indicate datasets that can stress both low-complexity graph mining algorithms (e.g., centrality schemes or clustering) and high-complexity P, NP-complete, and NP-hard ones such as subgraph isomorphism.

So far, existing performance analyses on parallel graph algorithms focused on graphs with varying *sparsities* m/n (sparse and dense), *skews* in degree distribution (high and low skew), *diameters* (high and low), and *amounts of locality* that can be intuitively explained as the *number of inter-cluster* edges (many and few) [20]. In GMS, we recommend to use such graphs as well, as the above properties influence the runtimes of all described algorithms.

In Table 4, graphs with high degree distribution skews are indicated with large (relatively to n) maximum degrees Δ , which poses challenges for load balancing and others. Moreover, we list graphs with very high diameters (e.g., road networks) that stress iterative

algorithms where the runtime depends on the diameter. Next, to provide even more variability in the performance effects, we also consider graphs with relatively high diameters *and* with high skews in degree distributions, such as the youtube social network.

However, one of the insights that we gained with GMS is that *the higher-order structure, important for the performance of graph mining, can be little related to the above properties*. For example, in § 8.6, we describe two graphs with almost identical sizes, sparsities, and diameters, but very different performance characteristics for 4-clique mining. As we detail in § 8.6, this is because the origin of these graphs determines whether a graph has many cliques or dense (but mostly non-clique) clusters. Thus, we also explicitly recommend to use graphs of *different origins*. We provide details of this particular case in § 8.6 (cf. Livemocha and Flickr).

In addition, we explicitly consider the count of triangles T , as (1) it indicates clustering properties (and thus implies the amount of locality), and it gives hints on different higher-order characteristics (e.g., the more triangles per vertex, the higher a chance for having k -cliques for $k > 3$). Here, we also recommend using graphs that have *large differences in counts of triangles per vertex* (i.e., large T -skew). Specifically, a large difference between the *average* number of triangles per vertex T/n and the *maximum* T/n indicates that a graph may pose additional load balancing problems for algorithms that list cliques of possibly unbounded sizes, for example Bron-Kerbosch. We also consider such graphs, see Table 4.

Finally, GMS enables using synthetic graphs with the random uniform (the Erdős-Rényi model [93]) and power-law (the Kronecker model [139]) degree distributions. This is enabled by integrating the GMS platform with existing graph generators [20]. Using such synthetic graphs enables analyzing performance effects while *systematically* changing a specific *single* graph property such as n , m , or m/n , which is not possible with real-world datasets.

We stress that we refrain from prescribing concrete datasets as benchmarking input (1) for flexibility, (2) because the datasets themselves evolve and (3) the compute and memory capacities of architectures grow continually, making it impractical to stick to a fixed-sized dataset. Instead, in GMS, we analyze and discuss publicly available datasets in Section 8, making suggestions on their applicability for stressing performance of different algorithms.

4.3 Metrics

In GMS, we first use simple *running times* of algorithms (or their specific parts, for a *fine grained analysis*). Unless stated otherwise, we use all available CPU cores, to maximize utilization of the underlying system. We also consider *scalability analyses*, illustrating how the runtime changes with the increasing amount of parallelism (#threads). Comparison between the measured scaling behavior and the ideal speedup helps to identify potential scalability bottlenecks. Finally, we consider *memory consumption*.

We also assess the **machine-efficiency**, i.e., *how well a machine is utilized in terms of its memory bandwidth*. For this, we consider CPU core utilization, expressed with counts of stalled CPU cycles. One can measure this number easily with, for example, the established PAPI infrastructure [161] that enables gathering detailed performance data from hardware counters. As we will discuss in detail in Section 5, we seamlessly integrate GMS with PAPI, enabling gathering detailed data such as stalled CPU cycles but also more than that, for example cache misses and hits (L1, L2, L3, data vs. instruction, TLB), memory reads/writes, and many others.

Finally, we propose a new metric for measuring the **“algorithmic efficiency”** (**“algorithmic throughput”**). Specifically, we measure the *number of mined graph patterns in a time unit*. Intuitively, this metric indicates how efficient a given algorithm is in finding respective graph elements. An example such metric used in the past is *processed edges per second* (PEPS), used in the context of graph traversals and PageRank [143]. Here, we extend it to graph mining and to arbitrary graph patterns. In graph pattern matching, this metric is the number of the respective *graph subgraphs found per second* (e.g., maximal cliques per second). In graph learning, it is a count of vertex pairs with similarity derived per second (vertex similarity, link prediction), or the number of clusters/communities found per second (clustering, community detection). The algorithmic efficiency facilitates deriving performance insights associated with the structure of the processed graphs. By comparing relative throughput differences between different algorithms *for different input graphs*, one can conclude whether these differences consistently depend on *pattern* (e.g., *clique*) *density*.

The algorithmic efficiency metric may also be used to provide more compact results. As an example, consider two datasets, one – G_1 – with many small cliques, the other – G_2 – with few large cliques. Bron Kersbosch may be similar in both cases in its runtime, but its “clique efficiency” would be high for G_1 and low for G_2 . Thus, one could deduce based purely on the “clique throughput” that the best choice of algorithm depends on the number of cliques in the graph, because BK’s throughput suffers more when there are few cliques, but it has a high throughput when there are many of cliques. This cannot be deduced based purely on the run-time, but only using a combination of run-times *and* total clique counts.


4.4 Beyond The Scope of GMS

We fix GMB’s scope to include problems and algorithms related to “graph mining”, often also referred to as “graph analytics”, in the *offline (static)* setting, with a *single* input graph. Thus, we do not focus on streaming or dynamic graphs (as they usually come with vastly different design and implementation challenges [28]) and we do not consider problems that operate on multiple *different* input graphs. We leave these two domains for future work.

GMS also does not aim to cover advanced statistical methods that – for example – analyze power laws in input graphs. For this, we recommend to use specialized software, for example iGraph [77].

Finally, we also do not focus on many graph problems and algorithms traditionally researched in the *parallel programming* community and usually do not considered as part of graph *mining*. Examples are PageRank [167], Breadth-First Search [19], Betweenness Centrality [54, 148, 173, 194], and others [24, 27, 32, 36, 37, 39, 82, 100]. Many of these problems are addressed by abstractions such as vertex-centric [150], edge-centric [181], GraphBLAS [126] and the associated linear algebraic paradigm [126] with fundamental operations being matrix-matrix and matrix-vector products [35, 132, 133]. These works were addressed in detail in past analyses [26] and are included in existing suites such as GAPBS [20], Graph500 [162, 198], and GBBS [83]. Still, all the GMS modularity levels (1 – 5) can be used to extend the GMS platform with any of such algorithms.

5 GMS PLATFORM & SET ALGEBRA

We now detail the GMS platform and how it enables modularity, extensibility, and high performance. Details of using the platform are described in an extensive documentation (available at the provided link). There are six main ways in which one can experiment with a graph mining algorithm using the GMS platform, indicated in Figure 2 with 1 – 5 and a block .

First, the user can provide a *new graph representation* 1 and the associated *routines for accessing the graph structure* 2. By default, GMS uses CSR. A seamless integration of a new graph representation is enabled by a modular design of files and classes with the representation code, and a concise interface (checking the degree $d(v)$, loading neighbors $N(v)$, iterating over vertices V or edges E , and verifying if an edge (u, v) exists) between a representation and the rest of GMS. The GMS platform also supports compressed graph representations. While many compression schemes focus on minimizing the amount of used storage [48] and require expensive decompression, some graph compression techniques entail *mild* decompression overheads, and they can even lead to overall *speedups* due to lower pressure on the memory subsystem [40]. Here, we offer ready-to-go implementations of such schemes, including bit packing, vertex relabeling, $\text{Log}(\text{Graph})$ [40], and others.

Second, the user can seamlessly add *preprocessing routines* 3 such as the reordering of vertices. Here, the main motivation is that by applying a relevant vertex reordering (relabeling), one can reduce the amount of work to be done in the actual following graph mining algorithm. For example, the degeneracy order can significantly reduce the work done when listing maximal cliques [94]. The user runs a selected preprocessing scheme with a single function call that takes as its argument a graph to be processed.

Third, one can plug in a *whole new graph algorithm* ④, thanks to a simple code structure and easy access to methods for loading a graph from file, building representations, etc.. GMS also facilitates modifying *fine parts of an algorithm* ⑤, such as a scheduling policy of a loop. For this, we ensure a *modular structure of the respective implementations, and annotate code*.

Finally, we use the fact that many graph algorithms, for example Bron-Kerbosch [56] and others [1, 59–61, 78, 79, 91, 91, 107, 207, 211], are formulated with *set algebra* and use a small group of well-defined operations such as set intersection \cap . In GMS, we enable the user to provide their own implementation of such operations and of the data layout of the associated sets. This facilitates controlling the layout of a single auxiliary data structure or an implementation of a particular subroutine (indicated with ⑤+). Thus, one is able to break complex graph mining algorithms into simple building blocks, and work on these building blocks independently. We already implemented a wide selection of routines for \cap , \cup , \setminus , $|\cdot|$, and \in ; we also offer different set layouts based on integer arrays, bit vectors, and compressed variants of these two.

Set algebra building blocks in GMS are sets, set operations, set elements, and set algebra based graph representations. The first three are grouped together in the Set interface. The last one is a separate class that appropriately combines the instances of a given Set implementation. We now detail each of these parts.

5.1 Set Interface

The Set interface, illustrated in Listing 1, encapsulates the representation of an arbitrary set and its elements, and the corresponding set algorithms. By default, set elements are vertex IDs (modeled as integers) but other elements (i.e., integer tuples to model edges) can also be used. Then, there are three types of methods in Set.

First, there are methods implementing set basic set algebra operations, i.e., “union” for \cup , “intersect” for \cap , and “diff” for \setminus . To enable performance tuning, they come in variants. “_inplace” indicates that the calling object is being modified, as opposed to the default method variant that returns a new set (avoiding excessive data copying). “_count” indicates that the result is the size of the resulting set, e.g., $|A \cap B|$ instead of $A \cap B$ (avoiding creating unnecessary structures). Then, add and remove enable devising optimized variants of \cup and \setminus in which only one set element is inserted or removed from a set; these methods always modify the calling set.

GMS offers other methods for performance tuning. This includes constructors (e.g., a move constructor, a constructor of a single-element set, or constructors from an array, a vector, or an initializer list), and general methods such as clone, which is used because – by default – the copy constructor is disabled for sets to avoid accidental data copying. GMS also offers conversion of a set to an integer array to facilitate using established parallelization techniques.

5.2 Implementations of Sets & Set Algorithms

On one hand, a set A can be represented as a contiguous sparse array with integers modeling vertex IDs (“sparse” indicates that only non-zero elements are explicitly stored), of size $W \cdot |A|$, where W is the memory word size [bits]. This representation is commonly used to store vertex neighborhoods. However, one can also represent A with a dense bitvector of size n [bits], where the i -th set bit

```

1 // Set: a type for arbitrary sets.
2 // SetElement: a type for arbitrary set elements.
3
4 class Set {
5 public:
6 //In methods below, we denote "*"this" pointer with A
7 //(1) Set algebra methods:
8 Set diff(const Set &B) const; //Return a new set C=A\B
9 Set diff(SetElement b) const; //Return a new set C=A\b
10 void diff_inplace(const Set &B); //Update A=A\B
11 void diff_inplace(SetElement b); //Update A=A\b
12 Set intersect(const Set &B) const; //Return a new set C=A\cap B
13 size_t intersect_count(const Set &B) const; //Return |A\cap B|
14 void intersect_inplace(const Set &B); //Update A=A\cap B
15 Set union(const Set &B) const; //Return a new set C=A\cup B
16 Set union(SetElement b) const; //Return a new set C=A\cup b
17 Set union_count(const Set &B) const; //Return |A\cup B|
18 void union_inplace(const Set &B); //Update A=A\cup B
19 void union_inplace(SetElement b); //Update A=A\cup b
20 bool contains(SetElement b) const; //Return b in A ? true: false
21 void add(SetElement b); //Update A=A\cup b
22 void remove(SetElement b); //Update A=A\b
23 size_t cardinality() const; //Return set's cardinality
24
25 //(2) Constructors (selected):
26 Set(const SetElement *start, size_t count); //From an array
27 Set(std::vector<SetElement> &vec); //From a vector
28 //Set initialization with initializer list of elements:
29 Set(std::initializer_list<SetElement> &data);
30 Set(); Set(Set &&); //Default and Move constructors
31 Set(SetElement); //Constructor of a single-element set
32 static Set Range(int bound); //Create set {0,1,...,bound-1}
33
34 //(3) Other methods:
35 begin() const; //Return iterators to set's start
36 end() const; //Return iterators to set's end
37 Set clone() const; //Return a copy of the set
38 void toArray(int32_t *array) const; //Convert set to array
39 operator==; operator!=; //Set equality/inequality comparison
40
41 private:
42 using SetElement = GMS::NodeId; //(4) Define a set element
43 }

```

Algorithm 1: The set algebra interface provided by GMS.

means that a vertex $i \in A$ (“dense” indicates that all zero bits are explicitly stored). While being usually larger than a sparse array, a dense bitvector is more space-efficient when A is very large, which happens when some vertex connects to the majority of all vertices. Now, depending on A ’s and B ’s representations, $A \cap B$ can itself be implemented with different set algorithms. For example, if A and B are sorted sparse arrays with similar sizes ($|A| \approx |B|$), one prefers the “merge” scheme where one simply iterates through A and B , identifying common elements (taking $O(|A| + |B|)$ time). If one set (e.g., B) is represented as a bitvector, one may prefer a scheme where one iterates over the elements of a sparse array A and checks if each element is in B , which takes $O(1)$ time, giving the total of $O(|A|)$ time for the whole intersection.

Moreover, a bitvector enables insertion or deletion of vertices into a set in $O(1)$ time, which is useful in algorithms that rely on dynamic sets, for example Bron-Kerbosch [61, 79, 91, 207]. There are more set representations with other performance characteristics, such as sparse [1, 107] or compressed [34] bitvectors, or hashtables, enabling further performance/storage tradeoffs.

Importantly, using different *set representations* or *set algorithms* does *not* impact the formulations of *graph algorithms*. GMS exploits this fact to facilitate development and experimentation.

By default, GMS offers three implementations of Set interface:

- **RoaringSet** A set is implemented with a bitmap compressed using recent “roaring bitmaps” [64, 138]. A roaring bitmap offers diverse compression forms within the same bitvector. They offer mild

compression rates but do not incur expensive decompression. As we later show, these structures result in high performance of graph mining algorithms running on top of them.

- **SortedSet** GMS also offers sets stored as sorted vectors. This reflects the established CSR graph representation design, where each neighborhood is a sorted contiguous array of integers.
- **HashSet** Finally, GMS offers an implementation of Set with a hashtable. By default, we use the Robin Hood library [62].

5.3 Set-Centric Graph Representations

Sets are building blocks for a graph representation: one set implements one neighborhood. To enable using arbitrary set designs, GMS harnesses templates, typed by the used set definition, see Listing 2. GMS provides ready-to-go representations based on the RoaringSet, SortedSet, and HashSet set representations.

```

1 template <class TSet>
2 class SetGraph {
3 public:
4     using Set = TSet; int64_t num_nodes() const;
5     const Set& out_neigh(NodeId node) const;
6     int64_t out_degree(NodeId node) const;
7     /* Some functions omitted */;

```

Algorithm 2: A generic graph representation.

5.4 Pipeline Interface

Beyond the set algebra related interfaces, GMS also offers a dedicated API for easy experimenting with other parts of the processing pipeline (1 – 5). This API is illustrated in Listing 3. It enables separate testing of each particular stage, but also enables the user to define and analyze their own specific stages.

```

1 class MyPipeline : public GMS::Pipeline {
2 public:
3     //Any benchmark-specific arguments, including the input graph,
4     //are passed to the constructor
5     MyPipeline(const GMS::CLI::Args &a, const SortedSetGraph &g);
6     //Functions for the individual steps.
7     void convert(); //Potential conversion of g to another format
8     void preprocess(); //Needed preprocessing
9     void kernel(); //Desired graph mining algorithm
10 private:
11     /* Any state variables that are shared between steps */;

```

Algorithm 3: A generic graph representation.

5.5 PAPI Interface

GMS also uses the PAPI library for easy access to hardware performance counters, cf. § 4.3. Importantly, we support seamless gathering of the performance data from *parallel* code regions³ An example usage of PAPI in GMS is in Listing 4. All the details on how to use the GMS PAPI support are also available in the online documentation.

```

1 //Init PAPI for parallel use, measure CPU cycles
2 //stalled on memory accesses, and on any resources
3 GMS::PAPIW::INIT_PARALLEL(PAPI_MEM_SCY, PAPI_RES_STL);
4 GMS::PAPIW::START();
5 #pragma omp parallel
6 {
7     //Benchmarked parallel region
8 }
9 GMS::PAPIW::STOP();

```

Algorithm 4: Using PAPI for detailed performance measurements of a parallel region in GMS.

³We currently support OpenMP and plan to include other infrastructures such as Intel TBB.

6 HIGH-PERFORMANCE & SIMPLICITY

We now detail how using the GMS benchmarking platform leads to simple (i.e., programmable) *and* high-performance implementations of many graph mining algorithms.

We now use the GMS benchmarking platform to enhance existing graph mining algorithms. We provide consistent speedups (detailed in Section 8). Some new schemes also come with theoretical advancements (detailed in Section 7). The following descriptions focus on (1) how we ensure the *modularity* of GMS algorithms (for programmability), and (2) what GMS design choices ensure *speedups*. Selected modular parts are marked with the blue color and the type of modularity (1 – 5) (explained in § 3 and Figure 2). Marked set operations are implemented using the Set interface, see Listing 1. Whenever we use parallelization (“in parallel”), we ensure that it does not involve conflicting memory accesses. For clarity, we focus on *formulations* and we discuss implementation details (e.g., parallelization) in the next sections.

6.1 Use Case 1: Degeneracy Order & k -Cores

A *degeneracy* of a graph G is the smallest d such that every subgraph in G has a vertex of degree at most d . Thus, degeneracy can serve as a way to measure the graph sparsity that is “closed under taking a graph subgraph” (and thus more robust than, for example, the average degree). A *degeneracy ordering* (DGR) is an “ordering of vertices of G such that each vertex has d or fewer neighbors that come later in this ordering” [91]. DGR can be obtained by repeatedly removing a vertex of minimum degree in a graph. The derived DGR can be directly used to compute the k -core of G (a maximal connected subgraph of G whose all vertices have degree at least k). This is done by iterating over vertices in the DGR order, and removing vertices with out-degree less than k .

DGR, when used as a preprocessing routine, has been shown to accelerate different algorithms such as Bron-Kerbosch [91]. In the GMS benchmarking platform, we provide an implementation of DGR that is modular and can be seamlessly used with other graph algorithms as preprocessing (3). Moreover, we alleviate the fact that the default DGR is not easily parallelizable and takes $O(n)$ iterations even in a parallel setting. For this, GMS delivers a modular implementation of a recent $(2 + \epsilon)$ -approximate degeneracy order [24] (ADG), which has $O(\log n)$ iterations for any $\epsilon > 0$. Specifically, the strict degeneracy ordering can be relaxed by introducing the approximation (multiplicative) factor k that determines, for each vertex v , the *additional number of neighbors that can be ranked higher in the order than v* . Formally, in a k -approximate degeneracy ordering, every vertex v can have at most $k \cdot d$ neighbors ranked higher in this order. Deriving ADG is in Algorithm 5. It is similar to computing the DGR, which iteratively removes vertices of the smallest degree. The main difference is that one removes in parallel a *batch* of vertices with degrees smaller than $(1 + \epsilon)\widehat{d}_U$ (cf. set R and Line 7). The parameter $\epsilon \geq 0$ controls the accuracy of the approximation; \widehat{d}_U is the average degree in the induced subgraph $G(U, E[U])$, U is a “working set” that tracks changes to V . ADG relies on set cardinality and set difference, enabling the GMS set algebra modularity (5).

```

1 //Input: A graph G 1. Output: Approx. degeneracy order (ADG)  $\eta$ .
2 i = 1 // Iteration counter
3  $U = V$  //  $U$  is the induced subgraph used in each iteration  $i$ 
4 while  $U \neq \emptyset$  do:
5    $\widehat{d}_U = \left( \sum_{v \in U} |N_U(v)| \right) / |U|$  //Get the average degree in  $U$ 
6   //  $R$  contains vertices assigned priority in this iteration:
7    $R = \{v \in U : |N_U(v)| \leq (1 + \epsilon)\widehat{d}_U\}$ 
8   for  $v \in R$  in parallel 2 5 do:  $\eta(v) = i$  //assign the ADG order
9    $U = U \setminus R$  5 //Remove assigned vertices
10  i = i + 1

```

Algorithm 5: Deriving the approximate degeneracy order (ADG) in GMS. More than one number indicates that a given snippet is associated with more than one modularity type.

6.2 Use Case 2: Maximal Clique Listing

Maximal clique listing, in which one enumerates all *maximal* cliques (i.e., fully-connected subgraphs not contained in a larger such subgraph) in a graph, is one of core graph mining problems [61, 69, 79, 80, 88, 122, 129, 130, 141, 145, 149, 166, 186, 196, 199, 208, 214, 217, 223]. The recursive backtracking algorithm by Bron and Kerbosch (BK) [56] together with a series of enhancements [79, 91, 92, 207] (see Algorithm 6) is an established and, in practice, the most efficient way of solving this problem. Intuitively, in BK, one iteratively considers each vertex v in a given graph, and searches for all maximal cliques that contain v . The search process is conducted *recursively*, by starting with a single-vertex clique $\{v\}$, and augmenting it with v 's neighbors, one at a time, until a maximal clique is found. Still, the number of maximal cliques in a general graph, and thus BK's runtime, may be exponential [159].

Importantly, the order in which all the vertices are selected for processing (at the *outermost* level of recursion) may heavily impact the amount of work in the following iterations [79, 91, 92]. Thus, in GMS, we use different vertex orderings, integrated using the GMS preprocessing modularity **(3)**. One of our core enhancements is to use the ADG order (§ 6.1). As we will show, this brings theoretical (Section 7) and empirical (Section 8) advancements.

A key part are vertex sets P , X , and R . They together navigate the way in which the recursive search is conducted. P (“Potential”) contains candidate vertices that will be considered for belonging to the clique currently being expanded. X (“eXcluded”) are the vertices that are definitely *not* to be included in the current clique (X is maintained to avoid outputting the same clique more than once). R is a currently considered clique (may be non-maximal). In GMS, we extensively experimented with different set representations for P , X , and R , which was facilitated by the set algebra based modularity **(5)**. Our goal was to use representations that enable fast “bulk” set operations such as intersecting large sets (e.g., $X \cap N(v)$ in Line 23) but also efficient fine-grained modifications of such sets (e.g., $X = X \cup \{v\}$ in Line 28). For this, we use roaring bitmaps. As we will show (Section 8), using such bitvectors as representations of P , X , and R brings overall speedups of even more than $9\times$.

Now, at the outermost recursion level, for each vertex v_i , we have $R = \{v_i\}$ (Line 13). This means that the considered clique starts with v_i . Then, we have $P = N(v_i) \cap \{v_{i+1}, \dots, v_n\}$ and $X = N(v_i) \cap \{v_1, \dots, v_{i-1}\}$. This removes unnecessary vertices from P and X . As we proceed in a fixed order of vertices in the main loop, when starting a recursive search for $\{v_i\}$, we will definitely *not* include vertices $\{v_1, \dots, v_{i-1}\}$ in P , and thus we can limit P to $N(v_i) \cap \{v_{i+1}, \dots, v_n\}$ (a similar argument applies to R). Note that

these intersections may be implemented as simple splitting of the neighbors $N(v_i)$ into two sets, based on the vertex order. This is another example of the decoupling of general simple set algebraic formulations in GMS and the underlying implementations **(5)**.

In each recursive call of BK-Pivot, each vertex from P is added to R to create a new clique candidate R_{new} explored in the following recursive call. In this recursive call, P and X are respectively restricted to $P \cap N(v)$ and $X \cap N(v)$ (any other vertices besides $N(v)$ would not belong to the clique R_{new} anyway). After the recursive call returns, v is moved from P (as it was already considered) to X (to avoid redundant work in the future). The key condition for checking if R is a *maximal* clique is $P \cup X == \emptyset$. If this is true, then no more vertices can be added to R (including the ones from X that were already considered in the past) and thus R is maximal.

The BK variant in GMS also includes an additional important optimization called *pivoting* [207]. Here, for any vertex $u \in P \cup X$, only u and its *non* neighbors (i.e., $P \setminus N(u)$) need to be tested as candidates to be added to P . This is because any potential maximal clique must contain *either* u *or* one of its non-neighbors. Otherwise, a potential clique could be enlarged by adding u to it. Thus, when selecting u (Line 20), one may use any scheme that *minimizes* $|P \setminus N(u)|$ [207]. The advantage of pivoting is that it further prunes the search space and thus limits the number of recursive calls.

For further performance improvements, we also use roaring bitmaps to implement graph neighborhoods, exploiting the GMS modularity of representations and set algebra **(1, 2, 5)**.

An established way to derive the pivot vertex $u \in P \cup X$, introduced by Tomita et al. [207], is to find $u = \operatorname{argmin}_{v \in P \cup X} |P \cap N(v)|$. This approach minimizes the size of P before the associated recursive BK-Pivot call. Yet, it comes with a computational burden, because – to select u – one must conduct the set operation $|P \cap N(v)|$ as many as $|P \cup X|$ times. This issue was addressed by proposing to derive $|P \cap N_H(v)|$ instead of $|P \cap N(v)|$, where H is an induced subgraph of G , with the vertex set $P \cup X$ and the edge set $\{\{x, y\} \in E \mid x \in P \wedge y \in P \cup X\}$ [92]. Using $N_H(v)$ reduces the amount of work in each $|P \cap N_H(v)|$, because $N_H(v)$ is smaller than $N(v)$. Such subgraph H is *precomputed* before choosing u , and is then passed to the recursive BK-Pivot call, to accelerate precomputing subgraphs H at deeper recursion levels.

We observe that the precomputed subgraph H can be used not only to accelerate pivot selection, but also in several other set operations. First, one can use $P \setminus N_H(u)$ instead of $P \setminus N(u)$ to reduce the cost of set difference; note that this does not introduce more iterations in the following loop because no vertex in P is included in $N(u) \setminus N_H(u)$. Second, we can also use H to compute $P \cap N_H(v)$ and $X \cap N_H(v)$ instead of $P \cap N(v)$ and $X \cap N(v)$, also reducing the amount of work in set intersections.

We also investigated the impact of constructing the H subgraphs on each recursion level, as initially advocated [92], versus only at the outermost level. We observe that, while the latter always offers performance advantages due to the large reductions in work, the former often introduces overheads that outweigh gains, due to the memory cost (from maintaining many additional subgraphs) and the increase in runtimes (from constructing many subgraphs). In our final BK-ADG version, we only derive H at the outermost loop

iteration, once for each vertex v , and use a given H at each level of the search tree associated with v .

We also developed a variant of BK-ADG that, similarly to BK-DAS, uses nested parallelism *at each level of recursion*. This approach proved consistently slower than the version without this feature.

```

1 /* Input: A graph  $G$  1, Output: all maximal cliques. */
2
3 //Preprocessing: reorder vertices with DGR or ADG; see § 6.1.
4  $(v_1, v_2, \dots, v_n) = \text{preprocess}(V, /* \text{selected vertex order} */)$  3
5
6 //Main part: conduct the actual clique enumeration.
7 for  $v_i \in (v_1, v_2, \dots, v_n)$  do: //Iterate over  $V$  in a specified order
8 //For each vertex  $v_i$ , find maximal cliques containing  $v_i$ .
9 //First, remove unnecessary vertices from  $P$  (candidates
10 //to be included in a clique) and  $X$  (vertices definitely
11 //not being in a clique) by intersecting  $N(v_i)$  with vertices
12 //that follow and precede  $v_i$  in the applied order.
13  $P = N(v_i) \cap \{v_{i+1}, \dots, v_n\}$  5;  $X = N(v_i) \cap \{v_1, \dots, v_{i-1}\}$  5;  $R = \{v_i\}$ 
14
15 //Run the Bron-Kerbosch routine recursively for  $P$  and  $X$ .
16 BK-Pivot( $P, \{v_i\}, X$ )
17
18 BK-Pivot( $P, R, X$ ) //Definition of the recursive BK scheme
19 if  $P \cup X == \emptyset$  5: Output  $R$  as a maximal clique
20  $u = \text{pivot}(P \cup X)$  5 //Choose a "pivot" vertex  $u \in P \cup X$ 
21 for  $v \in P \setminus N(u)$  5: // Use the pivot to prune search space
22 //New candidates for the recursive search
23  $P_{\text{new}} = P \cap N(v)$  5;  $X_{\text{new}} = X \cap N(v)$  5;  $R_{\text{new}} = R \cup \{v\}$  5
24 //Search recursively for a maximal clique that contains  $v$ 
25 BK-Pivot( $P_{\text{new}}, R_{\text{new}}, X_{\text{new}}$ )
26 //After the recursive call, update  $P$  and  $X$  to reflect
27 //the fact that  $v$  was already considered
28  $P = P \setminus \{v\}$  5;  $X = X \cup \{v\}$  5

```

Algorithm 6: Enumeration of maximal cliques, a Bron-Kerbosch variant by Epstein et al. [92] with GMS enhancements.

6.3 Use Case 3: k -Clique Listing

GMS enabled us to enhance a state-of-the-art k -clique listing algorithm [78]. Our GMS formulation is shown in Algorithm 7. We reformulated the original scheme (without changing its time complexity) to expose the implicitly used set operations (e.g., Line 18), to make the overall algorithm more modular. The algorithm uses recursive backtracking. One starts with iterating over edges (2-cliques), in Lines 11–12. In each backtracking search step, the algorithm augments the considered cliques by one vertex v and restricts the search to neighbors of v that come after v in the used vertex order.

Two schemes marked with **3** indicate two preprocessing routines that appropriately reorder vertices and – for the obtained order – assign directions to the edges of the input graph G . Both are well-known optimizations that reduce the search space size [78]. For such a modified G , we denote *out-neighbors* of any vertex u with $N^+(u)$. Then, operations marked with **5** refer to accesses to the graph structure and different set operations that can be replaced with any implementation, as long as it preserves the semantics of set membership, set cardinality, and set intersection.

The modular design and using set algebra enables us to easily experiment with different implementations of C_i , $N^+(u) \cap C_i$, and others. For example, we successfully and rapidly redesigned the reordering scheme, reducing the number of pointer chasing and the total amounts of communicated data. We investigated the generated assembly code of the respective part; it has 22 x86 mov instructions,

```

1 /*Input: A graph  $G$  1,  $k \in \mathbb{N}$  Output: Count of  $k$ -cliques  $ck \in \mathbb{N}$ . */
2
3 //Preprocessing: reorder vertices with DGR or ADG; see § 6.1.
4 //Here, we also record the actual ordering and denote it as  $\eta$ 
5  $(v_1, v_2, \dots, v_n; \eta) = \text{preprocess}(V, /* \text{selected vertex order} */)$  3
6
7 //Construct a directed version of  $G$  using  $\eta$ . This is an
8 //additional optimization to reduce the search space:
9  $G = \text{dir}(G)$  3 //An edge goes from  $v$  to  $u$  iff  $\eta(v) < \eta(u)$ 
10  $ck = 0$  //We start with zero counted cliques.
11 for  $u \in V$  in parallel do: 2 //Count  $u$ 's neighboring  $k$ -cliques
12  $C_2 = N^+(u)$ ;  $ck += \text{count}(2, G, C_2)$ 
13
14 function count( $i, G, C_i$ ):
15 if ( $i == k$ ): return  $|C_k|$  5 //Count  $k$ -cliques
16 else:
17  $ci = 0$ 
18 for  $v \in C_i$  5 do: //search within neighborhood of  $v$ 
19  $C_{i+1} = N^+(v) \cap C_i$  5 //  $C_i$  counts  $i$ -cliques.
20  $ci += \text{count}(i+1, G, C_{i+1})$ 
21 return  $ci$ 

```

Algorithm 7: k -Clique Counting; see Listing 5 for the explanation of symbols.

compared to 31 before the design enhancement⁴. Moreover, we improved the memory consumption of the algorithm. The space allocated per subgraph C_i (e.g., **5**) is now upper bounded by $|C_i|^2$ (counted in vertices) instead of the default Δ^2 . When parallelizing over edges, this significantly reduces the required memory (for large maximum degrees Δ , even up to >90%). Finally, the modular approach taken by the GMS platform enables more concise (and thus less complex) algorithm formulation. Specifically, the original version had to use a separate routine for listing cliques for $k = 3$, while the GMS's reformulation enables all variants for $k \geq 3$.

6.4 Use Case 4: Subgraph Isomorphism

GMS ensured speedups in the most recent parallel variant of the VF3 subgraph isomorphism algorithm [59, 60]. Here, the GMS platform facilitates plugging in arbitrary variants of algorithms without having to modify other parts of the toolchain (**4** – **5**) (Listing is in the extended report). First, example used optimizations in the baseline are *work splitting* combined with *work stealing*. Specifically, threads receive lists of vertices from which they start recursive backtracking. However, due to diverse graph structure, this search can take a variable amount of time (because there is more backtracking for some vertices) so some threads finish early. To combat this, we use a lockfree queue, where idling threads steal work from other threads. The queue element is the ID of a vertex from where to begin backtracking. The thread performs a compare-and-swap (CAS) atomic to retrieve a vertex from its queue. Idle threads select threads (that they steal from) uniformly at random. We also use a *precompute* scheme: during runtime, we gather information about possible mappings between vertices with their neighborhoods, and certain specific query graphs. This can accelerate, for example, searching through certain parts of the target graph.

6.5 Use Case 5: Vertex Similarity & Clustering

We include vertex similarity and clustering in GMS. Vertex similarity measures heavily use \cap . For example, the well-known Jaccard

⁴We used "compiler explorer" (<https://godbolt.org/>) for assembly analysis

k -Clique Listing Node Parallel [78]	k -Clique Listing Edge Parallel [78]	★ k -Clique Listing with ADG (§ 6.3)	ADG (Section 6)	Max. Cliques Eppstein et al. [91]	Max. Cliques Das et al. [79]	★ Max. Cliques with ADG (§ 7.3)	Subgr. Isomorphism Node Parallel [58, 75]	Link Prediction†, JP Clustering
Work $O\left(mk \left(\frac{d}{2}\right)^{k-2}\right)$	$O\left(mk \left(\frac{d}{2}\right)^{k-2}\right)$	$O\left(mk \left(d + \frac{\epsilon}{2}\right)^{k-2}\right)$	$O(m)$	$O\left(dm^3 d^{1/3}\right)$	$O\left(3^{n/3}\right)$	$O\left(dm^3^{(2+\epsilon)d/3}\right)$	$O\left(n\Delta^{k-1}\right)$	$O(m\Delta)$
Depth $O\left(n + k \left(\frac{d}{2}\right)^{k-1}\right)$	$O\left(n + k \left(\frac{d}{2}\right)^{k-2} + d^2\right)$	$O\left(k \left(d + \frac{\epsilon}{2}\right)^{k-2} + \log^2 n + d^2\right)$	$O\left(\log^2 n\right)$	$O\left(dm^3 d^{1/3}\right)$	$O(d \log n)$	$O\left(\log^2 n + d \log n\right)$	$O\left(\Delta^{k-1}\right)$	$O(\Delta)$
Space $O(nd^2 + K)$	$O(md^2 + K)$	$O(md^2 + K)$	$O(m)$	$O(m + nd + K)$	$O(m + pd\Delta + K)$	$O(m + pd\Delta + K)$	$O(m + nk + K)$	$O(m\Delta)$

Table 5: Work, depth, and space for some graph mining algorithms in GMS. d is the graph degeneracy, K is the output size, Δ is the maximum degree, p is the number of processors, k is the number of vertices in the graph that we are mining for, n is the number of vertices in the graph that we are mining, and m is the number of edges in that graph. † Link prediction and the JP clustering complexities are valid for the Jaccard, Overlap, Adamic Adar, Resource Allocation, and Common Neighbors vertex similarity measures. ★ Algorithms derived in this work. Additional bounds for BK are in Table 6

and overlap similarities of $u, v \in V$ are defined as $\frac{|N(u) \cap N(v)|}{|N(u) \cup N(v)|}$ and $\frac{|N(u) \cap N(v)|}{\min(|N(u)|, |N(v)|)}$. We provide a modular implementation in the GMS platform (5), where one can use different set representations (bitvectors, compressed bitvectors, integer arrays, others) and two different routines for \cap : (1) simple merging of sorted sets (taking $O(|N(v)| + |N(u)|)$ time) and a “galloping” variant where, for each element x from a smaller set $N(v)$, one uses binary search to check if $x \in N(u)$ (taking $O(|N(v)| \log |N(u)|)$ time). This enables fine tuning performance.

6.6 Use Case 6: k -Clique-Star Listing

A k -clique-star is a k -clique with additional neighboring vertices that are connected to all the vertices in the clique. k -clique-stars were proposed as graph motifs that relax the restrictive nature of k -cliques [117] (large cliques are expected to be rare because every vertex in a clique, regardless of the clique size, must be connected to all other vertices in this clique). Our observation is that those extra vertices that are connected to the k -clique actually form a $(k + 1)$ -clique (together with this k -clique). Thus, to find k -clique-stars, we first mine $(k + 1)$ -cliques. Then, we find k -clique-stars within each $(k + 1)$ -clique using set union, membership, and difference.

6.7 Use Case 7: Link Prediction

Here, one is interested in developing schemes for predicting whether two non-adjacent vertices can become connected in the future. There exist many schemes for such prediction [10, 142, 146, 202] and for assessing the accuracy of a specific link prediction scheme [211]. We start with some graph with *known* links (edges). We derive $E_{sparse} \subseteq E$, which is E with random links removed; $E_{sparse} =$

$E \setminus E_{rndm}$. $E_{rndm} \subseteq E$ are randomly selected *missing* links from E (links to be predicted). We have $E_{sparse} \cup E_{rndm} = E$ and $E_{sparse} \cap E_{rndm} = \emptyset$. Now, we apply the link prediction scheme S (that we want to test) to each edge $e \in (V \times V) \setminus E_{sparse}$. The higher a value $S(e)$, the more probable e is to appear in the future (according to S). Now, the effectiveness eff of S is computed by verifying how many of the edges with highest prediction scores ($E_{predict}$) actually are present in the original dataset E : $eff = |E_{predict} \cap E_{rndm}|$.

6.8 Developing Graph Representations

The right data layout is one of key enablers of high performance. There exists a plethora of graph representations, layouts, models, and compression schemes [34]. Different compression schemes may vastly differ in the compression ratio as well as the performance of accessing and mining a graph. For example, some graphs compressed with a combination of techniques implemented in the WebGraph framework [49] can use even below one bit per link. Yet, decompression overheads may significantly impact the performance of graph mining algorithms running on such compressed graphs. Then, a recent compressed Log(Graph) representation can deliver 20-35% space reductions with simple bit packing, while eliminating decompression overheads or even *delivering speedups due to reduced amounts of transferred data* [40]. Besides graph compression, there exist many other schemes related to representations, for example NUMA-awareness in graph storage [222]; they all impact performance of graph processing.

We consider the aspect of data layout and graph representation design in GMS and we enable the user to analyze *relationships between graph storage and the performance of graph algorithms*.

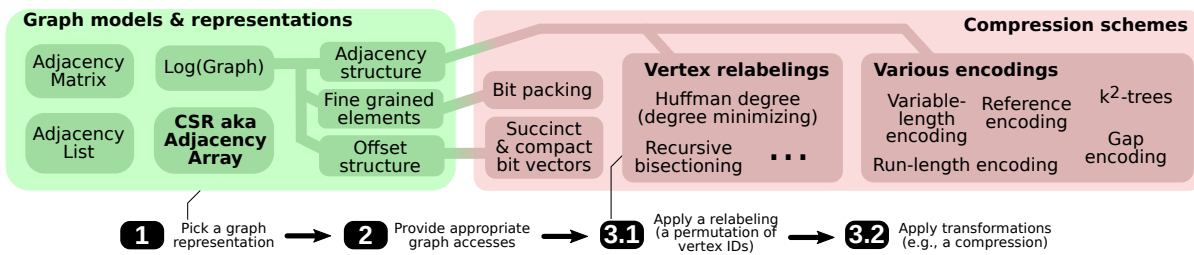


Figure 3: Selected storage schemes (graph models, representations, and graph compression methods) considered in the GMS platform. All the schemes are outlined and pictured in more detail in Figure 10 (in the Appendix) and described in detail in a recent survey [34]. Developing and using a specific representation in GMS corresponds to steps 1 – 3 in the pipelined GMS design (cf. Figure 2).

Specifically, the user can rapidly develop or use an existing storage scheme and analyze its space utilization and its impact on the performance of graph algorithms and graph queries. The considered storage schemes are illustrated in Figure 3, they include **graph models and representations** (e.g., Log(Graph) [40]), and **graph compression schemes** (e.g., difference encoding [34], k^2 -trees [55], bit packing [34], or succinct offsets [34, 101]). All these schemes offer different trade-offs between the required storage space and the delivered performance. The current version of GMS implements many of these schemes, but it also offers an intuitive and extensive interface that facilitates constructing new ones.

We provide more details of these storage schemes, and how to use them, in the Appendix (Section B). In general, in GMS, one first selects a model to be used (by default, it is the Adjacency List Model) and its specific implementation. GMS uses a simple Compressed Sparse Row (CSR) by default. Other available schemes include Log(Graph) with its bit packing of vertex IDs or succinct and compact offsets to neighborhoods [40]. Then, one must provide the implementation of graph accesses (fetching neighbors or a given vertex, checking the degree of a given vertex, verifying if a given edge exists). After that, one may apply additional preprocessing. First, one can relabel vertex IDs (i.e., apply a permutation of vertex IDs), for example the Huffman degree relabeling. Second, one may provide a transformation of each (permuted) neighborhood, for example encoding neighborhoods using Varint compression.

Further details on the permutations and transformations of vertex neighborhoods can also be found in the Log(Graph) paper [40].

7 CONCURRENCY ANALYSIS

In this part of GMS, we show how to assess a priori the properties of parallel graph mining algorithms, *reducing time spent on algorithm design and development and providing performance insights that are portable across machines that differ in certain ways (e.g., in the sizes of their caches) and independent of various implementation details*. We first broadly discuss the approach and the associated trade-offs. Second, as use cases, we pick k -clique and maximal clique listing, and we *enhance state-of-the-art algorithms addressing these problems*. Table 5 summarizes the GMS theoretical results; many of these bounds are novel.

7.1 Methodology, Models, Tools

We use the established *work-depth analysis* for bounding run-times of parallel algorithms. Here, the total number of instructions performed by an algorithm (over all number of processors for a given input size) is the *work* of the algorithm. The longest chain of sequential dependencies (for a given input size) is the *depth* of an algorithm [42, 45]. This approach is used in most recent formal analyses of parallel algorithms in the shared-memory setting [83, 111]. Overall, we consider *four* aspects of a parallel algorithm: (1) the overhead compared to a sequential counterpart, quantified with work, (2) the scalability, which is illustrated by depth, (3) the space usage, and – when applicable – (4) the approximation ratio. These four aspects often enable different tradeoffs.

7.2 Discussion On Trade-Offs

For many problems, there is a **tradeoff between work, depth, space**, and sometimes **approximation ratio** [78, 125, 157]. Which algorithm is the best choice hence depends on the available number of processors and the available main memory. For today’s shared memory machines, typically the number of processors/cores is relatively small (e.g., 18 on our machines) and main memory is not much bigger than the graphs we would like to process (e.g., 64GiB or 768GiB on our machines, see Section 8). Thus, *reducing work (and maintaining close to linear space in the input plus output)* is a high priority to obtain good performance in practice [83].

An algorithm with a work that is much larger than the best sequential algorithm will require many processors to be faster than the latter. An algorithm with large depth will stop scaling for a small number of processors. An estimate of the runtime of an algorithm with work W and depth D on p processors is $W/p + D$. This estimate is optimistic as it neglects the cost for scheduling threads and caching issues (e.g., false sharing). Yet, it has proven a useful model in developing efficient graph algorithms in practice [83].

The space used by a parallel algorithm limits the largest problem that can be solved on a fixed machine. This is crucial for graph mining problems with exponential time complexities where we want the *space to be close to the input size plus the output size*.

We illustrate a work / depth / space tradeoff with k -clique listing [78] (§ 6.3). All following designs are pareto-optimal in terms of the work / depth / space tradeoff and they are useful in different circumstances (for different machines).

First, consider a **naive** algorithm variant. Starting from every vertex, one spawns parallel recursive searches to complete the current clique. The advantage of this approach is that it has low depth $O(k)$, but the work and space is $\Theta(n\Delta^{k-1})$, which can be prohibitive.

This approach can be enhanced by using the DGR order to guide the search as described in § 6.3 (the “**Node Parallel**” variant). Here, one invokes a parallel search starting from each vertex for cliques that contain this vertex as the first vertex in the order. This reduces the space to almost linear $\Theta(nd^2)$, where d is the degeneracy of the graph. The depth is increased to $\Theta(n + k(d/2)^{k-1})$. This design was reported to have poor scalability in practice [78].

One can also invoke a parallel search for every *edge* (“**Edge Parallel**”) and try to find a clique that contains it (and follows the DGR order). The depth decreases by a factor of d to $\Theta(n + k(d/2)^{k-2} + d^2)$, but the space increases by a factor of $\frac{m}{n}$ to $O(md^2)$. This approach has a good work / depth / space tradeoff in practice [78].

7.3 Bounds for Graph Mining Algorithms

Table 5 presents work-depth and space bounds for considered graph mining algorithms. Here, we obtain *new better* bounds for maximal clique listing. The main idea is to combine existing corresponding algorithms [79, 92] with the ADG ordering. Specifically, the new maximal clique listing improves upon the Eppstein et al. [92] and Das et al. [79]: our depth is better than both while work is better than that of [92] and adds only a small factor to work in [92]. We also provide a new k -clique listing variant, again by using ADG. The variant scales better than Danisch et al. [78] (column 2) if n is

much bigger than kd^{k-2} . This variant matches a recent scheme by Shi et al. [191], which uses a similar approach.

7.4 Improving k -Clique Listing

Finally, one can use the **approximate degeneracy order** (ADG, cf. § 6.1) instead of DGR, which results in new performance bounds. Proceed as for the Edge Parallel variant, but use the $(2+\epsilon)$ -approximate parallel degeneracy order. This is easy to implement in the GMS benchmarking platform, as all one has to do is to change the preprocessing reordering routine from DGR to ADG. The depth becomes $\Theta(k(d + \frac{\epsilon}{2})^{k-2} + \log^2 n)$ and the work is increased to $\Omega(mk(d + \frac{\epsilon}{2})^{k-2})$. This design scales better if n is much bigger than kd^{k-2} and outperforms other variants in practice (see Section 8).

7.5 Improving Maximal Clique Listing

We now analyze our parallel maximal clique listing algorithm (cf. § 6.2). The key idea is to use ADG, the *relaxation* of the strict degeneracy order when processing vertices iteratively in the highest level of recursion in the BK algorithm. As in k -cliques, this is easy to implement with the GMS platform. This improves upon the Eppstein et al. [92] (BK-E) and Das et al. [79] (BK-DAS): our depth is better than both while work is better than that of BK-DAS and adds only a small factor to the work amount in BK-E.

For constant degeneracy graphs, our algorithm has linear work and poly-logarithmic depth, see Table 5 for a comparison with previous work. We note that, for many classes of sparse graphs, such as scale-free networks [16] and planar graphs [225], $\Delta \gg d$. Moreover, we often also have $\log n \ll \Delta$. Thus, the depth of BK-ADG is in such cases lower than that of BK-DAS.

We provide our bounds for the case where *nested parallelism* is employed. If only the outer loop which launches the calls to BK-Pivot and the construction of the arguments to BK-Pivot is parallelized, the depth is still $O((d\Delta)^{(2+\epsilon)d/3})$ and the space is only $O(m + nd + K)$ (where K is the output size).

	Work	Depth
Chiba and Nishizeki [69]	$O(d^2 n(n-d)3^{d/3})$	$O(d^2 n(n-d)3^{d/3})$.
Chiba and Nishizeki [69]	$O(nd^{d+1})$	$O(nd^{d+1})$.
Chrobak and Eppstein [72]	$O(nd^2 2^d)$	$O(nd^2 2^d)$.
Eppstein et al. [92]	$O(dm 3^{\frac{d}{3}})$	$O(dm 3^{\frac{d}{3}})$.
Das et al. [79]	$O(3^{\frac{n}{3}})$	$O(d \log n)$.
This Paper	$O(dm 3^{\frac{(2+\epsilon)d}{3}})$	$O(\log^2 n + d \log n)$.

Table 6: Additional bounds for enumerating all maximal cliques.

We first state the cost of computing the ADG order (cf. § 6.1), which is the key difference to the algorithm by Das et al. [79].

LEMMA 7.1. *Computing a $(2 + \epsilon)$ -approximate degeneracy order takes $O(m)$ work and $O(\log^2 n)$ depth, for any constant ϵ .*

Eppstein gave a generic work bound for an invocation of BK-Pivot(P, v_i, X) that we can use for our setting.

LEMMA 7.2 (EPPSTEIN [92]). *Excluding the work to report the found maximal cliques, BK-Pivot(P, v_i, X) takes $O((d|X|)3^{|P|/3})$ work.*

We combine Eppstein’s bound with the bounds on BK-Pivot and ADG to obtain work and depth bounds for BK-ADG.

LEMMA 7.3. *Finding all maximal cliques with BK-ADG takes $O(dm 3^{(2+\epsilon)d/3})$ work and $O(\log^2 n + d \log n)$ depth.*

PROOF. We first sketch the used parallel compute primitives. *Intersecting* two sets A and B takes $O(|A| |B|)$ work and $O(1)$ depth. Performing a *Reduction* over an array of n values (for example to compute their sum) takes $O(n)$ work and $O(\log n)$ depth.

Computing ADG takes $O(m)$ work and $O(\log^2 n)$ depth; see Lemma 7.1. Next, for all invocations of BK-Pivot, $|P| \leq (2 + \epsilon)d$, by the properties of the ADG order. Moreover, the size of the set $|X|$ in the invocation of BK-Pivot(P, v_i, X) is at most $\Delta(v_i)$. Hence, by using Lemma 7.2 for each invocation of BK-Pivot, we conclude that the work is $O(dm 3^{(2+\epsilon)d/3})$, excluding the cost to output the maximal cliques. Because there are $(n-d)3^{d/3}$ maximal cliques [92], the work is not dominated by the cost to report the maximal cliques.

The depth of BK-Pivot is $O(M \log n)$, where M is the size of the maximum clique [79]. As the size of the largest clique is bounded by the degeneracy (i.e., $M < d$), this is $O(d \log n)$. All the calls to BK-Pivot from BK-ADG can be launched simultaneously. \square

8 EVALUATION

We describe how GMS facilitates performance analysis of various aspects of graph mining, and accelerates the state of the art. *We focus on accelerating the four core mining problems from Section 6.*

8.1 Datasets, Methodology, Architectures

We first sketch the evaluation methodology. For measurements, we omit the first 1% of performance data as warmup. We derive enough data for the mean and 95% non-parametric confidence intervals. We use arithmetic means as summaries.

8.1.1 *Datasets.* We consider SNAP (S) [140], KONECT (K) [131], DIMACS (D) [82], Network Repository (N) [180], and WebGraph (W) [48] datasets. As explained in § 4.2, for flexibility, we do not fix specific datasets. Instead, we illustrate a wide selection of public datasets in Table 7, arguing *which parameters make them useful or challenging*. Details of these parameters are in § 4.2.

8.1.2 *Comparison Baselines.* For each considered graph mining problem, we compare different GMS variants to *the most optimized state-of-the-art algorithms available*. We compare to the original existing implementations. Details are stated in the following sections.

8.1.3 *Parallelism.* Unless stated otherwise, we use *full parallelism*, i.e., we run algorithms on *the maximum number of cores available on a given system*. We also analyzed *scaling* (how performance changes when varying number of used cores), the results show consistent advantages of GMS variants over other baselines.

8.1.4 *Architectures.* We used different systems for a *broad evaluation and to analyze and ensure performance portability* of our implementations. First, we use an in-house Einstein and Euler servers. Einstein is a Dell PowerEdge R910 with an Intel Xeon X7550 CPUs

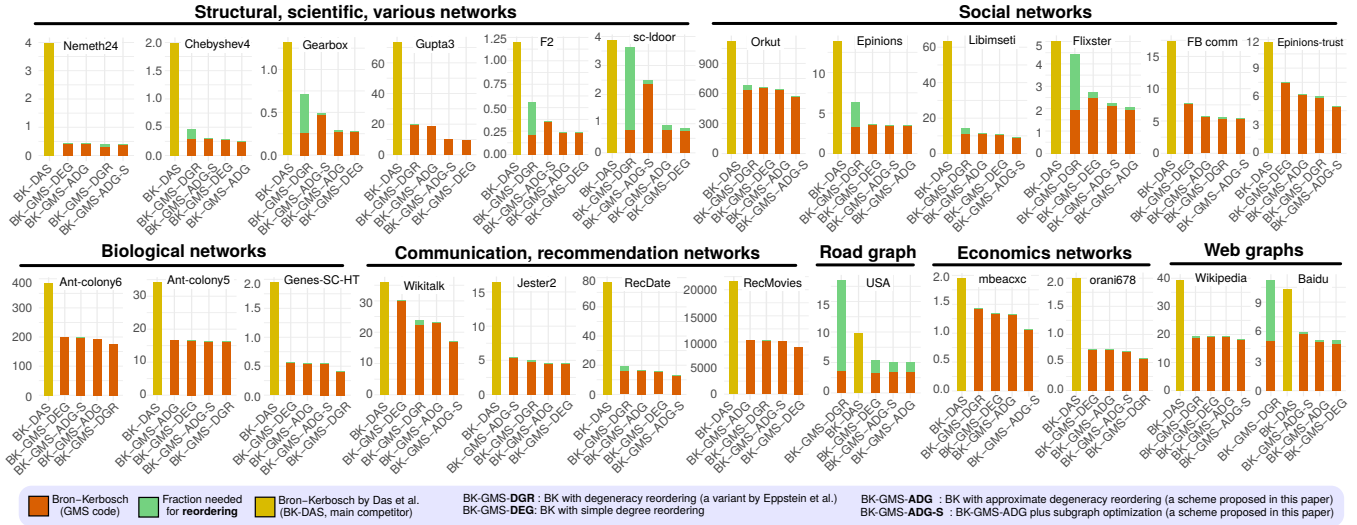


Figure 4: Speedups of the parallel GMS BK algorithm over a state-of-the-art implementation by Das et al. [79] (BK-DAS) and a recent algorithm by Eppstein et al. [91] (BK-GMS-DGR). System: Daint.

@ 2.00GHz with 18MB L3 cache, 1TiB RAM, and 32 cores per CPU (grouped in four sockets). Euler has an HT-enabled Intel Xeon Gold 6150 CPUs @ 2.70GHz with 24.75MB L3 cache, 64 GiB RAM, and 36 cores per CPU (grouped in two sockets). We also use servers from the CSCS supercomputing center, most importantly a compute server with Intel Xeon Gold 6140 CPU @ 2.30GHz, 768 GiB RAM, 18 cores, and 24.75MB L3. Finally, we also used XC50 compute nodes in the Piz Daint Cray supercomputer (one such node comes with 12-core Intel Xeon E5-2690 HT-enabled CPU 64 GiB RAM).

8.2 Faster Maximal Clique Listing

We start with our key result: GMS enabled us to outperform a state-of-the-art fastest available algorithm for maximal clique listing by Das et al. [79] (BK-DAS) by nearly an order of magnitude. The results are in Figure 4. We compare BK-DAS with several variants of BK developed in GMS as described in § 6. BK-GMS-DGR uses the degeneracy order and is a variant of the Eppstein’s scheme [91], enhanced in GMS. BK-GMS-DEG uses the simple degree ordering. BK-GMS-ADG and BK-GMS-ADG-S are two variants of a new BK algorithm proposed in this work, combining BK with the ADG ordering; the latter also uses the subgraph caching optimization (§ 6). We also compare to the original Eppstein scheme, it was always slower. GMS also enabled us to experiment with Intel Thread Building Blocks vs. OpenMP for threading in both the outermost loop and in inner loops (we exploit nested parallelism), we only show the OpenMP variants as they always outperform TBB.

Figure 4 shows consistent speedups of GMS variants over BK-DAS. We could quickly deliver these speedups by being able to plug in different set operations and optimizations in BK. Moreover, many plots show the large preprocessing overhead when using DGR. It sometimes helps to reduce the actual clique listing time (compared to ADG), but in most cases, “ADG plus clique listing” are faster than “DGR plus clique listing”: ADG is very fast and it

reduces the BK runtime to the level comparable to that achieved by DGR. This confirms the theoretical predictions of the benefits of BK-GMS-ADG over BK-GMS-DGR or BK-DAS. Finally, the comparably high performance (for many graphs) of BK-GMS-ADG, BK-GMS-ADG-S, and BK-GMS-DEG is due to the optimizations based on set algebra, for example using fast *and* compressed roaring bitmaps to implement neighborhoods and auxiliary sets P , X , and R (cf. § 6), which enables fast set operations heavily used in BK. Overall, BK-GMS is often faster than BK-DAS by $>50\%$, in some cases even $>9\times$.

We stress that the speedups of the implementations included in the GMS benchmarking platform are consistent over many graphs of different structural characteristics (cf. Table 7) that entail deeply varying load balancing properties. For example, some graphs are very sparse, with virtually no cliques larger than triangles (e.g., the USA road network) while others are relatively sparse with many triangles (and higher cliques), with low or moderate skews in triangle counts per vertex (e.g., Gearbox or F2). Finally, some graphs have large or even huge skews in triangle counts per vertex (e.g., Gupta3 or RecDate), which gives significant differences in the depths of the backtracking trees and thus load imbalance.

We also derived the **algorithmic efficiency** results, i.e., the number of maximal cliques found per second; selected data is in Figure 1. The results follow the run-times; the GMS schemes consistently outperform BK-DAS (the plots are in the technical report). These results show more distinctively that BK-GMS finds maximal cliques consistently better than BK-DAS, even if input graphs have vastly different clustering properties. For example, BK-GMS-ADG outperforms BK-DAS for Gupta3 (huge T -skew), F2 (medium T -skew), and ldoor (low T -skew).

8.3 Faster k -Clique Listing

GMS also enabled us to accelerate a very recent k -clique listing algorithm [78], see Figure 5. We were able to rapidly experiment with

Graph †	n	m	$\frac{m}{n}$	\widehat{d}_i	\widehat{d}_o	T	$\frac{T}{n}$	Why selected/special?
[so] (K) Orkut	3M	117M	38.1	33.3k	33.3k	628M	204.3	Common, relatively large
[so] (K) Flickr	2.3M	22.8M	9.9	21k	26.3k	838M	363.7	Large T but low m/n .
[so] (K) Libimseti	221k	17.2M	78	33.3k	25k	69M	312.8	Large m/n
[so] (K) Youtube	3.2M	9.3M	2.9	91.7k	91.7k	12.2M	3.8	Very low m/n and T
[so] (K) Flixster	2.5M	7.91M	3.1	1.4k	1.4k	7.89M	3.1	Very low m/n and T
[so] (K) Livemocha	104k	2.19M	21.1	2.98k	2.98k	3.36M	32.3	Similar to Flickr, but a lot fewer 4-cliques (4.36M)
[so] (N) Ep-trust	132k	841k	6	3.6k	3.6k	27.9M	212	Huge T -skew ($\widehat{T} = 108k$)
[so] (N) FB comm.	35.1k	1.5M	41.5	8.2k	8.2k	36.4M	1k	Large T -skew ($\widehat{T} = 159k$)
[wb] (K) DBpedia	12.1M	288M	23.7	963k	963k	11.68B	961.8	Rather low m/n but high T
[wb] (K) Wikipedia	18.2M	127M	6.9	632k	632k	328M	18.0	Common, very sparse
[wb] (K) Baidu	2.14M	17M	7.9	97.9k	2.5k	25.2M	11.8	Very sparse
[wb] (N) WikiEdit	94.3k	5.7M	60.4	107k	107k	835M	8.9k	Large T -skew ($\widehat{T} = 15.7M$)
[st] (N) Chebyshev4	68.1k	5.3M	77.8	68.1k	68.1k	445M	6.5k	Very large T and T/n and T -skew ($\widehat{T} = 5.8M$)
[st] (N) Gearbox	154k	4.5M	29.2	98	98	141M	915	Low \widehat{d} but large T ; low T -skew ($\widehat{T} = 1.7k$)
[st] (N) Nemeth25	10k	751k	75.1	192	192	87M	9k	Huge T but low $\widehat{T} = 12k$
[st] (N) F2	71.5k	2.6M	36.5	344	344	110M	1.5k	Medium T -skew ($\widehat{T} = 9.6k$)
[sc] (N) Gupta3	16.8k	4.7M	280	14.7k	14.7k	696M	41.5k	Huge T -skew ($\widehat{T} = 1.5M$)
[sc] (N) ldoor	952k	20.8M	21.5	76	76	567M	595	Very low T -skew ($\widehat{T} = 1.1k$)
[re] (N) MovieRec	70.2k	10M	142.4	35.3k	35.3k	983M	14k	Huge T and $\widehat{T} = 4.9M$
[re] (N) RecDate	169k	17.4M	102.5	33.4k	33.4k	286M	1.7k	Enormous T -skew ($\widehat{T} = 1.6M$)
[bi] (N) sc-ht (gene)	2.1k	63k	30	472	472	4.2M	2k	Large T -skew ($\widehat{T} = 27.7k$)
[bi] (N) AntColony6	164	10.3k	62.8	157	157	1.1M	6.6k	Very low T -skew ($\widehat{T} = 9.7k$)
[bi] (N) AntColony5	152	9.1k	59.8	150	150	897k	5.9k	Very low T -skew ($\widehat{T} = 8.8k$)
[co] (N) Jester2	50.7k	1.7M	33.5	50.8k	50.8k	127M	2.5k	Enormous T -skew ($\widehat{T} = 2.3M$)
[co] (K) Flickr (photo relations)	106k	2.31M	21.9	5.4k	5.4k	108M	1019	Similar to Livemocha, but many more 4-cliques (9.58B)
[ec] (N) mbeaccx	492	49.5k	100.5	679	679	9M	18.2k	Large T , low $\widehat{T} = 77.7k$
[ec] (N) orani678	2.5k	89.9k	35.5	1.7k	1.7k	8.7M	3.4k	Large T , low $\widehat{T} = 80.8k$
[ro] (D) USA roads	23.9M	28.8M	1.2	9	9	1.3M	0.1	Extremely low m/n and T

Table 7: Some considered real-world graphs. **Graph class/origin:** [so]: social network, [wb]: web graph, [st]: structural network, [sc]: scientific computing, [re]: recommendation network, [bi]: biological network, [co]: communication network, [ec]: economics network, [ro]: road graph. **Structural features:** m/n : graph sparsity, \widehat{d}_i : maximum in-degree, \widehat{d}_o : maximum out-degree, T : number of triangles, T/n : average triangle count per vertex, T -skew: a skew of triangle counts per vertex (i.e., the difference between the smallest and the largest number of triangles per vertex). Here, \widehat{T} is the maximum number of triangles per vertex in a given graph. **Dataset:** (W), (S), (K), (D), (C), and (N) refer to the publicly available datasets, explained in § 8.1. For more details, see § 4.2.

different variants, such as node parallel and edge parallel schemes, described in § 6.3 and in Section 7. Our optimizations from § 6.3 (e.g., a memory-efficient layout of C_i) ensure consistent speedups of up to 10% for different parameters (e.g., clique size k), input graphs, and reordering routines. Additionally, we show that using the ADG order brings further speedups over DEG or DGR.

8.4 Faster Degeneracy Reordering and k -Cores

We also analyze in more detail the performance of different reordering routines (DEG, DGR, and ADG) and their impact on graph mining algorithms in GMS (cf. § 6). We also show their impact on the run-time of BK maximal clique listing by Eppstein et al. [91] (BK-E). The results are in Figure 6. ADG, due to its beneficial scalability properties (cf. Section 7), *outperforms the exact DGR*. At the same time, it *similarly* reduces the runtime of BK-E [91] (cf. leftmost and rightmost bars). The $2 + \varepsilon$ approximation ratio has mild influence on performance. Specifically, the lower ε is, the more (mild) speedup is observed. This is because larger ε enables more parallelism, but then less accurate degeneracy ordering may incur more work when listing cliques. Moreover, ADG combined with BK-E cumulatively *outperforms the simple DEG reordering*: the latter

is also fast, but its impact on the Bron-Kerbosch run-time is lower, ultimately failing to provide comparable speedups. We note that the results of BK+DGR being slower than BK+DEG are consistent with independent results in a recent BK paper [79]. This additionally highlights our insight that using ADG over DGR and DEG is the fastest of all three variants. We were able to rapidly experiment with different reorderings as – *thanks to GMS’s modularity* – we could seamlessly integrate them with BK-E [91].

8.5 Faster Subgraph Isomorphism

GMS enabled us to accelerate a very recent parallel VF3-Light subgraph isomorphism baseline by 2.5 \times . The results are in Figure 7 (we use the same dataset as in the original work [60]). We illustrate the impact from different optimizations outlined in § 6. We were also able to use SIMD vectorization in the binary search part of the algorithms, leading to additional 1.1 \times speedup.

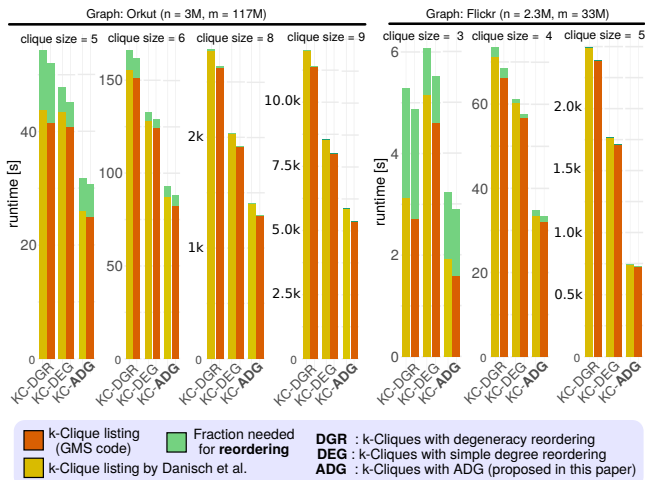


Figure 5: Speedups of (1) the GMS implementation of k -clique listing over a state-of-the-art algorithm [78], and (2) of the ADG reordering in k -clique listing over DGR/DEG. System: Daint.

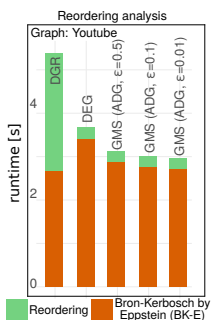


Figure 6: Speedups of ADG for different ϵ over DEG/DGR, details in § 8.4. System: Ault.

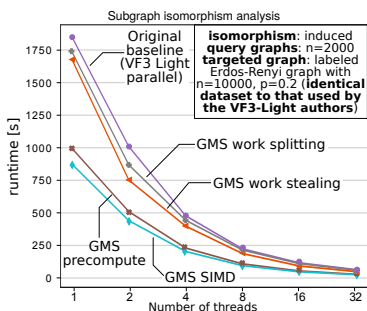


Figure 7: Speedups of different GMS variants of subgraph isomorphism over the state-of-the-art parallel VF3-Light baseline [60]. Details in § 8.5. System: Euler.

8.6 Subtleties of Higher-Order Structure

Subtleties of Higher-Order Structure One of the insights that we gained with GMS is that graphs similar in terms of n , m , sparsity m/n , and degree distributions, may have very different characteristics in their higher-order structure. For example, a graph of photo relations in Flickr and a Livemocha social network (see Table 7 for details) are similar in the above properties, but the former has 9,578,965,096 4-cliques while the latter has only 4,359,646 4-cliques. This is because, while a in a social network 4-cliques of friendships may be only *relatively common*, they should *occur very often* in a network where photos are related if they share *some* metadata (e.g., location). Thus, one should *carefully select input datasets* to properly evaluate respective graph mining algorithms, as seemingly similar graphs may have very different higher-order characteristics, which may vastly impact performance and conclusions when developing a new algorithm.

8.7 Analysis of Synthetic Graphs

We illustrate example results for synthetic graphs, see Figure 8a (with BK-GMS-DGR). Using power-law Kronecker graphs enable us to study the performance impact from varying the graph sparsity m/n while fixing all other parameters. For very sparse graphs, the cost of mining cliques is much lower than that of vertex reordering during preprocessing. However, as m/n increases, reordering begins to dominate. This is because Kronecker graphs in general do not have large cliques, which makes the mining process finish relatively fast, while reordering costs grow proportionally to m/n .

8.8 Machine Efficiency Analysis

We show example analysis of CPU utilization, using the PAPI interface in GMS, see Figure 8b. The plots illustrate the flattening of speedups with the increasing #threads, accompanied by the steady growth of stalled CPU cycles (both total counts and ratios), showing that maximal clique listing is memory bound [68, 89, 118, 221, 223].

8.9 Memory Consumption Analysis

We illustrate example memory consumption results in Figure 8c; we compare the size of three GMS set-centric graph representations, showing both *peak* memory usage when constructing a representation (bars) and sizes of ready representations (all in GB). Interestingly, while the latter are similar (except for v-usa), peak memory usage is visibly highest for RoaringSet. We also compare to the representation used by Das et al. [79], it always comes with the highest peak storage costs.

8.10 Algorithmic Throughput Analysis

The advantages of using algorithmic throughput can be seen by comparing Figure 1 and 4. While plain runtimes illustrate which algorithm is faster for which graph, the algorithmic throughput also *enables combining this outcome with the input graph structure*. For example, the GMS variants of BK have relatively lower benefits over BK by Das et al. [79] *whenever the input graph has a higher density of maximal cliques*. This motivates using the GMS BK especially for very sparse graphs without large dense clusters. One can derive analogous insights for any other patterns such a k -cliques.

8.11 GMS and Graph Processing Benchmarks

There is very little overlap with GMS and existing graph processing benchmarks, see Section 1 and Table 1. The closest one is GBBS [83], which supports the exact same variant of mining k -cliques. We compare GBBS to GMS in Figure 9; we also consider the edge-based very recent implementation by Danisch et al. [78]. GMS offers consistent advantages for different graphs and large clique sizes.

8.12 GMS and Pattern Matching Frameworks

There is also little overlap between GMS and pattern matching frameworks, cf. Table 1. While they support mining patterns, they focus on patterns of fixed sizes (e.g., k -cliques). We compare GMS to two very recent frameworks that, similarly to GMS, target shared-memory parallelism, Peregrine [118] and RStream [210]. Peregrine can only list k -cliques. It does not offer a native scheme for maximal clique listing and we implement it by iterating over k -cliques of

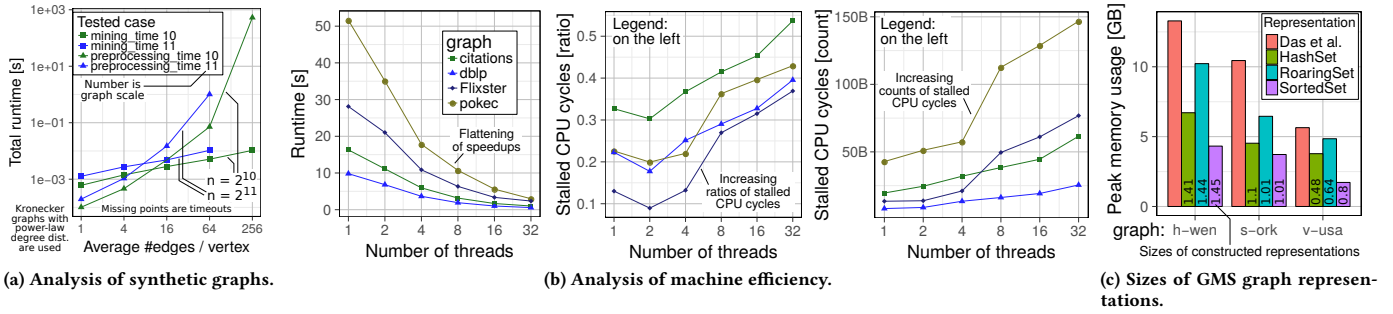


Figure 8: Additional analyses for the parallel GMS BK algorithm (BK-GMS-DGR). System: Daint.

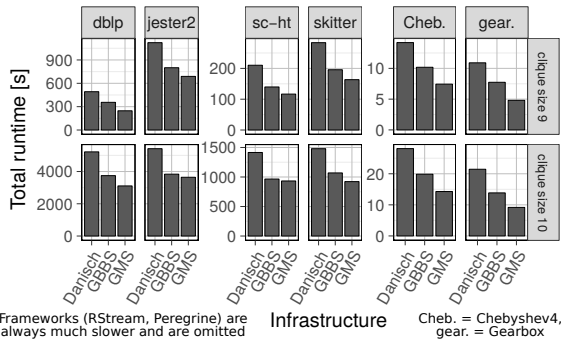


Figure 9: Comparison between GMS and the GBBS benchmark, for mining k -cliques. We also include Danisch’s algorithm (edge-centric) for additional reference. System: Einstein (full parallelism).

different sizes (we consult the authors of Peregrine to find the best scheme). RStream is only able to find k -cliques. Overall, GMS is much faster in all considered schemes (10-100× over Peregrine and more than 100× over RStream). This is because these systems focus on programming abstractions, which improves programmability but comes with performance overheads. GMS enables maximizing performance of parallel algorithms targeting specific problems.

9 RELATED WORK & DISCUSSION

We already exhaustively analyzed a large body of works related to graph processing benchmarks [26, 32, 147] and graph pattern matching frameworks, see Section 1 and Tables 1 and 2. General graph processing is summarized in a recent overview [182]. In general, GMS complements these works by delivering the first benchmarking suite that specifically targets graph mining.

While in the current GMS version we focus on the parallel shared-memory setting, GMS could be extended into multiple directions as future work. This includes moving into distributed processing [102, 103] and incorporating high-performance techniques [95, 104, 205] such as Remote Direct Memory Access [30, 33, 98, 99, 185, 187] combined with using general high-performance networks that work well with communication-intensive workloads [25, 31, 38, 85]. We are also working on variants of graph mining algorithms in GMS that harness the capabilities of the underlying hardware, such as low-diameter on-chip networks [29, 105, 160], NUMA and focus on data locality [187, 203], near- and in-memory processing [7, 8,

52, 112, 114, 134, 135, 163, 164, 170, 188–190], various architecture-related compression techniques [171, 172], and others [81, 128]. One could incorporate various forms of recently proposed lossy graph compression and summarization [34, 41, 144], and graph neural networks [22, 215, 216].

10 CONCLUSION

We introduce GraphMineSuite (GMS), the first benchmarking suite for graph mining algorithms. GMS offers an extensive *benchmark specification* and taxonomy that distill more than 300 related works and can aid in selecting appropriate comparison baselines. Moreover, GMS delivers a highly modular *benchmarking platform*, with dozens of parallel implementations of key graph mining algorithms and graph representations. Unlike *frameworks for pattern matching* which focus on abstractions and programming models for expressing mining specific patterns, GMS simplifies designing high-performance *algorithms for solving specific graph mining problems* from a *wide* graph mining area. Extending GMS towards distributed-memory systems or dynamic workloads are interesting future lines of work. Third, GMS’ *concurrency analysis* illustrates theoretical tradeoffs between time, work, storage, and accuracy, of several representative problems in graph mining; it can be used as a guide when rapidly analyzing the scalability of a planned graph mining scheme, or to obtain performance insights independent of implementation details. Finally, we show GMS’ potential by using it to *enhance state-of-the-art graph mining algorithms*, leading to theoretical and empirical advancements in maximal clique listing (speedups by >9× and better work-depth bounds over the *fastest known Bron-Kerbosch baseline*), degeneracy reordering and core decomposition (speedups by >2×), k -clique listing (speedups by up to 1.1× and better bounds), and subgraph isomorphism (speedups by 2.5×).

Acknowledgements: We thank Hussein Harake, Colin McMurtrie, Mark Klein, Angelo Mangili, and the whole CSCS team granting access to the Ault and Daint machines, and for their excellent technical support. We thank Timo Schneider for immense help with computing infrastructure at SPCL. We thank Maximilien Danisch, Oana Balalau, Mauro Sozio, Apurba Das, Seyed-Vahid Sanei-Mehri, and Srikanta Tirhappura for providing us with the implementations of their algorithms for solving k -clique and maximal clique listing. We thank Dimitrios Lekkas, Athina Sotiropoulou, Foteini Strati, Andreas Triantafyllos, Kenza Amara, Chia-I Hu, Ajaykumar Unagar, Roger Baumgartner, Severin Kistler, Emanuel Peter, and Alain Senn for helping with the implementation in the early stages of the project.

APPENDIX

We now provide extensions of several sections.

A DETAILS OF PROBLEMS AND ALGORITHMS IN GRAPH MINING

We additionally provide more details of considered graph mining problems and the associated algorithms.

- **Maximal Cliques Listing** For finding maximal cliques, we use the established Bron-Kerbosch (BK) algorithm [56], a recursive backtracking algorithm often used in practice, with well-known pivoting and degeneracy optimizations [61, 91, 151, 207].

- **k -Clique Listing** GMS considers listing k -cliques. We select a state-of-the-art algorithm by Danisch et al. [78]. The algorithm is somewhat similar to Bron-Kerbosch in that it is also recursive backtracking. The difference is that its work is polynomial. We also separately consider Triangle Counting as it comes with a plethora of specific studies [9, 184, 193].

- **Dense Non-Clique Subgraph Discovery** We also incorporate a problem of discovering dense *non-clique* subgraphs. Relevant classes of subgraphs are quasi-cliques, k -cores, k -plexes, kd -cliques, k -clubs, k -clans, $dalks$, $damks$, dks , k -clique-stars, and others [117, 136]. Here, we implemented a very recent algorithm for listing k -clique-stars [117]; k -clique-stars are dense subgraphs that combine the characteristics of cliques *and stars* (relaxing the restrictive nature of k -cliques) [117]. GMS also implements an exact and an approximate algorithm for k -core decomposition [94].

- **Subgraph Isomorphism** Subgraph isomorphism (SI) is an important NP-Complete problem, where one finds all embeddings of a certain *query graph* H in another *target graph* G . SI can be *non-induced* and *induced*; we consider both. Consider a case where an embedding of H is found in G , but there are some additional edges in G that connect some vertices that belong to this embedding. In the non-induced variant, this situation is permitted, unlike in the induced variant, where the found embedding cannot have such additional edges. In GMS, we consider recent algorithms: VF2 [75] and adapted Glasgow [155] for induced SI, and VF3-Light [59] and TurboISO [108] for non-induced SI. Our selection covers different approaches for solving SI: VF2 and VF3-Light represent backtracking as they descent from the well-known ULLMAN algorithm [209]. TurboISO uses a graph indexing as opposed to backtracking while Glasgow incorporates implied constraints.

- **Frequent Subgraph Mining** We separately consider the Frequent Subgraph Mining (FSM) problem [120], in which one finds all subgraphs that occur more often than a specified threshold. An FSM algorithm consists of (1) a strategy for exploring the tree of candidate subgraphs, and (2) a subroutine where one checks if a candidate is included in the processed graph. (2) usually solves the subgraph isomorphism problem, covered above. (1) uses either a BFS-based or a DFS-based exploration strategy.

- **Vertex Similarity** Vertex similarity measures can be used on their own, for example in graph database queries [179], or as a building block of more complex algorithms such as clustering [119]. We consider seven measures: Jaccard, Overlap, Adamic Adar, Resource Allocation, Common Neighbors, Total Neighbors, and Preferential Attachment measures [137, 179]. All these measures associate (in

different ways) the degree of similarity between v and u with the number of common neighbors of vertices v and u .

- **Link Prediction** Here, one is interested in developing schemes for predicting whether two non-adjacent vertices can become connected in the future. There exist many schemes for such prediction that are based on variations of vertex similarity [10, 142, 146, 202]. We provide them in GMS, as well as a simple algorithm for assessing the accuracy of a specific link prediction scheme [211], which assesses how well a given prediction scheme works.

- **Clustering and Community Detection** We consider graph clustering and community detection, a widely studied problem. We pick Jarvis-Patrick clustering (JP) [119], a scheme that uses similarity of two vertices to determine whether these two vertices are in the same cluster. Moreover, we consider Label Propagation [175] and the Louvain method [46], two established methods for detecting communities that, respectively, use the notions of *label dominance* and *modularity* in assigning vertices to communities.

- **Approximate Degeneracy Ordering** We also consider an easily parallelizable algorithm to compute an *approximate* degeneracy order (the algorithm has $O(\log n)$ iterations for any constant $\epsilon > 0$ and has an approximation ratio of $2 + \epsilon$ [24]). The algorithm is based on a streaming scheme for large graphs [94] and uses set cardinality and difference. The derived degeneracy order can be directly used to compute the k -core of G (a maximal connected subgraph of G whose all vertices have degree at least k). This is done by iterating over vertices in the degeneracy order and removing all vertices with out-degree less than k (in the oriented graph).

- **Optimization Problems** Third, we also consider some problems from a family of *optimization problems*, also deemed important in the literature [5]. Here, we focus on graph coloring (GC), considering several graph coloring algorithms that represent different approaches: Jones and Plassmann’s [123] and Hasenplaugh et al.’s [110] heuristics based on appropriate vertex orderings and vertex prioritization, Johansson’s [121] and Barenboim’s [17] randomized palette-based heuristics that use conflict resolution, and Elkin et al.’s [90] and sparse-dense decomposition [109] that are examples of state-of-the-art distributed algorithms.

B NAVIGATING THE MAZE OF GRAPH REPRESENTATIONS

The right data layout is one of key enablers of high performance. For this, we now overview the most relevant graph representations and compression schemes. We picture key designs in Figure 10.

B.1 Graph Representations

We consider several graph representations.

B.1.1 Adjacency List and Adjacency Array. AL is a very popular representation that uses $O(m \log n + n \log m)$ space, with many implementations and variants. AM uses $O(n^2)$ space and is thus rarely directly used. However, several interesting compression schemes are based on AM, for example k^2 -trees or some succinct graphs [34].

B.1.2 CSR aka Adjacency Array. Compressed Sparse Row (CSR), also referred to as Adjacency Array (AA), usually consists of n arrays that contain neighborhoods of graph vertices. Each array is usually sorted by vertex IDs. AA also contains a structure with

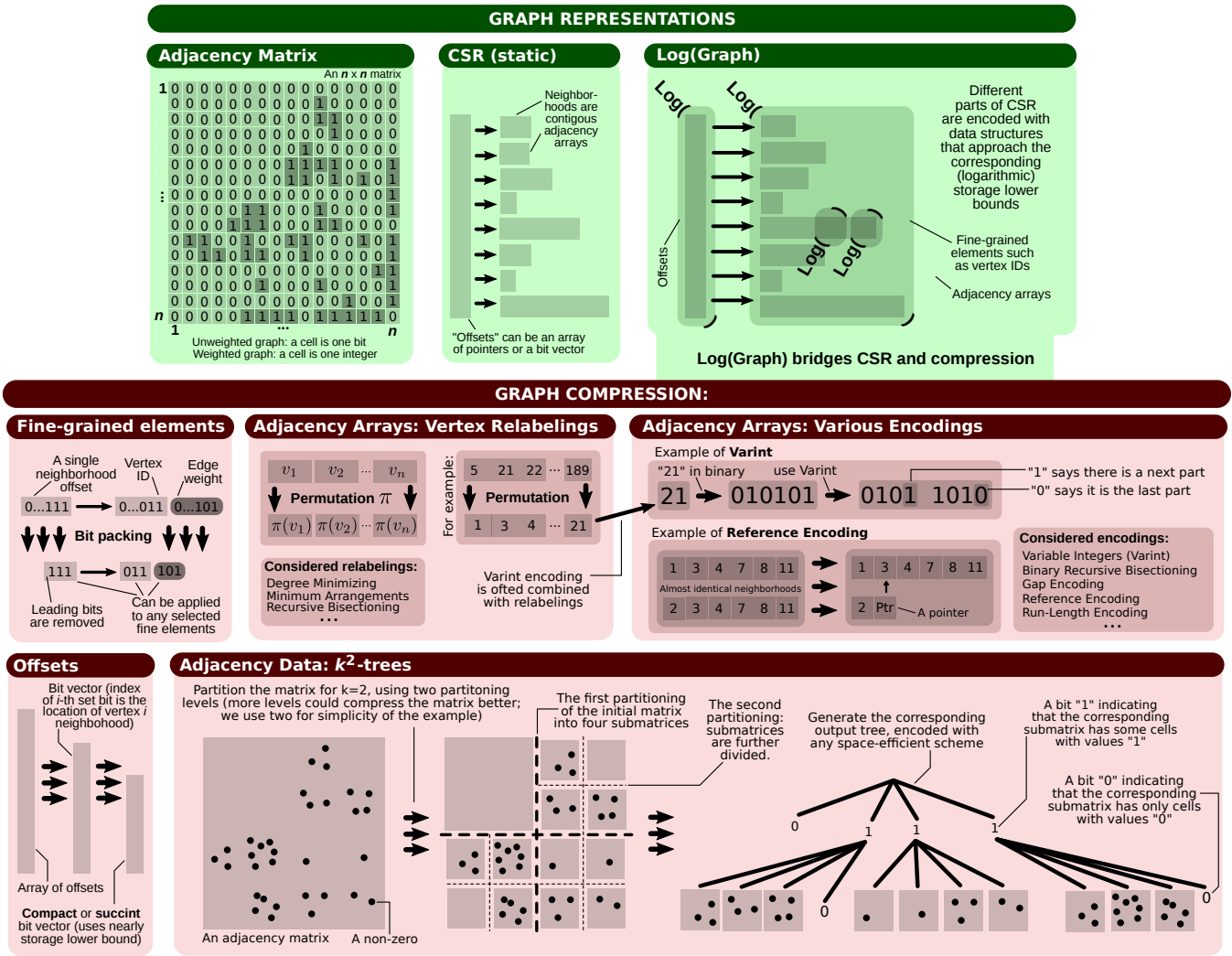


Figure 10: An overview of models, representations, and compression schemes. The GMS logo indicates the ones provided in the current GMS platform. A recent survey provides full details on all the representations [34]

offsets (or pointers) to each neighborhood array. AA is very popular in processing engines for static graphs [40]. Due to its simplicity, it offers very low latency of accesses. Moreover, its variants are used in graph streaming settings [28], where edges and vertices may be inserted or deleted over time.

B.1.3 Log(Graph). Log(Graph) [40] is a recently proposed variant of AA, in which one separately compresses fine elements of the representation (vertex IDs, edge weights, etc.) as well as coarse parts, such as a whole offset array. The main compression method in Log(Graph) is encoding each considered graph element using a data structure that approaches the corresponding logarithmic storage lower bounds while simultaneously enabling fast accesses. An example such structure used in Log(Graph) are succinct bit vectors. One Log(Graph) advantage is low-overhead decompression. Another benefit is a tunable storage-performance tradeoff: one can choose to

compress more aggressively at the cost of more costly decompression, and vice versa. Third, Log(Graph) is modular: the user can select which parts of AA are compressed.

B.2 Graph Compression Schemes

In a survey on lossless graph compression and space-efficient graph representations [34], we illustrate that the majority of graph compression methods fall into two major families of methods: relabelings (permutations) and transformations.

One type of the considered compression schemes for adjacency data are **relabelings** that permute vertex IDs. Different permutations enable more or less efficient compression of vertex IDs (e.g., when combining permutations with a Varint compression and gap encoding [40]). Established examples use shingle ordering [70], recursive bisection [43], degree minimizing [40], and Layered Label Propagation [47].

Algorithm	AL (sorted)	AM	EL (unsorted)	EL (sorted)
Node Iterator (TC)	$O\left(n + m^{3/2} \log \Delta\right)^*$	$O\left(n + m^{3/2}\right)$	$O\left(n + m^{3/2} (\Delta + \log m)\right)$	$O\left(n + m^{5/2}\right)$
Rank Merge (TC)	$O\left(n + n\Delta + m^{3/2}\right)$	$O\left(n + n\Delta + m^{3/2}\right)$	$O\left(n + n\Delta + m^{3/2}\right)$	$O\left(n + n\Delta + m^{3/2}\right)$
BFS, top-down	$\Theta(n + m)$	$\Theta(n + m)$	$O(n \log m + m)$	$O(nm + n + m)$
PageRank, pushing	$O\left(n + m^{3/2} \log \Delta\right)^*$	$O\left(n + m^{3/2}\right)$	$O\left(n + m^{3/2} (\Delta + \log m)\right)$	$O\left(n + m^{5/2}\right)$
D -Stepping (SSSP)	$O\left(n + m + \frac{L}{D} + n_D + m_D\right)$	$O\left(n^2 + \frac{L}{D} + nn_D + m_D\right)$	$O\left(nm + \frac{L}{D} + n_D (\log m + \Delta) + m_D\right)$	$O\left(nm + m + \frac{L}{D} + n_D m + m_D\right)$
Bellman-Ford (SSSP)	$O(n^2 + nm)$	$O(n^3)$	$O(n + nm)$	$O(n + nm)$
Boruvka (MST)	$O(m \log n)$	$O(n^2 \log n)$	$O(nm \log n \log m)$	$O(n^2 m)$
Boman (Graph Coloring)	$O(n + m)$	$O(n^2)$	$O(n^2)$	$O(n + nm)$
Betweenness Centrality	$O(nm)$	$O(n^3)$	$O(nm \log m)$	$O(nm^2)$

Table 8: Time complexity of graph algorithms for different graph representations. “*” indicates that the $\log \Delta$ terms becomes Δ when the used AL representation is unsorted. D is a parameter of the Delta-Stepping algorithm that controls the “bucket size” and thus the amount of parallelism (for $D = 1$ one obtains Dijkstra’s algorithm while for $D = \infty$ one obtains the Bellman-Ford algorithm). L is the maximum length of a shortest path between any two vertices.

Graph query	AL	AM	EL (unsorted)	EL (sorted)
Iterate over all vertices	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Iterate over all edges	$\Theta(n + m)$	$\Theta(n^2)$	$\Theta(m)$	$\Theta(m)$
Iterate over a neighborhood	$\Theta(\Delta)$	$\Theta(n)$	$\Theta(m)$	$\Theta(\log m + \Delta)^{\#}$
Check vertex’ degree	$\Theta(n)^*$	$\Theta(n)^*$	$\Theta(m)^*$	$\Theta(\log m + \Delta)^{**}$
Check edge’s existence	$O(\log \Delta)$	$O(1)$	$O(m)$	$O(\log m)$
Check edge’s weight	$O(1)$	$O(n)$	$\Theta(m)$	$\Theta(\log m + \Delta)^{\#}$

Table 9: Time complexity of graph queries for different graph representations. “*” indicates that a given complexity can be reduced to $O(1)$ with $\Theta(m + n)$ preprocessing and n auxiliary storage. “#” indicates that a given complexity assumes that each edge (u, v) is present twice in the edge list (both in u ’s and in v ’s neighborhoods), which requires double storage but does not increase the preprocessing complexity.

Another type of compression schemes provided by GMS are **transformations** that apply a certain function to the adjacency data [40] after a relabeling is used. Here, in addition to the above-mentioned Varint and gap encoding, GMS considers k^2 -trees [55], run-length and reference encoding [40], and implementation of certain schemes proposed in WebGraph [49].

B.3 Theoretical Analysis

We offer a brief theoretical analysis on the impact of different representations on the performance of graph queries and algorithms. The analysis for the former can be found in Table 9. For the latter, Table 8 provides time complexities of Triangle Counting (Node-Iterator and Rank-Merge schemes [193]), Page Rank (pushing and pulling [26]), BFS, Betweenness Centrality (Proutzos et al.’ algorithm [173]), Single Source Shortest Paths (Δ -Stepping [156] and Bellman-Ford [21]), Minimum Graph Coloring (Boman et al.’ algorithm [50]), and Minimum Spanning Tree (Boruvka’s algorithm [53]).

C ADDITIONAL RESULTS

Figure 11 shows additional results for the performance of various Bron-Kerbosch variants, when measured in the mined cliques per time unit.

REFERENCES

- [1] C. R. Aberger, A. Lamb, S. Tu, A. Nötzli, K. Olukotun, and C. Ré. Emptyheaded: A relational engine for graph processing. *ACM Transactions on Database Systems (TODS)*, 42(4):1–44, 2017.

- [2] M. Adedoyin-Olowe, M. M. Gaber, and F. Stahl. A survey of data mining techniques for social media analysis. *arXiv preprint arXiv:1312.4617*, 2013.
- [3] C. C. Aggarwal and H. Wang. Graph data management and mining: A survey of algorithms and applications. In *Managing and mining graph data*, pages 13–68. Springer, 2010.
- [4] C. C. Aggarwal and H. Wang. A survey of clustering algorithms for graph data. In *Managing and mining graph data*, pages 275–301. Springer, 2010.
- [5] C. C. Aggarwal, H. Wang, et al. *Managing and mining graph data*, volume 40. Springer, 2010.
- [6] M. Ahmad, F. Hijaz, Q. Shi, and O. Khan. Crono: A benchmark suite for multithreaded graph algorithms executing on futuristic multicores. In *2015 IEEE International Symposium on Workload Characterization*, pages 44–55. IEEE, 2015.
- [7] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi. A scalable processing-in-memory accelerator for parallel graph processing. In *Proceedings of the 42nd Annual International Symposium on Computer Architecture*, pages 105–117, 2015.
- [8] J. Ahn, S. Yoo, O. Mutlu, and K. Choi. Pim-enabled instructions: A low-overhead, locality-aware processing-in-memory architecture. In *2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*, pages 336–348. IEEE, 2015.
- [9] M. Al Hasan and V. S. Dave. Triangle counting in large networks: a review. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 8(2):e1226, 2018.
- [10] M. Al Hasan et al. Link prediction using supervised learning. In *SDM*, 2006.
- [11] M. Al Hasan and M. J. Zaki. A survey of link prediction in social networks. In *Social network data analytics*, pages 243–275. Springer, 2011.
- [12] K. Ammar and M. T. Özsu. Wgb: Towards a universal graph benchmark. In *Advancing Big Data Benchmarks*, pages 58–72. Springer, 2013.
- [13] T. G. Armstrong, V. Ponnkanti, D. Borthakur, and M. Callaghan. Linkbench: a database benchmark based on the facebook social graph. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 1185–1196, 2013.
- [14] K. Asanovic et al. A view of the parallel computing landscape. *CACM*, 2009.
- [15] D. A. Bader and K. Madduri. Design and implementation of the hpcc graph analysis benchmark on symmetric multiprocessors. In *International Conference on High-Performance Computing*, pages 465–476. Springer, 2005.
- [16] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- [17] L. Barenboim et al. The locality of distributed symmetry breaking. *JACM*, 2016.
- [18] S. Bassini, M. Danelutto, and P. Dazzi. *Parallel Computing is Everywhere*, volume 32. IOS Press, 2018.
- [19] S. Beamer, K. Asanović, and D. Patterson. Direction-optimizing breadth-first search. *Scientific Programming*, 21(3-4):137–148, 2013.
- [20] S. Beamer, K. Asanović, and D. Patterson. The gap benchmark suite. *arXiv preprint arXiv:1508.03619*, 2015.
- [21] R. Bellman. On a routing problem. *Quarterly of applied mathematics*, 16(1):87–90, 1958.
- [22] T. Ben-Nun, M. Besta, S. Huber, A. N. Ziogas, D. Peter, and T. Hoefler. A modular benchmarking infrastructure for high-performance and reproducible deep learning. In *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 66–77. IEEE, 2019.
- [23] P. Berkhin. A survey of clustering data mining techniques. In *Grouping multidimensional data*, pages 25–71. Springer, 2006.
- [24] M. Besta, A. Carigiet, Z. Vonarburg-Shmaria, K. Janda, L. Gianinazzi, and T. Hoefler. High-performance parallel graph coloring with strong guarantees on work, depth, and quality. In *ACM/IEEE Supercomputing*, 2020.



Figure 11: Additional data from performance analysis related to mining maximal cliques; the Y axis plots the number of maximal cliques mined per time unit. The naming of the schemes is the same as in the main body of the document, with the following exceptions: "BK-GMS-SG" indicates the subgraph optimization ("BK-GMS-ADG-S"), "BK-GMS-DG" indicates the degeneracy ordering ("BK-GMS-DGR"), and "BK-GMS-DE" as well as "BK-TBB-DE" are the degree orderings in GMS and TBB, respectively.

- [25] M. Besta, J. Domke, M. Schneider, M. Konieczny, S. Di Girolamo, T. Schneider, A. Singla, and T. Hoefler. High-performance routing with multipathing and path diversity in ethernet and hpc networks. *IEEE Transactions on Parallel and Distributed Systems*, 32(4):943–959, 2020.
- [26] M. Besta et al. To push or to pull: On reducing communication and synchronization in graph computations. In *ACM HPDC*, 2017.
- [27] M. Besta, M. Fischer, T. Ben-Nun, D. Stanojevic, J. D. F. Licht, and T. Hoefler. Substream-centric maximum matchings on fpga. *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, 13(2):1–33, 2020.
- [28] M. Besta, M. Fischer, V. Kalavri, M. Kapralov, and T. Hoefler. Practice of streaming processing of dynamic graphs: Concepts, models, and systems. *arXiv preprint arXiv:1912.12740*, 2019.
- [29] M. Besta, S. M. Hassan, S. Salamanchili, R. Ausavarungnirun, O. Mutlu, and T. Hoefler. Slim noc: A low-diameter on-chip network topology for high energy efficiency and scalability. In *ACM SIGPLAN Notices*, 2018.
- [30] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. In *ACM HPDC*, pages 37–48, 2014.
- [31] M. Besta and T. Hoefler. Slim Fly: A Cost Effective Low-Diameter Network Topology. Nov. 2014. ACM/IEEE Supercomputing.
- [32] M. Besta and T. Hoefler. Accelerating irregular computations with hardware transactional memory and active messages. In *ACM HPDC*, 2015.
- [33] M. Besta and T. Hoefler. Active access: A mechanism for high-performance distributed data-centric computations. In *ACM ICS*, 2015.
- [34] M. Besta and T. Hoefler. Survey and taxonomy of lossless graph compression and space-efficient graph representations. *arXiv preprint arXiv:1806.01799*, 2018.
- [35] M. Besta, R. Kanakagiri, H. Mustafa, M. Karasikov, G. Rättsch, T. Hoefler, and E. Solomonik. Communication-efficient jaccard similarity for high-performance distributed genome comparisons. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 1122–1132. IEEE, 2020.
- [36] M. Besta, F. Marending, E. Solomonik, and T. Hoefler. Slimsell: A vectorizable graph representation for breadth-first search. In *Parallel and Distributed Processing Symposium (IPDPS), 2017 IEEE International*, pages 32–41. IEEE, 2017.
- [37] M. Besta, E. Peter, R. Gerstenberger, M. Fischer, M. Podstawski, C. Barthels, G. Alonso, and T. Hoefler. Demystifying graph databases: Analysis and taxonomy of data organization, system designs, and graph queries. *arXiv preprint arXiv:1910.09017*, 2019.
- [38] M. Besta, M. Schneider, K. Cynk, M. Konieczny, E. Henriksson, S. Di Girolamo, A. Singla, and T. Hoefler. Fatpaths: Routing in supercomputers and data centers when shortest paths fall short. *ACM/IEEE Supercomputing*, 2020.
- [39] M. Besta, D. Stanojevic, J. D. F. Licht, T. Ben-Nun, and T. Hoefler. Graph processing on fpgas: Taxonomy, survey, challenges. *arXiv preprint arXiv:1903.06697*, 2019.
- [40] M. Besta, D. Stanojevic, T. Zivic, J. Singh, M. Hoerold, and T. Hoefler. Log (graph): a near-optimal high-performance graph representation. In *Proceedings of the 27th International Conference on Parallel Architectures and Compilation Techniques*, page 7. ACM, 2018.
- [41] M. Besta, S. Weber, L. Gianinazzi, R. Gerstenberger, A. Ivanov, Y. Oltchik, and T. Hoefler. Slim graph: practical lossy graph compression for approximate graph processing, storage, and analytics. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, page 35. ACM, 2019.
- [42] G. Bilardi and A. Pietracaprina. *Models of Computation, Theoretical*, pages 1150–1158. Springer US, Boston, MA, 2011.
- [43] D. K. Blandford, G. E. Blelloch, and L. A. Kash. Compact representations of separable graphs. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '03*, pages 679–688, Philadelphia, PA, USA, 2003. Society for Industrial and Applied Mathematics.
- [44] G. E. Blelloch. Problem based benchmark suite, 2011.
- [45] G. E. Blelloch and B. M. Maggs. *Parallel Algorithms*, page 25. Chapman & Hall/CRC, 2 edition, 2010.
- [46] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
- [47] P. Boldi, M. Rosa, M. Santini, and S. Vigna. Layered label propagation: A multiresolution coordinate-free ordering for compressing social networks. In *Proceedings of the 20th international conference on World wide web*, pages 587–596. ACM, 2011.
- [48] P. Boldi and S. Vigna. The WebGraph framework I: Compression techniques. In *World Wide Web Conf. (WWW)*, pages 595–601, 2004.
- [49] P. Boldi and S. Vigna. The webgraph framework i: compression techniques. In *Proceedings of the 13th international conference on World Wide Web*, pages 595–602. ACM, 2004.
- [50] E. G. Boman et al. A scalable parallel graph coloring algorithm for distributed memory computers. In *Euro-Par*, pages 241–251. 2005.
- [51] P. Boncz. LDBC: benchmarks for graph and RDF data management. In *IDEAS*, 2013.
- [52] A. Boroumand, S. Ghose, M. Patel, H. Hassan, B. Lucia, K. Hsieh, K. T. Malladi, H. Zheng, and O. Mutlu. Lazyrim: An efficient cache coherence mechanism for processing-in-memory. *IEEE Computer Architecture Letters*, 16(1):46–50, 2016.
- [53] O. Boruvka. O jistém problému minimálním. 1926.
- [54] U. Brandes. A faster algorithm for betweenness centrality. *J. of Math. Sociology*, 25(2):163–177, 2001.
- [55] N. R. Brisaboa, S. Ladra, and G. Navarro. k 2-trees for compact web graph representation. In *International Symposium on String Processing and Information Retrieval*, pages 18–30. Springer, 2009.
- [56] C. Bron and J. Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *CACM*, 1973.
- [57] M. Burtscher, R. Nasre, and K. Pingali. A quantitative study of irregular programs on gpus. In *2012 IEEE International Symposium on Workload Characterization (IISWC)*, pages 141–151. IEEE, 2012.
- [58] V. Carletti et al. Introducing vf3: A new algorithm for subgraph isomorphism. In *Springer GBRPR*, 2017.
- [59] V. Carletti et al. The VF3-light subgraph isomorphism algorithm: when doing less is more effective. In *Springer S+SSPR*, 2018.
- [60] V. Carletti et al. A parallel algorithm for subgraph isomorphism. In *Springer GBRPR*, 2019.
- [61] F. Cazals and C. Karande. A note on the problem of reporting maximal cliques. *Theoretical Computer Science*, 407(1-3):564–568, 2008.
- [62] P. Celis. *Robin hood hashing*. University of Waterloo, 1986.
- [63] D. Chakrabarti and C. Faloutsos. Graph mining: Laws, generators, and algorithms. *ACM CSUR*, 2006.
- [64] S. Chambi, D. Lemire, O. Kaser, and R. Godin. Better bitmap performance with roaring bitmaps. *Software: practice and experience*, 46(5):709–719, 2016.
- [65] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, and K. Skadron. Rodinia: A benchmark suite for heterogeneous computing. In *2009 IEEE international symposium on workload characterization (IISWC)*, pages 44–54. Ieee, 2009.
- [66] H. Chen, M. Liu, Y. Zhao, X. Yan, D. Yan, and J. Cheng. G-miner: an efficient task-oriented graph mining system. In *Proceedings of the Thirteenth EuroSys Conference*, page 32. ACM, 2018.
- [67] X. Chen, R. Dathathri, G. Gill, and K. Pingali. Pangolin: An efficient and flexible graph mining system on cpu and gpu. *arXiv preprint arXiv:1911.06969*, 2019.
- [68] J. Cheng, L. Zhu, Y. Ke, and S. Chu. Fast algorithms for maximal clique enumeration with limited memory. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1240–1248, 2012.
- [69] N. Chiba and T. Nishizeki. Arboricity and subgraph listing algorithms. *SIAM Journal on computing*, 14(1):210–223, 1985.
- [70] F. Chierichetti, R. Kumar, S. Lattanzi, M. Mitzenmacher, A. Panconesi, and P. Raghavan. On compressing social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 219–228. ACM, 2009.
- [71] A. Ching et al. One trillion edges: Graph processing at facebook-scale. *VLDB*, 2015.
- [72] M. Chrobak and D. Eppstein. Planar orientations with low out-degree and compaction of adjacency matrices. *Theoretical Computer Science*, 86(2):243–266, 1991.
- [73] D. J. Cook and L. B. Holder. *Mining graph data*. John Wiley & Sons, 2006.
- [74] M. Copik, G. Kwasniewski, M. Besta, M. Podstawski, and T. Hoefler. Sebs: A serverless benchmark suite for function-as-a-service computing. *arXiv preprint arXiv:2012.14132*, 2020.
- [75] L. P. Cordella et al. A (sub) graph isomorphism algorithm for matching large graphs. *IEEE TPAMI*, 2004.
- [76] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT press, 2009.
- [77] G. Csardi and T. Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695(5):1–9, 2006.
- [78] M. Danisch et al. Listing k-cliques in sparse real-world graphs. In *WWW*, 2018.
- [79] A. Das et al. Shared-memory parallel maximal clique enumeration. In *IEEE HiPC*, 2018.
- [80] A. Das, M. Svendsen, and S. Tirthapura. Change-sensitive algorithms for maintaining maximal cliques in a dynamic graph. *CoRR*, vol. abs/1601.06311, 2016.
- [81] J. de Fine Licht, S. Meierhans, and T. Hoefler. Transformations of high-level synthesis codes for high-performance computing. *arXiv preprint arXiv:1805.08288*, 2018.
- [82] C. Demetrescu, A. V. Goldberg, and D. S. Johnson. *The Shortest Path Problem: Ninth DIMACS Implementation Challenge*, volume 74. American Math. Soc., 2009.
- [83] L. Dhulipala et al. Theoretically efficient parallel graph algorithms can be fast and scalable. In *ACM SPAA*, 2018.
- [84] L. Dhulipala, J. Shi, T. Tseng, G. E. Blelloch, and J. Shun. The graph based benchmark suite (gbbs). In *Proceedings of the 3rd Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)*, pages 1–8, 2020.
- [85] S. Di Girolamo, K. Taranov, A. Kurth, M. Schaffner, T. Schneider, J. Beránek, M. Besta, L. Benini, D. Roweth, and T. Hoefler. Network-accelerated non-contiguous memory transfers. *arXiv preprint arXiv:1908.08590*, 2019.

- [86] V. Dias et al. Fractal: A general-purpose graph pattern mining system. In *ACM SIGMOD*, 2019.
- [87] R. Diestel. *Graph theory*. Springer, 2018.
- [88] N. Du, B. Wu, L. Xu, B. Wang, and X. Pei. A parallel algorithm for enumerating all maximal cliques in complex network. In *Sixth IEEE International Conference on Data Mining-Workshops (ICDMW'06)*, pages 320–324. IEEE, 2006.
- [89] J. D. Eblen, C. A. Phillips, G. L. Rogers, and M. A. Langston. The maximum clique enumeration problem: algorithms, applications, and implementations. In *BMC bioinformatics*, volume 13, page S5. Springer, 2012.
- [90] M. Elkin et al. $(2\delta - 1)$ -edge-coloring is much easier than maximal matching in the distributed setting. In *ACM-SIAM SODA*, 2014.
- [91] D. Eppstein et al. Listing all maximal cliques in sparse graphs in near-optimal time. In *SAAC*, 2010.
- [92] D. Eppstein and D. Strash. Listing all maximal cliques in large sparse real-world graphs. In *International Symposium on Experimental Algorithms*, pages 364–375. Springer, 2011.
- [93] P. Erdős and A. Rényi. On the evolution of random graphs. *Selected Papers of Alfred Rényi*, 1976.
- [94] M. Farach-Colton and M. Tsai. Computing the degeneracy of large graphs. In *LATIN*, 2014.
- [95] J. S. Firoz, M. Zalewski, A. Lumsdaine, and M. Barnas. Runtime scheduling policies for distributed graph algorithms. In *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 640–649. IEEE, 2018.
- [96] H. Fu, J. Liao, J. Yang, L. Wang, Z. Song, X. Huang, C. Yang, W. Xue, F. Liu, F. Qiao, et al. The sunway taishan supercomputer: system and applications. *Science China Information Sciences*, 59(7):072001, 2016.
- [97] B. Gallagher. Matching structure and semantics: A survey on graph-based pattern matching. In *AAAI Fall Symposium: Capturing and Using Patterns for Evidence Detection*, pages 45–53, 2006.
- [98] R. Gerstenberger, M. Besta, and T. Hoefler. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. In *ACM/IEEE Supercomputing, SC '13*, pages 53:1–53:12, 2013.
- [99] R. Gerstenberger, M. Besta, and T. Hoefler. Enabling highly scalable remote memory access programming with mpi-3 one sided. *Communications of the ACM*, 61(10):106–113, 2018.
- [100] L. Gianinazzi et al. Communication-avoiding parallel minimum cuts and connected components. In *ACM PPOPP*, pages 219–232. ACM, 2018.
- [101] S. Gog, T. Beller, A. Moffat, and M. Petri. From theory to practice: Plug and play with succinct data structures. In *International Symposium on Experimental Algorithms*, pages 326–337. Springer, 2014.
- [102] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: distributed graph-parallel computation on natural graphs. In *OSDI*, volume 12, page 2, 2012.
- [103] J. E. Gonzalez, R. S. Xin, A. Dave, D. Crankshaw, M. J. Franklin, and I. Stoica. Graphx: Graph processing in a distributed dataflow framework. In *OSDI*, volume 14, pages 599–613, 2014.
- [104] D. Gregor and A. Lumsdaine. The parallel bgl: A generic library for distributed graph computations. *Parallel Object-Oriented Scientific Computing (POOSC)*, 2:1–18, 2005.
- [105] B. Grot, J. Hestness, S. W. Keckler, and O. Mutlu. Kilo-noc: A heterogeneous network-on-chip architecture for scalability and service guarantees. In *2011 38th Annual International Symposium on Computer Architecture (ISCA)*, pages 401–412. IEEE, 2011.
- [106] M. Han, K. Daudjee, K. Ammar, M. T. Özsu, X. Wang, and T. Jin. An experimental comparison of pregel-like graph processing systems. *Proceedings of the VLDB Endowment*, 7(12):1047–1058, 2014.
- [107] S. Han, L. Zou, and J. X. Yu. Speeding up set intersections in graph algorithms using simd instructions. In *Proceedings of the 2018 International Conference on Management of Data*, pages 1587–1602. ACM, 2018.
- [108] W.-S. Han et al. Turbo iso: towards ultrafast and robust subgraph isomorphism search in large graph databases. In *ACM SIGMOD/PODS*. ACM, 2013.
- [109] D. G. Harris, J. Schneider, and H.-H. Su. Distributed $(\delta + 1)$ -coloring in sublogarithmic rounds. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 465–478, 2016.
- [110] W. Hasenplaugh et al. Ordering heuristics for parallel graph coloring. In *Proceedings of the 26th ACM symposium on Parallelism in algorithms and architectures*, pages 166–177, 2014.
- [111] W. Hasenplaugh, T. Kaler, T. B. Schardl, and C. E. Leiserson. Ordering heuristics for parallel graph coloring. In *26th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '14, Prague, Czech Republic - June 23 - 25, 2014*, pages 166–177, 2014.
- [112] H. Hassan, G. Pekhimenko, N. Vijaykumar, V. Seshadri, D. Lee, O. Ergin, and O. Mutlu. Chargecache: Reducing dram latency by exploiting row access locality. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 581–593. IEEE, 2016.
- [113] T. Horváth et al. Cyclic pattern kernels for predictive graph mining. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2004.
- [114] K. Hsieh, S. Khan, N. Vijaykumar, K. K. Chang, A. Boroumand, S. Ghose, and O. Mutlu. Accelerating pointer chasing in 3d-stacked memory: Challenges, mechanisms, evaluation. In *2016 IEEE 34th International Conference on Computer Design (ICCD)*, pages 25–32. IEEE, 2016.
- [115] M. Injadat, F. Salo, and A. B. Nassif. Data mining techniques in social media: A survey. *Neurocomputing*, 214:654–670, 2016.
- [116] A. P. Iyer, Z. Liu, X. Jin, S. Venkataraman, V. Braverman, and I. Stoica. {ASAP}: Fast, approximate graph pattern mining at scale. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, pages 745–761, 2018.
- [117] S. Jabbour, N. Mhadhbi, B. Raddaoui, and L. Sais. Pushing the envelope in overlapping communities detection. In *International Symposium on Intelligent Data Analysis*, pages 151–163. Springer, 2018.
- [118] K. Jamshidi, R. Mahadasa, and K. Vora. Peregrine: a pattern-aware graph mining system. In *Proceedings of the Fifteenth European Conference on Computer Systems*, pages 1–16, 2020.
- [119] R. A. Jarvis and E. A. Patrick. Clustering using a similarity measure based on shared near neighbors. *IEEE Transactions on computers*, 100(11):1025–1034, 1973.
- [120] C. Jiang, F. Coenen, and M. Zito. A survey of frequent subgraph mining algorithms. *The Knowledge Engineering Review*, 28(1):75–105, 2013.
- [121] Ö. Johansson. Simple distributed $\delta + 1$ -coloring of graphs. *Information Processing Letters*, 70(5), 1999.
- [122] D. S. Johnson, M. Yannakakis, and C. H. Papadimitriou. On generating all maximal independent sets. *Information Processing Letters*, 27(3):119–123, 1988.
- [123] M. T. Jones and P. E. Plassmann. A parallel graph coloring heuristic. *SIAM Journal on Scientific Computing*, 14(3):654–669, 1993.
- [124] A. Joshi, Y. Zhang, P. Bogdanov, and J.-H. Hwang. An efficient system for subgraph discovery. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 703–712. IEEE, 2018.
- [125] D. R. Karger and C. Stein. A new approach to the minimum cut problem. *Journal of the ACM (JACM)*, 43(4):601–640, 1996.
- [126] J. Kepner, P. Aaltonen, D. Bader, A. Buluç, F. Franchetti, J. Gilbert, D. Hutchison, M. Kumar, A. Lumsdaine, H. Meyerhenke, et al. Mathematical foundations of the graphblas. In *High Performance Extreme Computing Conference (HPEC)*, 2016 IEEE, pages 1–9. IEEE, 2016.
- [127] W. Khaouid et al. K-core decomposition of large networks on a single pc. *Proceedings of the VLDB Endowment*, 9(1):13–23, 2015.
- [128] Y. Kim, V. Seshadri, D. Lee, J. Liu, and O. Mutlu. A case for exploiting subarray-level parallelism (salp) in dram. In *2012 39th Annual International Symposium on Computer Architecture (ISCA)*, pages 368–379. IEEE, 2012.
- [129] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. *Theoretical Computer Science*, 250(1-2):1–30, 2001.
- [130] F. Kose, W. Weckwerth, T. Linke, and O. Fiehn. Visualizing plant metabolomic correlation networks using clique-metabolite matrices. *Bioinformatics*, 17(12):1198–1208, 2001.
- [131] J. Kunegis. Konect: the koblenz network collection. In *Proc. of Intl. Conf. on World Wide Web (WWW)*, pages 1343–1350. ACM, 2013.
- [132] G. Kwasniewski, T. Ben-Nun, A. N. Ziogas, T. Schneider, M. Besta, and T. Hoefler. On the parallel i/o optimality of linear algebra kernels: near-optimal lu factorization. In *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 463–464, 2021.
- [133] G. Kwasniewski, M. Kabić, M. Besta, J. VandeVondele, R. Solcà, and T. Hoefler. Red-blue pebbling revisited: near optimal parallel matrix-matrix multiplication. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–22, 2019.
- [134] B. C. Lee, P. Zhou, J. Yang, Y. Zhang, B. Zhao, E. Ipek, O. Mutlu, and D. Burger. Phase-change technology and the future of main memory. *IEEE micro*, 30(1):143–143, 2010.
- [135] D. Lee, Y. Kim, V. Seshadri, J. Liu, L. Subramanian, and O. Mutlu. Tiered-latency dram: A low latency and low cost dram architecture. In *2013 IEEE 19th International Symposium on High Performance Computer Architecture (HPCA)*, pages 615–626. IEEE, 2013.
- [136] V. E. Lee, N. Ruan, R. Jin, and C. Aggarwal. A survey of algorithms for dense subgraph discovery. In *Managing and Mining Graph Data*, pages 303–336. Springer, 2010.
- [137] E. A. Leicht et al. Vertex similarity in networks. *Physical Review E*, 73(2):026120, 2006.
- [138] D. Lemire, O. Kaser, N. Kurz, L. Deri, C. O'Hara, F. Saint-Jacques, and G. Ssi-Yan-Kai. Roaring bitmaps: Implementation of an optimized software library. *Software: Practice and Experience*, 48(4):867–895, 2018.
- [139] J. Leskovec et al. Kronecker graphs: An approach to modeling networks. *J. of Machine Learning Research*, 11(Feb):985–1042, 2010.
- [140] J. Leskovec and A. Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [141] B. Lessley, T. Perciano, M. Mathai, H. Childs, and E. W. Bethel. Maximal clique enumeration with data-parallel primitives. In *2017 IEEE 7th Symposium on Large Data Analysis and Visualization (LDAV)*, pages 16–25. IEEE, 2017.

- [142] D. Liben-Nowell and J. Kleinberg. The link-prediction problem for social networks. *Journal of the American society for information science and technology*, 58(7):1019–1031, 2007.
- [143] H. Lin et al. Shentu: processing multi-trillion edge graphs on millions of cores in seconds. In *ACM/IEEE Supercomputing*. IEEE Press, 2018.
- [144] Y. Liu, T. Safavi, A. Dighe, and D. Koutra. Graph summarization methods and applications: A survey. *ACM Computing Surveys (CSUR)*, 51(3):1–34, 2018.
- [145] L. Lu, Y. Gu, and R. Grossman. dmaximalcliques: A distributed algorithm for enumerating all maximal cliques and maximal clique distribution. In *2010 IEEE International Conference on Data Mining Workshops*, pages 1320–1327. IEEE, 2010.
- [146] L. Lü and T. Zhou. Link prediction in complex networks: A survey. *Physica A: statistical mechanics and its applications*, 390(6):1150–1170, 2011.
- [147] A. Lumsdaine, D. Gregor, B. Hendrickson, and J. W. Berry. Challenges in Parallel Graph Processing. *Par. Proc. Let.*, 17(1):5–20, 2007.
- [148] K. Madduri, D. Ediger, K. Jiang, D. A. Bader, and D. Chavarria-Miranda. A faster parallel algorithm and efficient multithreaded implementations for evaluating betweenness centrality on massive datasets. In *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–8. IEEE, 2009.
- [149] K. Makino and T. Uno. New algorithms for enumerating all maximal cliques. In *Scandinavian workshop on algorithm theory*, pages 260–272. Springer, 2004.
- [150] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 135–146. ACM, 2010.
- [151] G. Manoussakis. An output sensitive algorithm for maximal clique enumeration in sparse graphs. In *12th International Symposium on Parameterized and Exact Computation (IPEC 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [152] D. W. Matula and L. L. Beck. Smallest-last ordering and clustering and graph coloring algorithms. *JACM*, 1983.
- [153] D. Mawhirter, S. Reinehr, C. Holmes, T. Liu, and B. Wu. Graphzero: Breaking symmetry for efficient graph mining. *arXiv preprint arXiv:1911.12877*, 2019.
- [154] D. Mawhirter and B. Wu. Automine: harmonizing high-level abstraction and high performance for graph mining. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pages 509–523. ACM, 2019.
- [155] C. McCreesh and P. Prosser. A parallel, backjumping subgraph isomorphism algorithm using supplemental graphs. In *CP*. Springer, 2015.
- [156] U. Meyer and P. Sanders. Δ -stepping: a parallelizable shortest path algorithm. *Journal of Algorithms*, 49(1):114–152, 2003.
- [157] G. L. Miller et al. Improved parallel algorithms for spanners and hopsets. In *ACM SPAA*. ACM, 2015.
- [158] Z. Ming, C. Luo, W. Gao, R. Han, Q. Yang, L. Wang, and J. Zhan. Bdgs: A scalable big data generator suite in big data benchmarking. In *Advancing Big Data Benchmarks*, pages 138–154. Springer, 2013.
- [159] J. W. Moon and L. Moser. On cliques in graphs. *Israel journal of Mathematics*, 3(1):23–28, 1965.
- [160] T. Moscibroda and O. Mutlu. A case for bufferless routing in on-chip networks. In *Proceedings of the 36th annual international symposium on Computer architecture*, pages 196–207, 2009.
- [161] P. J. Mucci, S. Browne, C. Deane, and G. Ho. Papi: A portable interface to hardware performance counters. In *Proceedings of the department of defense HPCMP users group conference*, volume 710, 1999.
- [162] R. C. Murphy et al. Introducing the graph 500. *Cray User’s Group (CUG)*, 2010.
- [163] O. Mutlu. Memory scaling: A systems architecture perspective. In *2013 5th IEEE International Memory Workshop*, pages 21–25. IEEE, 2013.
- [164] O. Mutlu and L. Subramanian. Research problems and opportunities in memory systems. *Supercomputing frontiers and innovations*, 1(3):19–55, 2015.
- [165] L. Nai, Y. Xia, I. G. Tanase, H. Kim, and C.-Y. Lin. Graphbig: understanding graph computing in the context of industrial solutions. In *SC’15: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–12. IEEE, 2015.
- [166] T. J. Ottosen and J. Vomlel. Honour thy neighbour—clique maintenance in dynamic graphs. *on Probabilistic Graphical Models*, page 201, 2010.
- [167] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [168] C. H. Papadimitriou and K. Steiglitz. *Combinatorial optimization: algorithms and complexity*. Courier Corporation, 1998.
- [169] S. Parthasarathy, S. Tatikonda, and D. Ucar. A survey of graph mining techniques for biological datasets. In *Managing and mining graph data*, pages 547–580. Springer, 2010.
- [170] A. Pattnaik, X. Tang, A. Jog, O. Kayiran, A. K. Mishra, M. T. Kandemir, O. Mutlu, and C. R. Das. Scheduling techniques for gpu architectures with processing-in-memory capabilities. In *Proceedings of the 2016 International Conference on Parallel Architectures and Compilation*, pages 31–44, 2016.
- [171] G. Pekhimenko, V. Seshadri, Y. Kim, H. Xin, O. Mutlu, P. B. Gibbons, M. A. Kozuch, and T. C. Mowry. Linearly compressed pages: A low-complexity, low-latency main memory compression framework. In *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 172–184, 2013.
- [172] G. Pekhimenko, V. Seshadri, O. Mutlu, M. A. Kozuch, P. B. Gibbons, and T. C. Mowry. Base-delta-immediate compression: Practical data compression for on-chip caches. In *2012 21st international conference on parallel architectures and compilation techniques (PACT)*, pages 377–388. IEEE, 2012.
- [173] D. Proutzios and K. Pingali. Betweenness centrality: algorithms and implementations. In *ACM SIGPLAN Notices*, volume 48, pages 35–46. ACM, 2013.
- [174] A. Quamar, A. Deshpande, and J. Lin. Nscale: neighborhood-centric large-scale graph analytics in the cloud. *The VLDB Journal—The International Journal on Very Large Data Bases*, 25(2):125–150, 2016.
- [175] U. N. Raghavan, R. Albert, and S. Kumara. Near linear time algorithm to detect community structures in large-scale networks. *Physical review E*, 76(3):036106, 2007.
- [176] T. Ramraj and R. Prabhakar. Frequent subgraph mining algorithms—a survey. *Procedia Computer Science*, 47:197–204, 2015.
- [177] S. U. Rehman, A. U. Khan, and S. Fong. Graph mining: A survey of graph mining techniques. In *Seventh International Conference on Digital Information Management (ICDIM 2012)*, pages 88–92. IEEE, 2012.
- [178] P. Ribeiro, P. Paredes, M. E. Silva, D. Aparicio, and F. Silva. A survey on subgraph counting: Concepts, algorithms and applications to network motifs and graphlets. *arXiv preprint arXiv:1910.13011*, 2019.
- [179] I. Robinson, J. Webber, and E. Eifrem. *Graph databases*. " O’Reilly Media, Inc.", 2013.
- [180] R. A. Rossi and N. K. Ahmed. The network data repository with interactive graph analytics and visualization. In *AAAI*, 2015.
- [181] A. Roy, I. Mihailovic, and W. Zwaenepoel. X-stream: Edge-centric graph processing using streaming partitions. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 472–488. ACM, 2013.
- [182] S. Saker, A. Bonifati, H. Voigt, A. Iosup, K. Ammar, R. Angles, W. Aref, M. Arenas, M. Besta, P. A. Boncz, et al. The future is big graphs! a community view on graph processing systems. *arXiv preprint arXiv:2012.06171*, 2020.
- [183] S. E. Schaeffer. Graph clustering. *Computer science review*, 1(1):27–64, 2007.
- [184] T. Schank. Algorithmic aspects of triangle-based network analysis. *Phd in computer science, University Karlsruhe*, 3, 2007.
- [185] P. Schmid, M. Besta, and T. Hoefler. High-performance distributed RMA locks. In *ACM HPDC*, pages 19–30, 2016.
- [186] M. C. Schmidt, N. F. Samatova, K. Thomas, and B.-H. Park. A scalable, parallel algorithm for maximal clique enumeration. *Journal of Parallel and Distributed Computing*, 69(4):417–428, 2009.
- [187] H. Schweizer, M. Besta, and T. Hoefler. Evaluating the cost of atomic operations on modern architectures. In *IEEE PACT*, pages 445–456, 2015.
- [188] V. Seshadri, K. Hsieh, A. Boroum, D. Lee, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. Fast bulk bitwise and and or in dram. *IEEE Computer Architecture Letters*, 14(2):127–131, 2015.
- [189] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, M. A. Kozuch, et al. Rowclone: Fast and energy-efficient in-dram bulk data copy and initialization. In *Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 185–197, 2013.
- [190] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry. Ambit: In-memory accelerator for bulk bitwise operations using commodity dram technology. In *2017 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 273–287. IEEE, 2017.
- [191] J. Shi, L. Dhulipala, and J. Shun. Parallel clique counting and peeling algorithms. *arXiv preprint arXiv:2002.10047*, 2020.
- [192] J. Shun and G. E. Blelloch. Ligma: a lightweight graph processing framework for shared memory. In *ACM SIGPLAN Notices*, volume 48, pages 135–146, 2013.
- [193] J. Shun and K. Tangwongsan. Multicore triangle computations without tuning. In *Data Engineering (ICDE), 2015 IEEE 31st International Conference on*, pages 149–160. IEEE, 2015.
- [194] E. Solomonik, M. Besta, F. Vella, and T. Hoefler. Scaling betweenness centrality using communication-efficient sparse matrix multiplication. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–14, 2017.
- [195] C. L. Staudt and H. Meyerhenke. Engineering parallel algorithms for community detection in massive networks. *IEEE Transactions on Parallel and Distributed Systems*, 27(1):171–184, 2015.
- [196] V. Stix. Finding all maximal cliques in dynamic graphs. *Computational Optimization and applications*, 27(2):173–186, 2004.
- [197] J. A. Stratton, C. Rodrigues, I.-J. Sung, N. Obeid, L.-W. Chang, N. Anssari, G. D. Liu, and W.-m. Hwu. Parboil: A revised benchmark suite for scientific and commercial throughput computing. *Center for Reliable and High-Performance Computing*, 127, 2012.
- [198] T. Suzumura et al. Performance characteristics of Graph500 on large-scale distributed environment. In *Workload Char. (IISWC), IEEE Intl. Symp. on*, pages 149–158, 2011.
- [199] M. Svendsen, A. P. Mukherjee, and S. Tirthapura. Mining maximal cliques from a large graph using mapreduce: Tackling highly uneven subproblem sizes.

- Journal of Parallel and distributed computing*, 79:104–114, 2015.
- [200] L. Tang and H. Liu. Graph mining applications to social network analysis. In *Managing and Mining Graph Data*, pages 487–513. Springer, 2010.
- [201] Y. Tang. Benchmarking graph databases with cyclone benchmark. 2016.
- [202] B. Taskar et al. Link prediction in relational data. In *NIPS*, pages 659–666, 2004.
- [203] A. Tate, A. Kamil, A. Dubey, A. Größlinger, B. Chamberlain, B. Goglin, C. Edwards, C. J. Newburn, D. Padua, D. Unat, et al. Programming abstractions for data locality. 2014.
- [204] C. H. Teixeira et al. Arabesque: a system for distributed graph mining. In *Proceedings of the 25th Symposium on Operating Systems Principles*, pages 425–440. ACM, 2015.
- [205] L. Thebault. *Scalable and efficient algorithms for unstructured mesh computations*. PhD thesis, 2016.
- [206] L. T. Thomas, S. R. Valluri, and K. Karlapalem. Margin: Maximal frequent subgraph mining. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 4(3):1–42, 2010.
- [207] E. Tomita, A. Tanaka, and H. Takahashi. The worst-case time complexity for generating all maximal cliques and computational experiments. *Theor. Comput. Sci.*, 363(1):28–42, 2006.
- [208] S. Tsukiyama, M. Ide, H. Ariyoshi, and I. Shirakawa. A new algorithm for generating all the maximal independent sets. *SIAM Journal on Computing*, 6(3):505–517, 1977.
- [209] J. R. Ullmann. An algorithm for subgraph isomorphism. *Journal of the ACM (JACM)*, 23(1):31–42, 1976.
- [210] K. Wang, Z. Zuo, J. Thorpe, T. Q. Nguyen, and G. H. Xu. Rstream: marrying relational algebra with streaming for efficient graph mining on a single machine. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, pages 763–782, 2018.
- [211] L. Wang, K. Hu, and Y. Tang. Robustness of link-prediction algorithm based on similarity and application to biological networks. *Current Bioinformatics*, 9(3):246–252, 2014.
- [212] L. Wang, J. Zhan, C. Luo, Y. Zhu, Q. Yang, Y. He, W. Gao, Z. Jia, Y. Shi, S. Zhang, et al. Bigdatabench: A big data benchmark suite from internet services. In *2014 IEEE 20th international symposium on high performance computer architecture (HPCA)*, pages 488–499. IEEE, 2014.
- [213] T. Washio and H. Motoda. State of the art of graph-based data mining. *Acm Sigkdd Explorations Newsletter*, 5(1):59–68, 2003.
- [214] B. Wu, S. Yang, H. Zhao, and B. Wang. A distributed algorithm to enumerate all maximal cliques in mapreduce. In *2009 Fourth International Conference on Frontier of Computer Science and Technology*, pages 45–51. IEEE, 2009.
- [215] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [216] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- [217] Y. Xu, J. Cheng, and A. W.-C. Fu. Distributed maximal clique computation and management. *IEEE Transactions on Services Computing*, 9(1):110–122, 2015.
- [218] Z. Xu, X. Chen, J. Shen, Y. Zhang, C. Chen, and C. Yang. Gardenia: A graph processing benchmark suite for next-generation accelerators. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 15(1):1–13, 2019.
- [219] D. Yan, H. Chen, J. Cheng, M. T. Özsu, Q. Zhang, and J. Lui. G-thinker: big graph mining made easier and faster. *arXiv preprint arXiv:1709.03110*, 2017.
- [220] D. Yan, W. Qu, G. Guo, and X. Wang. Prefixpm: A parallel framework for general-purpose frequent pattern mining. In *Proceedings of the 36th IEEE International Conference on Data Engineering (ICDE) 2020*, 2020.
- [221] P. Yao, L. Zheng, Z. Zeng, Y. Huang, C. Gui, X. Liao, H. Jin, and J. Xue. A locality-aware energy-efficient accelerator for graph mining applications. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 895–907. IEEE, 2020.
- [222] K. Zhang, R. Chen, and H. Chen. NUMA-aware graph-structured analytics. In *Proceedings of the 20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 183–193. ACM, 2015.
- [223] Y. Zhang, F. N. Abu-Khzam, N. E. Baldwin, E. J. Chesler, M. A. Langston, and N. F. Samatova. Genome-scale computational approaches to memory-intensive applications in systems biology. In *SC'05: Proceedings of the 2005 ACM/IEEE Conference on Supercomputing*, pages 12–12. IEEE, 2005.
- [224] C. Zhao, Z. Zhang, P. Xu, T. Zheng, and X. Cheng. Kaleido: An efficient out-of-core graph mining system on a single machine. *arXiv preprint arXiv:1905.09572*, 2019.
- [225] X. Zhou and T. Nishizeki. Edge-coloring and f-coloring for various classes of graphs. In *Algorithms and Computation, 5th International Symposium, ISAAC '94, Beijing, P. R. China, August 25-27, 1994, Proceedings*, pages 199–207, 1994.