



A Single Machine System for Querying Big Graphs with PRAM

Yang Liu
Beihang University, China
ly_act@buaa.edu.cn

Wenfei Fan
Shenzhen Institute of Computing
Sciences, China
Beihang University, China
University of Edinburgh, UK
wenfei@inf.ed.ac.uk

Shuhao Liu*
Shenzhen Institute of Computing
Sciences, China
shuhao@sics.ac.cn

Xiaoke Zhu
Beihang University, China
zhuxk@buaa.edu.cn

Jianxin Li
Beihang University, China
lijx@buaa.edu.cn

ABSTRACT

This paper develops Planar (Plug and play PRAM), a single-machine system for graph analytics by reusing existing PRAM algorithms, without the need for designing new parallel algorithms. Planar supports both out-of-core and in-memory analytics. When a graph is too big to fit into the memory of a machine, Planar adapts PRAM to limited resources by extending a fixpoint model with multi-core parallelism, using disk as memory extension. For an in-memory task, it dedicates all available CPU cores to the task, and allows parallelly scalable PRAM algorithms to retain the property, *i.e.*, the more cores are available, the less runtime is taken. We develop a graph partitioning and work scheduling strategy to accommodate subgraph I/O, balance memory usage and reduce runtime, beyond traditional partitioners for multi-machine systems. Using real-life graphs, we empirically verify that Planar outperforms SOTA in-memory and out-of-core systems in efficiency and scalability.

PVLDB Reference Format:

Yang Liu, Wenfei Fan, Shuhao Liu, Xiaoke Zhu, and Jianxin Li. A Single Machine System for Querying Big Graphs with PRAM. PVLDB, 18(3): 756-769, 2024.
doi:10.14778/3712221.3712240

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/SICS-Fundamental-Research-Center/Planar>.

1 INTRODUCTION

A host of single-machine systems have been developed for graph analytics via multi-core parallelism, *e.g.*, [6, 31, 49, 56, 57, 63, 68, 71, 77, 80, 85, 86, 91]. These systems typically adopt a *vertex-centric* (VC) or *edge-centric* (EC) parallel model. A VC (resp. EC) program pivots computation around each vertex (resp. edge); it may only directly access its local data and adjacent edges (resp. endpoints), but it has to exchange information with “remote” entities via message passing. VC/EC is often inefficient in programming and execution

[6, 26, 35, 74]. It is nontrivial to program for problems that are constrained by “joint” conditions on multiple vertices, *e.g.*, subgraph isomorphism [28]. Moreover, the local scope of VC/EC operations often incurs redundant computation [74], and its message-passing model introduces extra synchronization complexity [5]. These overheads are often excessive, leading to limited scalability and high COST [60] of a graph system under a shared-memory architecture.

On the other hand, parallel models have been studied for shared-memory architectures for decades, notably PRAM (Parallel Random Access Machine) [27, 32, 76]. PRAM allows multiple processors to work in parallel via *single-instruction-multiple-data* (SIMD), and synchronize via shared memory. A large number of PRAM algorithms are already in place, and many of them are provably work-time optimal [38] or parallelly scalable, *i.e.*, they guarantee that the more processors are used, the less parallel runtime is taken [48].

Is it possible to develop a single-machine graph system in which one can plug existing PRAM algorithms, and the system executes the algorithms to make the most of multi-core parallelism and the shared memory of the machine? This way, the users do not have to think like a vertex/edge and develop new parallel algorithms starting from scratch; instead, they can simply leverage the decades of work on PRAM and make effective use of the well-developed PRAM algorithms, capitalizing on their scalability and efficiency.

It is, however, nontrivial to run PRAM algorithms in a single-machine system. The PRAM model assumes that the memory is large enough to load the entire dataset at once, and there are a polynomial number of processors [7, 34]. In contrast, a single-machine system has a fixed number of CPU cores, limited memory capacity and disk I/O bandwidth. One has to simulate a polynomial number of cores assumed by PRAM. Worse yet, for *out-of-core* processing of graphs that are too large to fit into the main memory of a machine at once, it needs to use disk as memory extension [6, 31, 49, 57, 68, 80, 91]. Such systems have to carefully partition graphs and schedule I/O and CPU operations so that their CPUs do not have to wait for long for subgraphs to be loaded into the memory.

Contributions & Organization. This paper develops Planar (Plug and play PRAM), a single-machine system for running PRAM algorithms for graph analytics. Underlying Planar are the following.

(1) *Parallel model* (Section 3). Planar proposes a unified parallel model for both in-memory and out-of-core tasks. For a query class Q , it takes as input an existing PRAM algorithm \mathcal{A} and a graph G . When G fits into the memory of a single machine, it executes

*Corresponding author.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 18, No. 3 ISSN 2150-8097.
doi:10.14778/3712221.3712240

algorithm \mathcal{A} on G with all available cores. When the graph is too big, it partitions G into subgraphs such that each subgraph can fit into the memory, and uses the secondary storage as memory extension. It loads the subgraphs into memory one by one, and runs SIMD for multi-core parallelism on each in-memory subgraph with all available cores. It iterates the computation over all subgraphs until it reaches a fixpoint, adapting the graph-centric model (GC) [26].

The parallel model makes the first effort to adapt PRAM to physical machines in the real world. It simplifies parallel graph programming by reusing existing PRAM algorithms, and retains their parallel scalability for in-memory tasks. For out-of-core tasks, it extends the data partitioning parallel model of GC [26] with multi-core parallelism and shared-memory synchronization.

(2) Partitioning and scheduling (Section 4). We study a new problem, which aims to overlap CPU and I/O operations, balance the use of limited memory and cope with the dynamic behaviors of iterative rounds; these new challenges are not encountered by graph partitioners for multi-machine systems. We show the intractability of the problem, and develop an effective joint partitioning and scheduling strategy. As preprocessing, we partition the graph into small blocks; at runtime, we adapt to available memory and system bottleneck dynamically via grouped block processing.

(3) Experimental evaluation (Section 5). Using real-life and synthetic graphs, we empirically find the following. For weakly connected components (WCC), single-source shortest path (SSSP), PageRank (PR), vertex coloring (Coloring), minimum spanning tree (MST), and random walk (RW), (a) Planar outperforms the state-of-the-art (SOTA) out-of-core systems by 34.42 \times on average, up to 302.01 \times . (b) On average it is 5.62 \times faster than the SOTA in-memory system. For parallelly scalable PRAM algorithms, it beats the SOTA by 5.91–9.58 \times ; with 6 \times cores, it speeds up by 3.36 \times . (c) Its adaptive partitioning and scheduling strategy improves performance by 1.87–2.12 \times . (d) It performs as well as the SOTA multi-machine systems that use 4–10 machines, saving the monetary cost by at least 81.7%.

We discuss PRAM in Section 2 and future work in Section 6. We defer proofs and Planar programs to [2] for the lack of space.

Related work. We categorize the related work as follows.

Parallel models. Several models are in place for graph analytics. (1) *Vertex-centric* (VC) [58, 63, 71] and *edge-centric* (EC) [57, 68, 91] models parallelize computation centered around graph neighborhoods, by programming from the perspective of a single vertex/edge. This makes some graph algorithms inefficient for, e.g., graph simulation [23]. (2) *Graph-centric model* (GC) [21, 26] parallelizes sequential graph algorithms across subgraphs, for user to think like a graph. (3) *Hybrid model* [92] adopts the data partitioning parallelism of GC and operation-level parallelism of VC at a single machine.

PRAM [27, 32, 76] supports SIMD parallelism for general computation. It facilitates interprocessor communication and synchronization via shared memory. However, several practical difficulties arise in mapping PRAM algorithms onto real-life physical machines [7], e.g., the fixed number of CPU cores and limited memory capacity. Some programming models, e.g., ICE [30] and XMTC [50], allow users to write lockstep programs similar to PRAM algorithms. These models, however, do not support out-of-core computing.

Planar proposes a parallel model to fit single-machine shared-

memory parallelism. As opposed to message passing-based VC/EC, it supports subgraph-based processing beyond neighborhood, and synchronizes via shared memory, reducing redundant work and I/O. Moreover, it simplifies parallel graph programming by reusing existing PRAM algorithms and retaining their parallel scalability.

Planar extends GC in the following. GC was designed for multi-machine systems that load all subgraphs just once to different machines and process the subgraphs simultaneously via message passing. It supports neither intra-subgraph parallelism nor out-of-core computation. In contrast, Planar targets multi-core parallelism at a single machine. It separates (a) intra-subgraph parallelism via SIMD parallelism and shared-memory synchronization of PRAM, from (b) inter-subgraph parallelism by simulating the fixpoint model of GC and partitioning/scheduling graphs for out-of-core tasks. It supports both in-memory and out-of-core computations.

Planar simplifies the VC programming and execution of hybrid [92] by reusing PRAM programs, retaining their parallel scalability, and demanding neither code revamp nor manual tuning of hybrid.

Single-machine systems. (1) *In-memory ones* [31, 56, 63, 71, 85, 86, 88] assume that a graph can be loaded entirely into memory. They adopt variants of VC/EC [63, 85, 86], and improve data locality via scheduling [85]. (2) *Semi-external systems* (Blaze [45] and [54, 87]) fit all vertex data in memory, and load (immutable) edge data from the secondary storage on-demand. They cannot handle graphs with a large number of vertices. (3) *Out-of-core systems* [6, 45, 49, 54, 57, 68, 77, 87, 91] use disk as memory extension, and focus on reducing the I/O cost of swapping data between the disk and memory. CLIP [6] adopts an asynchronous model to reduce redundant synchronization cost of EC, which may compromise the correctness. All the previous systems adopt VC/EC, except MiniGraph [92] that employs a hybrid model. (4) *Hardware-accelerated systems* [18, 44, 57, 79, 83, 89] leverage GPUs or FPGAs to accelerate graph computations. However, this hardware introduces additional costs and programming complexity, deviating from the cost-efficient and general-purpose design objectives. In contrast, Planar emphasizes affordability and simplicity, providing a competitive performance without relying on specialized hardware [2].

Planar supports a new parallel model to speed up in-memory and out-of-core graph computations, outperforming SOTA of both types (Section 5). For out-of-core execution in particular, it proposes an adaptive strategy for partitioning and scheduling, to cope with the dynamic runtime behavior, and reduce I/O and parallel runtime. These challenges are not encountered by in-memory systems.

Multi-machine systems [19–21, 26, 35, 55, 62, 74, 78, 81, 90] support big graph analytics by scaling out. Such systems adopt a shared-nothing architecture: they partition the input graph, and load the fragments to the machines at once; all workers process their local fragments in-memory in parallel, and synchronize via message passing. The communication cost and workload balancing among workers are thus two vital issues to performance. Rather than scaling out with multiple machines, Planar seeks cost-effective scaling up. It meets new challenges, e.g., limited memory and excessive I/O.

Graph partitioning. For multi-machine systems, the topic has been well studied (see [12, 14] for surveys). (1) Edge-cut [8, 40, 41, 47, 72] partitions vertices into disjoint sets and cuts edges. It promotes

Algorithm 1: Algorithm \mathcal{A} for WCC (Shiloach *et al.* [70]).

Status Declaration: $S_V = \{\bar{p}\}$ where $p(v) = v$ for each $v \in V$;
Input: Graph $G = (V, E, L)$. /*vertex represented by numeric ID. */
Output: The number of weakly connected components in G .
 1 **while** E_i not empty **do**
 2 **parallel for each** $e = \langle u, v \rangle \in E$ **do** Graft($\langle u, v \rangle$);
 3 **parallel for each** $v \in V$ **do** PointerJump(v);
 4 **parallel for each** $e = \langle u, v \rangle \in E$ **do** Contract($\langle u, v \rangle$);
 5 **return** the number of distinct values in \bar{p} ;
Procedure Graft($\langle u, v \rangle$):
 6 **if** $p(u) \neq p(v)$ **then**
 7 swap u and v if $p(u) > p(v)$; /* ensure $p(u) \leq p(v)$. */
 8 $p(p(v)) := p(u)$; /* graft the pseudo-tree of v to that of u . */
Procedure PointerJump(v):
 9 **repeat** $u := p(v)$; $p(v) := p(u)$;
 10 **until** $p(u) = u$; /* halt if u is the root of its pseudo-tree. */
Procedure Contract($\langle u, v \rangle$):
 11 **if** $p(u) = p(v)$ **then** remove edge $\langle u, v \rangle$;

locality but may lead to imbalanced fragments [33]. (2) Vertex-cut [18, 37, 59, 65, 66, 84, 89] partitions edges into disjoint sets and allows mirrored vertices. It balances partitions at the cost of locality [15]. Recent out-of-core systems, *e.g.*, [57, 91], employ a 2D partitioner, which enables fast indexing with massive border vertices. (3) Hybrid [9, 15, 17, 22, 52, 90] strikes a balance by combining the two. Representative heuristics include MDBGP [9] and an application-driven partitioner [22], designed for VC and GC, respectively.

The conventional partitioners mostly aim to reduce the replication factor and the balancing ratio. As will be seen in Section 4, these are not of primary concerns for out-of-core systems; in contrast, Planar tackles a unique joint partitioning and scheduling problem. (1) It develops a partitioner by advocating connectivity among subgraphs and locality within a subgraph, not the balancing ratio. (2) It conducts subgraph grouping and scheduling decisions adaptively at runtime, a mechanism not considered by prior partitioners.

2 PRELIMINARIES

This section reviews basic notations and PRAM algorithms.

Graphs. Assume a countably infinite alphabet Ω for labels. Consider graph $G = (V, E, L)$, directed or undirected, where V is a finite set of vertices; $E \subseteq V \times \Omega \times V$ is a finite set of edges, such that each edge is labeled with a label in Ω ; moreover, each vertex v in V carries a label $L(v) \in \Omega$ to represent properties.

Partition strategies. Given a graph G and a number m , a graph partitioner \mathcal{P} partitions G into *fragments* $\mathcal{F} = (F_1, \dots, F_m)$ such that each $F_i = (V_i, E_i, L_i)$ is a subgraph of G , $E = \bigcup_{i \in [1, m]} E_i$ and $V = \bigcup_{i \in [1, m]} V_i$. We use $F_i.B$ to denote the set of *border entities* (vertices and edges) that are shared by at least two subgraphs.

PRAM. PRAM is a theoretical model that simplifies the design and analysis of parallel algorithms, particularly in a shared-memory environment. It abstracts the complexities of hardware, assuming the availability of a large number of processors and unlimited shared memory. It supports SIMD parallelism among all processors in synchronization, and allows memory random access in constant time.

Decades of research have developed a rich set of PRAM graph algorithms. Compared to algorithms designed for message-passing

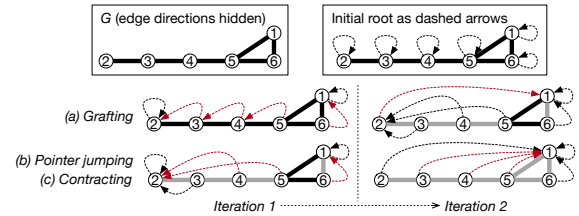


Figure 1: WCC computation over sample graph G with Algorithm 1.

models, they are often provably scalable and more efficient, and make a better fit to multi-core parallelism with shared memory of a single machine. For a class \mathcal{Q} of graph queries, a PRAM algorithm \mathcal{A} takes as input a query $Q \in \mathcal{Q}$ and graph G . It is specified by the following for data, parallelism and computation logic [27, 32, 76].

(1) *Status declaration.* Given G as a set of data arrays for (V, E, L) , algorithm \mathcal{A} declares and initializes a set of *status variables*, which are auxiliary data structures used by \mathcal{A} that represent the intermediate states of G . These include (a) variables associated to individual graph elements, *i.e.*, a set of *vertex* (resp. *edge*) *status*, denoted by S_V (resp. S_E), one for each vertex (resp. edge) in G ; and (b) variables that represent the overall state of G , *i.e.*, a set S_G of *global status*, which is particularly important for coordinating the overall control flow, *e.g.*, maintaining counters or flags that influence algorithm termination. We use $S(G)$ to denote tuple (S_V, S_E, S_G) , referred to as the *status* of G *w.r.t.* Q , keeping track of the computation.

(2) *Processor allocation.* With intermediate state $R(G) = (G, S(G))$ in shared memory, \mathcal{A} conceptually assigns each processor to a unique memory location (*e.g.*, a vertex or an edge) for SIMD parallelism, allowing read, compute or write operations on $R(G)$. Note that \mathcal{A} assumes n processors where n is a polynomial in $|V|$ and $|E|$.

(3) *Lockstep.* Algorithm \mathcal{A} specifies its logic of computation in a sequence of operations for SIMD parallelism, possibly with conditionals and loops. Each operation is called a *lockstep* executing SIMD instructions with synchronization enforced by a barrier at the end of each lockstep. The locksteps produce the final state $R'(G) = \mathcal{A}(Q, G)$ at the end of the sequence, which yields $Q(G)$.

Example 1: Consider WCC. Given a graph $G = (V, E, L)$, it counts the number of maximum subgraphs of G in which all vertices are connected to each other via a path, regardless of the edge direction.

The PRAM algorithm \mathcal{A} of [70] is shown Algorithm 1. Using $|E| + |V|$ processors, it computes WCC of G in $O(\log |V|)$ time. It maintains a disjoint set of *pseudo-trees*. A pseudo-tree rooted at vertex r , denoted by $\Lambda(r)$, is a tree-like structure where each vertex v has a parent $p(v)$ that points to another vertex in $\Lambda(v)$, except that the root r is its own parent, *i.e.*, $r = p(r)$. Algorithm \mathcal{A} iteratively merges these pseudo-trees based on edge connectivity, maintaining the invariant that all vertices in the same pseudo-tree are weakly connected. On graph G of Figure 1, it works as follows.

(1) *Status.* Algorithm \mathcal{A} declares a *parent pointer* $p(v)$ for each $v \in V$, initialized to itself and represented as dashed arrows in Figure 1. The pointers are stored as an array \bar{p} within the vertex status S_V .

(2) *Processors.* It virtually allocates a processor to each edge in E and each vertex in V , allowing all to be processed in parallel.

(3) *Locksteps.* \mathcal{A} maintains and updates parent pointers $\bar{p}(v)$ in a loop that comprises three lockstep operations. Figure 1 illustrates

changing parent pointers as the red dashed arrows in each lockstep, and removed edges in light colors. (a) *Graft* (Line 2). All edges $e = \langle u, v \rangle$ are checked in parallel. If u and v belong to different pseudo-trees $\Lambda(r)$ and $\Lambda(r')$, respectively, a merge is performed by grafting $\Lambda(r)$ onto $\Lambda(r')$, thus forming a larger pseudo-tree. In Figure 1, Iteration 1 of grafting results in two pseudo-trees rooted at 1 and 2, while Iteration 2 further merges them into one. (b) *Pointer jump* (Line 3). The parent pointers for all vertices are updated in parallel. For each vertex v in a pseudo-tree $\Lambda(r)$, the parent pointer $p(v)$ is updated so that $p(v) = r$. This ensures that all vertices in a pseudo-tree point directly to the root (see Figure 1). (c) *Contract* (Line 4). All edges are checked in parallel, to remove ones internal to a pseudo-tree.

Algorithm \mathcal{A} continues to iterate through these locksteps until no edges remain in G (Line 1). At this point, it returns the number of distinct pseudo-tree roots as WCC of G (Line 5). \square

As shown above, compared to VC/EC programs, PRAM algorithms (1) support beyond-neighborhood direct memory accesses, not restricting the operations within the neighborhood of each vertex; (2) exploit shared memory for synchronization, not via message passing; and (3) feature inherent load balancing among processors. Moreover, unlike VC/EC that may require nontrivial efforts in algorithm development, PRAM is backed by a rich set of existing algorithms, which simplifies both programming and debugging.

3 A PARALLEL COMPUTATION MODEL

This section introduces the parallel model of Planar, including how to program (Section 3.1) and execute (Section 3.2) a graph algorithm. We also discuss its benefits over prior parallel models (Section 3.3).

3.1 Programming with Planar

We aim to “plug” existing PRAM algorithms in a single-machine system. However, this is nontrivial. PRAM assumes (a) unlimited memory and a unit access cost, and (b) a polynomial number of cores such that one can process all edges with different cores in parallel. These are beyond the reach of a real-life machine; *e.g.*, when graphs are too large to fit into its memory, we have to use secondary storage as memory extension; with this comes I/O cost.

To close this gap, we propose a parallel model to make practical use of PRAM graph algorithms. We (1) decompose out-of-core computation into in-memory tasks; and (2) run a PRAM algorithm on each in-memory task to leverage multiple cores and shared memory. This is enabled by a simple, high-level programming abstraction.

For a query class \mathcal{Q} , the user needs to provide three functions: (1) PEval, an existing batch PRAM algorithm that evaluates queries $Q \in \mathcal{Q}$ on a subgraph; (2) IncEval, an existing incremental PRAM algorithm that refines partial results with border updates; and (3) Assemble, an aggregator for final results. To further simplify programming, each of these functions can be implemented using high-level primitives specialized for PRAM, *i.e.*, concurrent data structures and PRAM operators. The handling of partitioned graphs and synchronization is mostly hidden from users; the only addition is the definition of functions to resolve conflicts in border updates.

PEval. Batch function PEval takes as input a query $Q \in \mathcal{Q}$ and a subgraph F_i of G ($i \in [1, m]$); it computes the partial result $Q(F_i)$ at state $R_i = \mathcal{A}(Q, F_i)$ on F_i by PRAM algorithm \mathcal{A} . More specifically, PEval initializes *partial status* $S(F_i) = (S_{V_i}, S_{E_i}, S_G)$,

Algorithm 2: Planar program for WCC.

Status Declaration: $S_V = \{\bar{p}\}$ where $p(v) = v$ for each $v \in V$;
StatusAggr: $(p(v), p'(v)) \Rightarrow \min\{p(v), p'(v)\}$.
Function PEval (subgraph $F_i = (V_i, E_i, L_i)$) :
1 **while** E_i not empty **do**
2 EApply($(\forall e \in E_i) \Rightarrow \text{Graft}(e)$);
3 VApply($(\forall v \in V_i) \Rightarrow \text{PointerJump}(v)$);
4 EApply($(\forall e \in E_i) \Rightarrow \text{Contract}(e)$);
5 **return** state R_i on F_i ;
Function IncEval (partial result R_i , updates $\Psi[F_i]$) :
6 VApply($(\forall v \in \Psi[F_i](\bar{p})) \Rightarrow p(p(v)) := \Psi[F_i](p(v))$);
7 VApply($(\forall v \in V_i) \Rightarrow \text{PointerJump}(v)$);
8 **return** updated partial state R_i ;
Function Assemble (partial results R_1, R_2, \dots, R_m) :
9 **return** the number of distinct values in \bar{p} ;

where S_{V_i} (resp. S_{E_i}) is the status variables associated with vertices (resp. edges) in F_i . The partial evaluation process evaluates Q over F_i ; it returns updated status $S(F_i)$ as part of state R_i , keeping track of the computation. It also extracts *round updates*, which consist of changes to the border vertices/edges and their status. As PEval concludes on subgraph F_i , round updates are aggregated into a global cache Ψ , which resolves the conflicting updates to border entities.

PEval may implement an existing batch PRAM algorithm \mathcal{A} for Q . One only needs to extend it with the following.

(1) *Declare status in \mathcal{A} .* Function PEval declares the status $S(G)$ of G , by making use of concurrent data types. Status $S(G)$ includes

- vertex status S_V (resp. edge status S_E) in an array of length $|V|$ (resp. $|E|$), indexed by vertex (resp. edge) identifiers; and
- global status S_G as variables that are globally accessible.

Given subgraph F_i , PEval initializes *partial status* $S(F_i) = (S_{V_i}, S_{E_i}, S_G)$, where S_{V_i} (resp. S_{E_i}) is the subset of variables in S_V (resp. S_E) that are associated with vertices/edges in F_i . It maintains the *partial state* $R_i = (F_i, S(F_i))$, which is persisted onto secondary storage upon task completion to keep track of the computation.

(2) *Specify aggregators.* PEval on different subgraphs may make conflicting updates to status variables of border entities. To resolve the conflicts, PEval specifies two aggregation functions: (a) StatusAggr for status variables, and (b) MutateAggr for edge mutations. At the end of PEval, Planar aggregates border updates by applying both.

(3) *Implement function body.* Function PEval is essentially PRAM algorithm \mathcal{A} . It applies a sequence of synchronized parallel operations over F_i and initial status $S(F_i)$ to produce state R_i , where direct and concurrent memory accesses are granted for each operation.

The lockstep operations of \mathcal{A} are directly streamlined as if users were programming sequentially, using PRAM operators:

- (a) VApply(f_V): vertex parallel operator, which applies a function f_V to a set of vertices in parallel. Here f_V is the procedure of \mathcal{A} for processing each vertex v . It can access any variables relevant to v in partial status $S(F_i)$, and may mutate the graph (see below).
- (b) EApply(f_E): edge parallel operator, similar to VApply.

Both parallel operators are synchronized; we place an implicit synchronization barrier after each invocation of VApply and EApply. Within each parallel operator, all reads to R_i precede any write.

To support PRAM algorithms that transform the topological structure, functions f_V and f_E may invoke an additional primitive:

- (c) Mutate(e, e'): replace an existing edge e with a new edge e' .

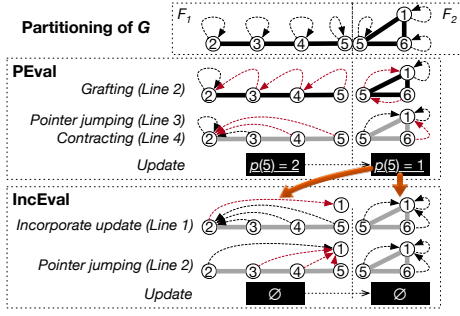


Figure 2: WCC computation over partitioned G with Planar.

Example 2: Now we show the PEval program for \mathcal{A} of Example 1. As shown in Algorithm 2, PEval (1) declares status along the same lines; (2) defines a StatusAggr function for S_V ; and (3) implements the three locksteps using VApply and EApply (Lines 1–4).

Figure 2 depicts the execution of PEval over G . Assume that G is partitioned into subgraphs F_1 and F_2 via vertex-cut (by “cutting” through vertex v_5), such that either can be processed completely in memory. Planar first processes F_1 , produces a single pseudo-tree rooted at v_2 , and generates an update in a set Ψ indicating v_5 ’s parent to be v_2 ; it then processes F_2 along the same line, and aggregates updates to $p(v_5)$ by taking the minimum. It finishes with $\Psi = \{p_\Psi(v_5) : v_1\}$ and subgraphs as pseudo-trees of height 1 (indicated by \bar{p}); function IncEval will later take these as input. \square

IncEval. The incremental function takes as input query Q , subgraph F_i , stale partial state R_i and relevant subset $\Psi[F_i]$ of updates. It updates the partial state incrementally in-place via $R_i = \mathcal{A}(Q, R_i \oplus \Psi[F_i])$, where $R_i \oplus \Psi[F_i]$ denotes integration of $\Psi[F_i]$ with R_i ; it adopts incremental evaluation so as to make maximum reuse of the last-round computation. In the end, Planar resets stale values in Ψ to a clean slate for the next round of IncEval execution.

Function IncEval implements an incremental PRAM algorithm \mathcal{A}_Δ for query class \mathcal{Q} , sharing the status declaration and aggregators of PEval. We may deduce \mathcal{A}_Δ from \mathcal{A} following [25] such that \mathcal{A}_Δ guarantees correctness and minimal incrementalization cost.

Example 3: Continuing with Example 2, IncEval implements an incrementalized \mathcal{A} for Planar, where the last round partial results consist of only a set of height-1 pseudo-trees. Algorithm 2 shows IncEval works on F_i as follows: (1) it incorporates aggregated updates in $\Psi[F_i]$ in parallel, such that for each border vertex v in F_i , its aggregated status update $p_\Psi(v)$ overrides $p(v)$, the parent of v ’s parent (Line 6); and (2) a subsequent parallel pointer jumping (Line 7) maintains the pseudo-trees with height 1.

Figure 2 also illustrates the execution of IncEval. Working over F_1 , update $\{p_\Psi(v_5) : v_1\}$ sets $p(v_2) = v_1$; vertex v_1 then becomes a border vertex shared by F_1 and F_2 . Then, pointer jumping changes the parent of all vertices in F_1 to v_1 . IncEval generates no further update to S_V ; thus, it triggers Assemble upon completion. \square

Assemble takes partial state R_i ($i \in [1, m]$) and partitions \mathcal{F} as input, and combines R_i to get the final answer $Q(G)$. It is triggered when IncEval makes changes to neither F_i nor $S(F_i)$ ($i \in [1, m]$).

Example 4: Assemble in Algorithm 2 simply counts distinct roots of all pseudo-trees in \bar{p} , following Line 5 of Algorithm 1. \square

3.2 Parallel Model

Given a query $Q \in \mathcal{Q}$ and a graph G , the parallel model coordinates the execution of PEval, IncEval and Assemble, no matter whether the computation fits into the memory of a single machine or not.

Out-of-core computation. If graph G exceeds the memory capacity, Planar partitions G into m subgraphs $\mathcal{F} = (F_1, F_2, \dots, F_m)$ such that each F_i and its status $S(F_i)$, are small enough to be processed in-memory. While Planar may use any graph partitioners \mathcal{P} [8, 13, 22, 33, 42, 46], we will develop an “optimal” one in Section 4.

Planar executes algorithm \mathcal{A} by decomposing computation over large G into manageable, in-memory PRAM tasks over subgraphs in \mathcal{F} . The tasks are organized in iterative rounds, synchronizing via the shared memory. Each PRAM task is executed one at a time without overlapping computation with others, using all available CPU cores at once. Following the principle of GC [26], the process of task decomposition and synchronization is made transparent.

In-memory computation. This is a special case under the parallel model. When G and its status $S(G)$ for \mathcal{A} fit entirely into memory, Planar can work with partition $\mathcal{F} = (G)$ directly. It has a single PRAM task, by simulating \mathcal{A} over G with all available cores.

Task decomposition. For out-of-core computation, Planar decomposes it into in-memory PRAM tasks over subgraphs, as follows.

Iterative evaluation. Given a query $Q \in \mathcal{Q}$, Planar works iteratively towards query result $Q(G)$, carrying out computation over each subgraph. To simplify the discussion, we adopt the BSP model [75], which separates computation in supersteps (rounds). A round starts with each subgraph being evaluated locally, and concludes with a global synchronization step, where border updates of all subgraphs are aggregated. For $t \geq 1$, a new round $t+1$ cannot start until round t has completed; the updates generated in round t are accessible only in round $t+1$. This ensures that the computation across the entire graph stays synchronized until a fixpoint is reached.

The iterative process has three phases: partial evaluation (PEval), incremental computation (IncEval), and termination (Assemble). Each phase takes a different PRAM function, as follows.

(1) *Partial evaluation.* The first round computes partial result $Q(F_i)$ for each subgraph $F_i \in \mathcal{F}$ by executing function PEval in parallel using all available cores. It also extracts *round updates*, which consist of changes to the border vertices/edges and their status.

(2) *Incremental computation.* Starting from the second round, Planar iteratively carries out one or more incremental evaluation rounds over each subgraph $F_i \in \mathcal{F}$, by IncEval. Intuitively, an IncEval round maintains the partial result at each subgraph, by refining it incrementally in response to border updates of the last round.

(3) *Termination.* If an IncEval round ends up with no change to status variables, Planar triggers function Assemble, which aggregates partial results of all subgraphs and returns query answer $Q(G)$.

Each round. In each round, Planar iterates through all subgraphs and executes one of the three functions, with all available cores. It processes one subgraph in memory at a time without overlapping computations of multiple subgraphs, by executing PRAM on a physical machine with p cores with a size- p thread pool. That is, a round involves (at most) m “independent” tasks to cope with limited memory and CPU cores. Moreover, Planar loads and processes each

subgraph together with its associated status, overlapping computation with I/O to improve CPU and I/O bandwidth utilization. With a subgraph under processing, it preloads the next into memory for buffering and persists the previous one (with the partial state) onto disk (see [2] for more details). This effectively creates a checkpoint for the computation, providing some resilience for failure recovery.

Fixpoint. The out-of-core execution of Planar can be modeled as a fixpoint computation over partition (F_1, F_2, \dots, F_m) of G . Following [26], one can verify that the parallel model of Planar supports all graph computations, which are covered by PRAM. The convergence of the fixpoint computation is guaranteed under PIE model. An assurance theorem and proof are deferred to [2].

3.3 Planar vs. VC/EC and GC

Taking WCC as an example, we compare the parallel model of Planar with prior models. We also develop Planar programs for SSSP, PR, Coloring, MST and RW; each implements a well-established prior PRAM algorithm for the query class, preserving the correctness and efficiency while requiring little modification. Their analyses are consistent (deferred to [2] for the lack of space).

Parallel scalability. We adapt the parallel scalability of [48] to characterize the effectiveness of parallel algorithms. Consider a sequential algorithm \mathcal{A}' for a query class \mathcal{Q} , which takes $t_{\mathcal{A}'}(|Q|, |G|)$ time to answer a query $Q \in \mathcal{Q}$ over graph G ; and a parallel algorithm \mathcal{A} that takes $t_{\mathcal{A}}(|Q|, |G|, p)$ time using p cores. The *speedup* of \mathcal{A} over sequential \mathcal{A}' is $s(|Q|, |G|, p) = t_{\mathcal{A}'}(|Q|, |G|) / t_{\mathcal{A}}(|Q|, |G|, p)$. We say that \mathcal{A} is *parallelly scalable relative to \mathcal{A}'* if for any Q and G , $s(|Q|, |G|, p) / p \geq \epsilon$ for some constant $\epsilon > 0$. Intuitively, it guarantees speedup of \mathcal{A} relative to a “yardstick” sequential \mathcal{A}' . In principle, such \mathcal{A} is able to reduce the cost of \mathcal{A}' with more cores.

Comparison with VC/EC. We start with VC/EC.

Example 5: A common VC/EC algorithm for WCC is HashMin [82]. Given graph $G = (V, E)$, it assigns each vertex a unique ID, and propagates the lowest ID across each connected component via iterative message passing through edges. It takes $O((|V| + |E|)D)$ time when G fits in memory, where D denotes the diameter of G .

In contrast, Planar does $O((|V| + |E|) \log D)$ amount of work [70] (see Examples 2–4). This is because \mathcal{A} shrinks D by half after each round via topological mutations, reducing message propagation. Neither \mathcal{A} nor HashMin is parallelly scalable relative to sequential BFS, as both incur polylog amount of redundant work. This said, Planar guarantees linear speedup for up to $|V| + |E|$ cores [70], but VC does not due to contention over high-degree “hubs” nodes.

For large G with partition $\mathcal{F} = (F_1, \dots, F_m)$, Planar takes at most $\lceil \log \min\{m, D\} \rceil$ rounds, with beyond-neighborhood computation of GC and contracting subgraphs. In contrast, HashMin takes D rounds in the worst case, incurring much more I/O. \square

Benefits. The parallel model of Planar makes a better fit to single-machine graph processing than VC/EC, as demonstrated by Example 5 and analyses with other common graph algorithms [2].

For in-memory workloads, Planar may do less work than VC/EC by taking a more efficient PRAM algorithm \mathcal{A}' moreover, if \mathcal{A} has proven parallelly scalable. Planar retains the property. In contrast, VC/EC struggles with aggregating messages at high-degree “hub”

vertices, which become stragglers and reduce parallel speedup.

Moreover, Planar supports direct beyond-neighborhood data accesses, flexible control flows to skip unnecessary computation, and graph topology mutations. These reduce redundant disk I/O.

Comparison with GC. Designed for multi-machine systems, GC [26] itself is not a good fit for a single-machine system, because it (a) does not support intra-subgraph parallelism and (b) requires expensive message passing for synchronization. As remarked in Section 1, the parallel model of Planar extends GC to a shared-memory multi-core architecture by supporting (1) SIMD intra-subgraph parallelism, (2) out-of-core computation and (3) memory-based synchronization and graph partitioning/scheduling.

4 PARTITIONING AND SCHEDULING

This section develops a graph partitioning and scheduling strategy for Planar. We start with unique challenges introduced by single-machine systems (Section 4.1). We then formalize partitioning and scheduling as an optimization problem and show its intractability (Section 4.2), followed by our strategy for the problem (Section 4.3). We focus on out-of-core Planar computations in this section.

4.1 Challenges

To process a graph G that exceeds the memory capacity of a single machine, Planar partitions G into a set \mathcal{F} of subgraphs. Most conventional partitioning strategies are designed for multi-machine systems. They strive to minimize the replication factor and balancing ratio [15, 22, 26, 33, 58], reducing the communication cost, redundant work and stragglers. On the contrary, Planar synchronizes via the shared memory for which the communication cost is negligible; moreover, it serializes subgraph processing such that workload skewness can hardly slow down computation. This said, single-machine systems introduce the following unique challenges.

Dynamic behavior. The runtime behavior of a Planar program may vary substantially across rounds. A PEval round loads each subgraph F_i from disk, executes PRAM algorithm \mathcal{A} and produces the partial state R_i . In contrast, an IncEval round reads R_i of the last round, runs incremental \mathcal{A}_Δ , updates R_i and persists it on disk. The cost of IncEval depends heavily on the border updates from the last round; the I/O depends on the size of F_i . Moreover, the runtime may change substantially in different IncEval rounds. Given this, how should we partition graph G for the best performance?

In contrast, this is not an issue for multi-machine systems.

Scheduling for CPU- vs. I/O-bound computation. Planar works out-of-core by repeatedly swapping subgraphs in and out of memory. A *round of execution* can be CPU-bound or I/O-bound. It is *CPU-bound* if its total computational cost over all subgraphs is higher than the total I/O cost; in this case, we should prevent I/O operations from blocking computation. Otherwise, it is an *I/O-bound* round, *i.e.*, the I/O dominates the execution cost, for which we should maximize the I/O bandwidth utilization. Both cases require that we schedule the subgraph processing to maximally overlap the CPU and I/O operations, whichever the bottleneck is.

Scheduling is not an issue for a multi-machine system, which (1) amortizes the I/O cost by loading each subgraph just once; and (2) processes all subgraphs at the same time with different machines.

Granularity. Subgraphs should be large enough to improve locality, reduce synchronization and redundant computation. However, overly large subgraphs may consume too much memory, leaving insufficient space for overlapping I/O. Can we balance granularity to minimize overhead within the constraint of limited memory?

This tradeoff is not studied by prior partitioners, since multi-machine systems assume sufficient memory for each worker and partition the graph based on the number of available machines.

These challenges demand a joint optimization effort in both partitioning and scheduling, taking into account the overlapping of CPU and I/O operations as well as the memory constraint.

4.2 Partitioning and Scheduling Problem

Based on a general cost model for a Planar program, we formalize the partitioning and scheduling problem and show its intractability.

Cost model. Consider a partition \mathcal{F} of G . We formulate the round cost in terms of the computational and I/O costs for each F_i .

Round cost. For round j , denote by $C_{\mathcal{A}_j}(F_i, p)$ the computational cost over F_i using p processors, and by $\text{IO}(F_i)$ the I/O cost, which is proportional to the size of the partial state R_i . Assuming that Planar processes subgraphs $\mathcal{F} = (F_1, \dots, F_m)$ in order and overlaps CPU and I/O whenever possible, the round cost is

$$C_j(\mathcal{F}) = \text{IO}(F_1) + \sum_{i=2}^m \max\{\text{IO}(F_i), C_{\mathcal{A}_j}(F_{i-1}, p)\} + C_{\mathcal{A}_j}(F_m, p). \quad (1)$$

Intuitively, this model accounts for the sequential processing of subgraphs and the overlapping of I/O and computation in a pipeline. The longer duration between the computation and loading determines the subgraph cost. The round cost is thus the sum of all.

Peak memory usage. When processing the partial state R_i of subgraph F_i at runtime, let $M_{\mathcal{A}}(F_i)$ denote its peak memory usage. This can be estimated as a function $M_{\mathcal{A}}(F_i) = \mu_{\mathcal{A}}(|R_i|)$, where $\mu_{\mathcal{A}}$ is determined by the space complexity of algorithm \mathcal{A} .

Profiling. Once a Planar program is compiled, we feed in some graphs as profiling tests, to (1) train a binary classification profiler to determine *qualitatively* whether the PEval round is CPU- or I/O-bound on the given machine; and (2) train function $\mu_{\mathcal{A}}(|R_i|)$ as a regression model. The training samples include runtime parameters *e.g.*, subgraph sizes, degree skewness, and border updates.

The test inputs are small in size, of the same type as the real input. In our experiments we used 4 graphs from 0.14–30.14GB in size, and the entire profiling procedure takes 96s. The profiler can accurately predict the bottleneck in >95% cases.

Problem statement. Consider a Planar program \mathcal{A} . Given input graph G , p cores and memory capacity B , we formulate the joint partitioning and scheduling problem as an optimization problem to find a partition $\mathcal{F} = (F_1, \dots, F_m)$ of G , such that if subgraphs are processed in order, the round cost is minimized and memory capacity B is never exhausted during execution. The objective is

$$\begin{aligned} & \arg \min_{\mathcal{F}} C_j(\mathcal{F}), \\ & \text{subject to} \quad M_{\mathcal{A}}(F_i) + M_{\mathcal{A}}(F_{i+1}) \leq B, i \in [1, m-1]. \end{aligned}$$

Its decision problem, denoted by DPSP, is to decide, given Planar program \mathcal{A} , graph G , integer m , memory bound B , and cost threshold η , whether there exists a valid partition \mathcal{F} of G such that if subgraphs are processed in order, the round cost is at most η .

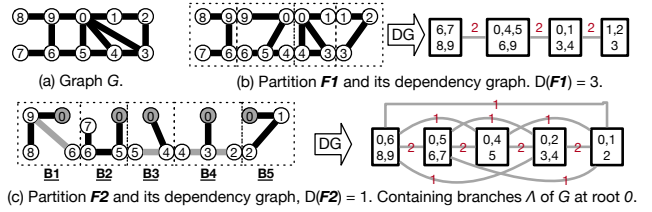


Figure 3: Partitions, branches and dependency graphs.

Theorem 1: DPSP is NP-hard. \square

Proof sketch: We show that DPSP is NP-hard for both I/O-bound and CPU-bound \mathcal{A} (see [2]), by reduction from the NP-complete 3-partition problem (cf. [28]). Given a set A of positive integers, we construct a graph G , an integer m , two positive numbers B and η , algorithm \mathcal{A} (with memory usage $M_{\mathcal{A}}(F_i)$ and round cost $C_j(\mathcal{F})$), such that the set A can be partitioned into disjoint subsets A_1, A_2 and A_3 of equal sum iff there exists a partition \mathcal{F} of G with m subgraphs that meets both the memory bound B and the cost threshold η . \square

4.3 Partitioning and Scheduling Strategies

Theorem 1 suggests that even with accurate cost estimations, an optimal partitioning and scheduling strategy still remains intractable. Hence, we seek an efficient heuristic method that intuitively optimizes performance based on observed workload characteristics.

Overview. We adopt a joint partitioning and scheduling strategy. The idea is to (1) partition G speculatively into small “blocks” during preprocessing, (2) group blocks into subgraphs that fit in memory at runtime, and (3) schedule adaptively based on the bottleneck. In other words, by processing graphs in manageable blocks, Planar efficiently adjusts partitioning/scheduling based on real-time measurements, circumventing the need for precise cost estimations.

At preprocessing, Planar decomposes G into a collection of small blocks, to reduce the runtime decision space and the complexity of grouping. It uses a new locality-aware branching technique optimized for block connectivity and locality, to speedup processing by facilitating update propagation and reducing border status.

To group blocks at runtime, Planar follows a state-of-the-art partitioner to minimize replication and boost locality by reducing border entities. As opposed to working with “static” subgraphs, Planar (1) makes block grouping decisions adaptively at runtime to cope with its dynamic behavior; (2) adopts different scheduling strategies for CPU- and I/O-bound rounds to accommodate the bottleneck; and (3) adjusts granularity based on runtime memory usage to balance the “on-stage” and “off-stage” working memory.

Speculative partitioning. We first develop a speculative partitioner that produces small blocks \mathcal{F} of G . To simplify the discussion, we adopt vertex-cut; it can be adapted to edge-cut or hybrid.

Motivation. Based on Equation 1, we identify key design objectives:

- *Connectivity among subgraphs.* Stronger connectivity boosts update propagation across the entire graph, which leads to faster convergence and fewer rounds in fixpoint computation.
- *Locality within a subgraph.* Better locality can (a) reduce the total size of partial states and hence the I/O cost; and (b) lower the IncEval round complexity. This is in principle consistent with minimizing the replication factor for multi-machine partitioners.

To this end, we develop a two-stage partitioning algorithm. In the first stage, we model block connectivity using a *dependency graph* (DG) and apply *branch decomposition* to minimize the diameter of DG(F). The second stage refines blocks through *greedy adjustment*, redistributing border edges and merging small blocks.

Dependency graph models the connectivity among blocks (subgraphs). Such a graph *w.r.t.* $\mathcal{F} = (F_1, \dots, F_m)$ is an undirected graph $DG(\mathcal{F}) = (V_{DG}, E_{DG}, L_{DG})$. It has m vertices, where $v_i \in V_{DG}$ denotes F_i for each $i \in [1, m]$. An edge e_{ij} exists between v_i and v_j if F_i and F_j share some entities; its weight $L_{DG}(e_{ij})$ denotes the number of shared entities. We will see how m is determined shortly.

Denote by $D(\mathcal{F})$ the diameter of $DG(\mathcal{F})$. Intuitively, the smaller $D(\mathcal{F})$ is, the stronger the connectivity is. It takes fewer steps to propagate an update throughout a small-diameter graph, benefiting label-setting algorithms [43], *e.g.*, WCC, SSSP and Coloring.

Example 6: Figures 3b and 3c show two ways of partitioning graph G (Figure 3a). Partition \mathcal{F}_1 has a dependency graph whose diameter is 3; \mathcal{F}_2 has $D(\mathcal{F}_2) = 1$ despite that it has 1 more subgraph. \square

(1) *Branch decomposition* employs a greedy algorithm, denoted by Decompose. It produces an arbitrary number of blocks. The idea is to cut the graph through a high-degree vertex r , boost the connectivity and minimize $D(\mathcal{F})$. It puts a few high-degree vertices on the border, and strives to keep others within a block. We will further reduce border entities via greedy adjustment and block grouping.

Consider *w.l.o.g.* connected $G = (V, E)$. As shown in Algorithm 3, after finding the highest-degree *root* vertex v_r in V , Decompose breaks $V \setminus \{v_r\}$ into a disjoint set Λ of branches with procedure SortBFSBranch (line 2). A branch λ in Λ is a set of vertices on the same branch in a BFS tree rooted at v , and can thus be reconstructed into a block via procedure Expand, by adding edges that are incident to its vertices. BFS traversal balances the height of each branch; this limits the diameter of each expanded block and reduces within-subgraph computation. The algorithm recurses if the expanded subgraph of a branch is too large to make a valid partition (lines 3–4).

Example 7: Continuing with Example 6, Decompose breaks G into five branches (Figure 3c). The edges in light color connect separate branches; they may be subject to adjustment later. \square

(2) *Greedy adjustment*. We next adjust \mathcal{F} by (a) redistributing border entities, and (b) merging blocks that are too small in size. It produces blocks whose number m is arbitrary yet more manageable.

With redistribution, our goal is to reduce border entities. More specifically, we find each edge e incident to a non-root border vertex v , and migrate e to another block with v tentatively. The changes are materialized if the total border entities are reduced.

Merging aims to (a) limit the number of blocks in \mathcal{F} to reduce the complexity of grouping, and (b) avoid small, scattered I/O requests and improve disk bandwidth utilization. We will merge blocks whose estimated memory usage is below a threshold. By default, we set it to 256MB for full bandwidth utilization [73] over various SSDs.

Remark. Speculative partitioning introduces a one-time preprocessing cost, higher than hash-based *e.g.*, VCut [84] and ECut [53]. Nevertheless, as we will show in Section 5, the cost can be amortized over a few rounds of computation. Moreover, we implement incremental partitioning [24] to cope with dynamic graphs.

Algorithm 3: Function Decompose.

Input: A connected graph G , memory budget τ .

Output: Vertex-cut subgraphs \mathcal{F} of G .

```

1  if  $M_{\mathcal{A}}(G) \leq B$  then return  $\{G\}$ ; else init  $\mathcal{F} := \emptyset$ ;
2  find max-degree vertex  $v_r$  in  $G$ ;  $\Lambda := \text{SortBFSBranch}(G, v_r)$ ;
3  while  $\Lambda$  has non-empty head  $\lambda$  and  $M_{\mathcal{A}}(\text{Expand}(\lambda)) > \tau$  do
4  |    $\mathcal{F} := \mathcal{F} \cup \text{BranchGroup}(\text{Expand}(\lambda))$ ; remove  $\lambda$  from  $\Lambda$ ;
5  return  $\mathcal{F}$ ;
```

Procedure Expand (λ):

```
6  |   return  $F = (V(\lambda), E_F)$ , where  $E_F = \{e \in G \mid e.\text{src} \in V(\lambda)\}$ ;
```

Procedure SortBFSBranch (G, v):

```
7  |   tree  $T := \text{BFS}(G, v)$ ; get  $T$ 's branches  $\Lambda := \{\lambda_1, \dots, \lambda_{|N(v)|}\}$ ;
8  |   return sorted  $\Lambda$ , decreasingly in  $M_{\mathcal{A}}(\text{Expand}(\lambda))$ ,  $\forall \lambda \in \Lambda$ ;
```

Block grouping. Given a speculative partition \mathcal{F} and memory budget τ , grouping selects a set of pending blocks to process as a single subgraph. The goal is to (1) promote locality within the group, and (2) limit memory usage to τ . Note that groupings are temporary; each subgraph is persisted as a separate file after computation.

At a high level, this is equivalent to partitioning the dependency graph $DG(\mathcal{F})$ via edge-cut, so as to minimize the total weight of border edges. To this end, we propose an algorithm called Grouping. It adapts the *neighbor expansion* heuristic of [84] to weighted dependency graphs, which guarantees an optimal replication factor.

More specifically, Grouping selects a set \tilde{F} of blocks, starting with a random vertex v_1 in $DG(\mathcal{F})$. It works iteratively until the memory budget τ is exhausted; each iteration adds one block (*i.e.*, a vertex in $DG(\mathcal{F})$) to \tilde{F} . In the i -th iteration, it finds v_i greedily based on

$$v_i = \arg \max_{v \in \mathcal{F} \setminus \tilde{F}} \sum_{j=1}^{i-1} L_{DG}(\langle v, v_j \rangle). \quad (2)$$

Intuitively, Grouping greedily adds blocks to \tilde{F} , to maximally hide border entities inside a group as “internal” entities. It is efficient. In our experiments over large graph `clueWeb` (see Table 1), \mathcal{F} has 328 blocks, and each grouping iteration takes a negligible ≤ 30 ms.

Example 8: Suppose that Planar takes 5-block \mathcal{F}_2 (Figure 3c) as input, with a memory budget τ of two blocks. If we start with block B1, the grouping strategy will group it with B2 and process both as a single subgraph, since the two share the most border entities. \square

Adaptive scheduling. To overlap CPU and I/O operations, we split the working memory into the off-stage area for buffering pending blocks, and the on-stage area for computation. With the strategies above, one question remains open: how can we balance the split, by setting a memory budget τ for off-stage area?

We use different strategies for CPU- and I/O-bound rounds. At each round, the profiler predicts the upcoming bottleneck utilizing the statistics of subgraphs and runtime parameters of previous rounds (if any). This prediction initializes the scheduling strategy, which is adapted to real-time measurements as the round proceeds.

CPU-bound rounds. The goal is to ensure that I/O operations do not block CPU computation. Thus, we attempt to keep CPUs busy at all times, and may allow some gaps in block loading.

To prevent CPUs from starving, we adopt greedy strategies. (1) With an empty off-stage area, we always load the smallest pending block. (2) When the current on-stage area concludes computation, we immediately group all blocks in the off-stage for computation. (3) For the off-stage area, we reserve an *upper bound* $\tau = B/2$, to

Table 1: Graph datasets.

Name	Type	$ V $	$ E $	Mean Distance	Data Size (GB)
friendster [1]	social network	65.6M	1.8B	5.1	28.9
web-sk [67]	Web	50M	1.9B	13.7	32.0
datagen [36]	synthetic	29M	2.6B	12.5	80.7
clueWeb [67]	Web	1.7B	7.9B	65.7	140.6
hyper12 [4]	Web	273M	9B	42.8	143.0

allow sufficient buffering for the next group of blocks. Intuitively, to improve CPU utilization in early rounds, we start with fine-grained groupings; to speed up computation, we fine-tune the granularity. Too small τ often leads to under-utilization of the working memory, while too large τ may result in oscillations in grouping sizes.

I/O-bound rounds. Our objective is to ensure computation not to block I/O. We strive to keep loading from the disk, and may tolerate idling CPUs. To this end, we attempt to load as many blocks as possible and process them as a group, until the aggregate memory usage of the group exceeds budget τ . In contrast to CPU-bound rounds, we set τ dynamically to maximize block grouping and memory utilization. More specifically, we start with a *lower bound* $\tau = B/2$, and gradually increase τ as long as the computation in the on-stage area concludes before the off-stage area is saturated.

Example 9: Continuing with Example 8, suppose that Planar has a working memory for 4 blocks. If the round is CPU-bound, it will start computation with the smallest block B5, buffering B3 and B4 in the off-stage area for grouping. If the round is I/O-bound, it will start a group of two blocks before commencing computation. \square

5 EXPERIMENTAL STUDY

Using real-life and synthetic graphs, we evaluated Planar for its (1) efficiency, (2) partitioning and scheduling strategy, (3) (parallel) scalability, and (4) performance vs. multi-machine systems.

Experimental setting. We start with the setups.

Datasets. We used five real-life and synthetic datasets, as described in Table 1. All graphs have been widely used in prior work, allowing us to have a comparison under similar conditions. Among these, datagen, clueWeb and hyper12 are among the few open-source graphs that cannot fit in the memory of our main testbed. Here datagen is a synthesized dataset in the LDBC [36] benchmark suite; hyper12 is a BFS sample [51] of hyperLink [4]. The two smaller graphs, friendster and web-sk, are similar in size but have different distributions; they can reveal interesting insights. Datasets were formatted and partitioned based on each system’s requirements; no partitioning is required for in-memory workloads.

To study the impact of graph characteristics on system performance, we generated a series of synthetic graphs using gMark [10], of type bibliography and uniprot. Each graph has 100M vertices, with an average degree ranging from 5–15, *i.e.*, 0.5B–1.5B edges.

Baselines. We evaluated four out-of-core systems: VC-based GraphChi [49], EC-based GridGraph [91] and Blaze [45], and MiniGraph [92] with a hybrid model. We omitted Mosaic [57] for its now-discontinued Intel Xeon Phi coprocessor [61], and CLIP [6] because it produces inconsistent results for multiple graph queries.

We tested SOTA in-memory systems Galois [63] and Ligma [71] for (parallel) scalability; we omitted CoroGraph [88] as it cannot handle graph with >4.3B edges. We also tested GPU systems

Table 2: Runtime statistics. Each round in Blaze is an EC superstep.

Query	Dataset	Metric	Planar	MiniGraph	Blaze	Planar _{static}	Planar _{rand}
WCC	web-sk	Time (s)	56.9	237.3 (4.17×)	66.0 (1.16×)	137.8 (2.42×)	67.9 (1.19×)
		# Rounds	2	6	24	10	2
		I/O (GB)	7.4	104.0 (14.05×)	28.5 (3.85×)	55.1 (7.45×)	7.4 (1.00×)
	friendster	Time (s)	54.7	130.3 (2.38×)	91.0 (1.66×)	142.6 (2.60×)	58.4 (1.07×)
		# Rounds	2	3	16	3	2
		I/O (GB)	7.0	54.5 (7.79×)	27.0 (3.86×)	26.1 (3.73×)	7.0 (1.00×)
SSSP	web-sk	Time (s)	19.2	366.5 (19.09×)	107.7 (5.61×)	156.7 (8.16×)	27.6 (1.44×)
		# Rounds	2	18	58	17	3
		I/O (GB)	11.0	201.1 (18.28×)	40.0 (3.64×)	85.8 (7.80×)	20.0 (1.82×)
	friendster	Time (s)	23.1	183.8 (7.96×)	96.8 (4.19×)	99.8 (4.32×)	30.3 (1.31×)
		# Rounds	2	8	31	8	3
		I/O (GB)	7.0	86.0 (12.29×)	36.4 (5.20×)	100.2 (14.31×)	12.0 (1.71×)

Subway [69] and CGGraph [16], and multi-machine Gluon [19] and GraphScope [21]. All systems were tested in default configurations.

We also tested four variants of Planar: (1) Planar_{static}, which disables block grouping (see Section 4.3); (2) Planar_{rand}, which groups blocks randomly, not based on the heuristic of Equation 2; (3) Planar_{par}, by allowing concurrent subgraphs processing; (4) Planar_{persist}, which persists border updates on disks at each round, and (5) Planar_{par+persist}, which enables both mechanism (3) and (4).

We evaluated four alternative partitioning strategies (Exp-2): (1) VCut, a state-of-the-art vertex-cut heuristic [84]. (2) ECut of [53], the edge-cut used by MiniGraph [92]. (3) 2DVCut, a vertex-cut partitioner used by GridGraph [91], Mosaic [57], and GraphScale [18]. (4) 1DVCut, the vertex-cut used by HitGraph [89]. They produce the same number m of blocks as our speculative partitioner does.

Algorithms. We evaluated Planar programs for WCC, PR, Coloring, SSSP, MST and RW (see [2] for implementation details), common graph queries included in various benchmarks [3, 11]. Among these, PR and Coloring are representative algorithms cast from VC/EC; the others have the best known asymptotic complexity for the query. For baselines, we used their out-of-box implementations if available. Since GridGraph does not support Coloring or MST out-of-box, we implemented their VC algorithms [29, 64].

Moreover, we tested various subgraph queries over bibliography, counting k -stars (*i.e.*, the number of papers with at least k authors) and k -hop paths (for paper impact analysis), where $k \in \{3, 4, 5\}$.

For SSSP, we randomly picked 10 vertices and used them as sources for each input graph. For RW, we initiated a walker at every vertex; each walker takes a 5-step walk.

We have validated the correctness and consistency of system outputs. We present the average of each experiment over 5 repetitions. We report results over some graphs; the other results are consistent.

Testbeds. Our main testbed is a workstation with a consumer-grade CPU and limited memory. It is powered by an Intel Core i9-7900X@3.30GHz CPU, with 13.75MB LLC and 20 cores, and 64GB of DDR4-2666 memory. The graphs were loaded from a 1TB WD Blue WDS100T2B0A SATA SSD, which has an average sequential read throughput of 560MB/s. To further test the parallel scalability for in-memory workloads (Exp-3), we used an enterprise-grade server with 512GB of DDR4-2933 memory and 4× Intel Xeon Gold 5320@2.20GHz CPUs, each with 39MB LLC and 26 cores. Unless noted otherwise, all system were tested with default configurations.

Experimental results. We next report our findings.

Exp-1: Efficiency. We first evaluated the efficiency and I/O of

Table 3: Out-of-core system performance (in seconds). GraphChi could not finish within 5 hours for all queries over clueWeb and hyper12.

Query	datagen					clueWeb				hyper12			
	Planar	MiniGraph	Blaze	GridGraph	GraphChi	Planar	MiniGraph	Blaze	GridGraph	Planar	MiniGraph	Blaze	GridGraph
WCC	60.9	85.8 (1.41×)	81.6 (1.34×)	104.4 (1.71×)	2688.6 (44.15×)	465.6	4219.4 (9.06×)	670.9 (1.44×)	7755.3 (16.66×)	203.2	7825.0 (38.51×)	340.5 (1.68×)	>5h (>88.58×)
SSSP	29.1	42.3 (1.45×)	55.2 (1.90×)	146.7 (5.04×)	2483.4 (85.34×)	229.7	1062.9 (4.63×)	506.8 (2.21×)	1415.9 (6.16×)	91.9	2348.9 (25.56×)	115.1 (1.25×)	>5h (>195.87×)
PR	61.0	100.7 (1.65×)	200.6 (3.29×)	185.6 (3.04×)	646.6 (10.60×)	529.3	2131.8 (4.03×)	623.8 (1.18×)	2537.2 (4.79×)	213.6	1175.6 (5.50×)	406.6 (1.90×)	1820.2 (8.52×)
Coloring	153.7	620.9 (4.04×)	190.4 (1.24×)	2105.3 (13.70×)	4381.5 (28.51×)	346.6	1719.1 (4.96×)	576.6 (1.66×)	2238.5 (6.46×)	218.6	15470.5 (70.77×)	346.1 (1.58×)	>5h (>82.34×)
MST	71.6	147.2 (2.06×)	84.8 (1.18×)	118.9 (1.66×)	907.7 (12.68×)	390.6	>5h (>46.08×)	958.9 (2.45×)	>5h (>46.08×)	252.4	>5h (>71.32×)	407.0 (1.61×)	>5h (>71.32×)
RW	19.6	174.1 (8.88×)	99.8 (5.09×)	207.1 (10.57×)	3177.2 (162.10×)	64.0	1941.4 (30.33×)	297.4 (4.65×)	8391.5 (131.12×)	59.6	2067.2 (34.68×)	114.5 (1.92×)	>5h (>302.01×)

Planar versus out-of-core baselines for various queries. Over two small graphs, we imposed a 16GB memory budget using cgroups to study the out-of-core behavior. Table 2 reports the runtime statistics of some algorithms compared with MiniGraph and Blaze, the best performing baselines supporting GC and VC/EC, respectively. Over large graphs, Table 3 reports the performance of all systems.

WCC. From Tables 2–3, we can see the following.

(1) On synthetic datagen (Table-3), Planar beats the four baselines by 1.34–44.15×. Over real-life clueWeb and hyper12, it outperforms MiniGraph by 9.06× and 38.51×, Blaze by 1.44× and 1.68× and GridGraph by 16.66× and >88.58×, respectively, while GraphChi cannot handle large graphs at this scale. Note that Planar is 2.29× faster on hyper12 than on clueWeb, a graph with fewer edges, since hyper12 has much fewer vertices (only 16.1% of clueWeb), allowing more graph contractions during PEval and faster IncEval rounds.

(2) Over two small graphs, Planar is 2.38–4.17× faster than MiniGraph, and 1.16–1.66× faster than Blaze (Table 2). Over web-sk, Planar only takes 2 rounds due to its beyond-neighborhood computation, and its partitioning strategy that promotes connectivity and hence reduces computation rounds. In contrast, MiniGraph and Blaze take 6 and 24 rounds (*i.e.*, supersteps), respectively. Despite the difference in the skewness of two graphs, Planar has a relative consistent performance; other systems vary greatly.

Here we omit the results of GridGraph and GraphChi; they take at least 3.06× longer than Planar over various workloads.

(3) Planar substantially reduces I/O. On web-sk and friendster (Table 2), its disk read is 90.9% and 74.1% less than MiniGraph and Blaze on average, respectively. Besides taking fewer rounds, Planar reduces I/O in IncEval rounds (Figure 4a), because (a) it contracts the graph by removing most edges in PEval and incurring little disk read in the following rounds over clueWeb; and (b) it reduces I/O further by skipping “inactive” subgraphs. This justifies speculative partitioning, in which a small subgraph is more likely to become inactive.

SSSP. As shown in Tables 2–3, Planar outperforms the all four competitors for SSSP over all graphs. (1) On the two large Web graphs, it is 4.63–25.56×, 1.25–2.21× and 6.16–195.87× faster than MiniGraph, Blaze and GridGraph, respectively. It does the job within 4 min, while GraphChi does not finish in 5 hours. (2) On the synthetic datagen, it beats the four baselines 1.45×–85.34×. (3) Planar takes fewer rounds and 91.9% less disk read than MiniGraph on friendster, leading to a 7.96× speedup. The I/O reductions are not as significant as with WCC, since Planar does not contract the graph during SSSP computation. The speedup is greater (19.09×) on web-sk. These verify that Planar is more effective on large-diameter graphs, since its partitioner strives to improve subgraph connectivity. It leads to faster convergence for GC computation; Planar takes only 2 rounds on both graphs, while MiniGraph takes 8–18

rounds. (4) Over friendster (resp. web-sk), Blaze generates 5.20× (resp. 3.64×) of the disk I/O of Planar and takes 4.19× (resp. 5.61×) longer, even though it leverages on-demand, fine-grained (4KB) I/O to reduce disk reads. This highlights the effectiveness of beyond-neighborhood computation in Planar, which substantially reduces redundant computation, leading to 93.5–96.6% fewer rounds.

PR. For PR over all large graphs, Planar beats the baselines consistently despite that they all execute the same algorithm. As shown in Table 3, it outperforms MiniGraph, Blaze, GridGraph, and GraphChi by at least 1.65×, 1.18×, 3.04× and 10.60×, respectively. This is because Planar maximally overlaps computation and I/O, optimized for shared-memory, multi-core concurrent data accesses.

Coloring. As shown in Table 3 and Figure 4b, (1) Planar is over 1.24× faster, takes at least 97% fewer rounds, and generates 51.1% less disk read, compared to EC/VC-based Blaze and GridGraph. The improvements stem from its beyond-neighborhood computation, which reduces redundant coloring fixes. (2) Planar beats MiniGraph by $\geq 4.04\times$, even though they both follow the same principle of GC and execute the same algorithm. Figure 4b reveals two reasons for this: (a) Planar takes 69.2% less I/O in the first round, benefiting from the compact storage format (see Section ??); and (b) it reduces more disk read in later rounds by skipping processing of more blocks.

MST and RW. As shown in Table 3, Planar consistently beats the four baselines for MST and RW. (1) For MST, Planar is at least 2.06×, 1.18×, 1.66× and 12.68× faster than MiniGraph, Blaze, GridGraph, and GraphChi, respectively. This is due mainly to its more efficient PRAM algorithm, which employs the graph-contraction mutations, like WCC. (2) For RW, Planar beats the best performing baseline by 1.92–5.09×, since Planar supports the random walk algorithm of [39], which is more efficient than that of VC/EC and Hybrid.

Subgraph counting. As shown in Figure 4m over bibliography, Planar answers all subgraph queries in 10s. For various star-counting queries, its performance is relatively consistent, beating Blaze by 1.13× on average. For path counting, Planar runs faster over simpler patterns, as expected. Its speedup over Blaze is 2.10–3.22×.

Exp-2: Ablation study. Next we tested the effectiveness of our partitioning and scheduling strategy, as well as other design choices.

Varying τ . We varied the memory τ reserved for I/O buffering. As shown in Figure 4c for WCC, setting $\tau = 0$ or $\tau = B$ effectively disables the overlapping of CPU and I/O operations, causing substantial slowdown. With $\tau = B/2$, Planar performs the best regardless of the partitioner, hence the V-shape of all lines. This justifies our adaptive scheduling strategy (Section 4.3).

Impact of partitioners. We tested Planar with different partitioners on clueWeb. As shown in Figure 4c for WCC with $\tau = B/2$, our partitioner beats the alternatives consistently. It speeds up VCut, ECut,

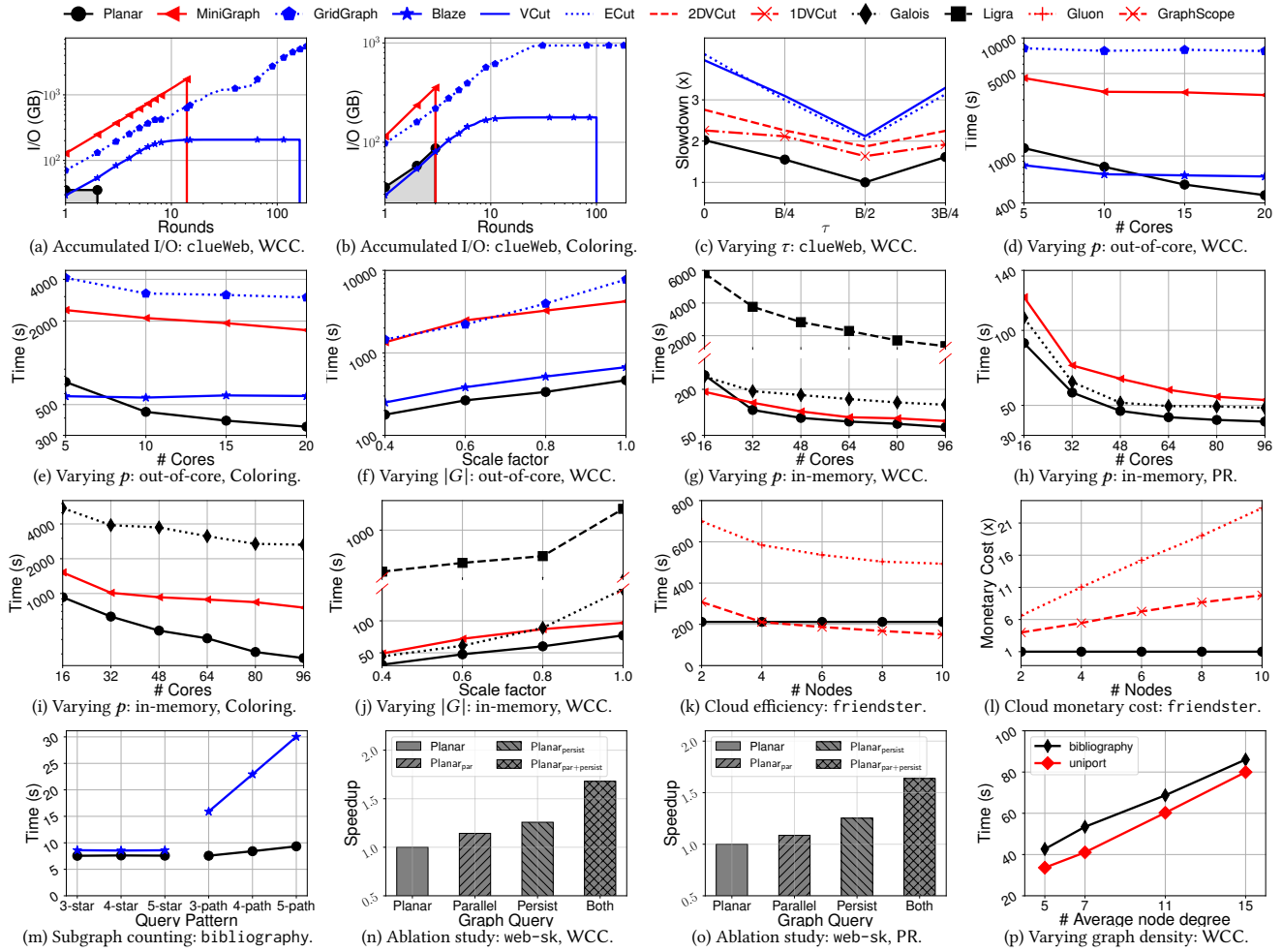


Figure 4: Efficiency, scalability, partitioning and scheduling of Planar, and its performance vs. multi-machine systems.

2DVCut, 1DVCut by 2.12 \times , 2.04 \times , 1.87 \times and 1.63 \times , respectively.

On the server testbed, Planar takes 51.9 min for preprocessing, while the other partitioners take 11.0–16.7 min. This is because speculative partitioning involves local graph traversals and greedy adjustments, which are more computationally intensive than the edge bucketing in other methods. This said, it is an one-time cost amortized over a few rounds of computation. Taking preprocessing into account, Planar beats MiniGraph after answering 1 WCC or 2 Coloring queries. This demonstrates that the partitioning strategy of Planar is effective and shows a quick return on investment.

Effectiveness of scheduling. Both grouped and ordered processing of subgraphs make Planar faster. In our experiments, the trained profiler can accurately predict the round bottleneck in >95% cases; it is 100% accurate for WCC and SSSP. Consider SSSP in Table 2. (1) On friendster (resp. web-sk), Planar_{static} took 4 \times (resp. 8.5 \times) more rounds, generated 14.31 \times (resp. 7.80 \times) more disk read and became 4.32 \times (resp. 8.16 \times) slower. This is because adaptive grouped processing elicits faster convergence, reducing redundant computation of repetitive border updates and allowing to skip some subgraphs in later rounds. (2) Planar beats Planar_{rand} by 1.37 \times on average. This justifies our ordering on subgraphs, which promotes the locality

for subgraph grouping. The results for the others are consistent.

Other techniques. We justified our design decisions to process subgraphs sequentially and cache border updates in-memory. Figures 4n–4o show that if we were to allow multiple subgraphs to be processed concurrently, Planar would take 1.09–1.14 \times longer; if we were to persist the updates to disk, it would be slowed by 1.26 \times . The two collectively can speedup Planar by 1.64–1.68 \times .

Exp-3: Scalability. For both out-of-core and in-memory computation, we evaluated the scalability of Planar with the number p of CPU cores and the graph size $|G|$. We report the results of WCC and Coloring; the results of the other queries are consistent.

Varying p : out-of-core. Scaling the number p of cores, we ran WCC (Figure 4d) and Coloring (Figure 4e) over c1ueWeb. (1) For WCC, Planar scales well with p . It gets a 2.50 \times speedup when p scales from 5 to 20 because its PEval round is CPU-bound. (2) It consistently beats MiniGraph and GridGraph, which barely scale with p due to the I/O-bound rounds. (3) It outperforms Blaze only when $p > 10$, because its PRAM algorithm is more computation-heavy than HashMin [82] (by a constant factor, see Example 5). However, the latter has limited parallelism; Planar speeds up substantially with

more cores, while Blaze barely improves when $p \geq 10$ since it cannot fully utilize the additional cores. (4) For Coloring, Planar is much faster than all baselines for $p > 5$ for similar reasons. It gets 2.10× faster with 4× cores, while other systems speedup <1.64×. Its I/O-efficient GC rounds can benefit more from higher parallelism.

Varying $|G|$: out-of-core. We sampled graphs G from large c1ueWeb using Edge Sampling [51], with a scale factor η for the fraction of edges to be sampled. As shown in Figure 4f when varying η from 0.4 to 1.0 for WCC, Planar scales well with $|G|$. It takes 2.61× longer, while it is 3.14× for MiniGraph, 2.68× for Blaze, and 5.33× for GridGraph. We omit the results of GraphChi, which is much slower.

Varying p : in-memory. Varying p from 16 to 96, Figures 4g–4i report the speedup of Planar and in-memory baselines. (1) For WCC, Planar is up to 1.95× and 28.39× faster than Galois and Ligra, respectively. (2) It is 28.3% slower than MiniGraph when $p = 16$; yet it beats MiniGraph by 1.18×–1.26× with more cores. Consistent with the out-of-core tests (Figure 4d), Planar scales better than MiniGraph with cores. (3) Using 6× more cores for WCC (resp. PR), it gets faster by 3.20× (resp. 2.33×), while it is only 1.98× (resp. 2.28×) for MiniGraph, 1.61× (resp. 2.24×) for Galois. Ligra scales as well as Planar, yet it is slower by magnitudes (omitted in Figure 4h). (4) For parallelly scalable Coloring, Planar speeds up by 3.36× with 6× more cores, much better than MiniGraph (2.02×) and Galois (2.07×). This verifies that when the PRAM algorithm is parallelly scalable, it retains the property for in-memory computations.

Varying $|G|$: in-memory. Using all 104 cores, we further tested in-memory systems by varying the size $|G|$ of graph G . As shown in Figure 4j for WCC, by varying sampling factor η from 0.4 to 1.0 over c1ueWeb, Planar scales better with $|G|$ than VC-based Ligra and Galois; it takes 2.46× longer, while it is 3.37× for Galois and 4.59× for Ligra. This justifies the parallel model of Planar, which supports PRAM algorithms with better asymptotic complexity *w.r.t.* $|G|$. Its scalability is comparable to MiniGraph (1.98×), since both implement efficient algorithms and can scale well with the graph size.

Sensitivity to graph topology. Over synthetic graphs for WCC, Figure 4p shows the performance of Planar under varying graph densities and types. It takes longer for denser graphs, as expected, scaling almost linearly with the number of edges. Given a similar graph size, it performs better over uniprot than more skewed bibliography.

Exp-4: Planar vs. multi-machine systems. We also evaluated the performance and cost effectiveness of Planar versus multi-machine systems GraphScope and Gluon. We deployed all systems in the cloud. More specifically, we ran Planar on a single 8-vCPU 32GB-memory instance. Gluon used instances of the same configuration; GraphScope used multiple 8-vCPU 64GB-memory instances because it requires more than 32 GB in all cases of our experiment.

Resource demands. Single-node Planar supports WCC computation over large graphs; in contrast, multi-machine systems easily run out-of-memory. Over c1ueWeb, for example, GraphScope (resp. Gluon) required at least eight 64GB-memory (resp. four 32GB-memory) nodes. This verifies that Planar lowers the bar of big graph analytics, where large memory capacities are no longer a necessity.

Execution time. As shown in Figure 4k for WCC on friendster, (1) using a single instance, Planar outperforms a 10-node Gluon

cluster by 2.33×, and performs comparably to a 4-node GraphScope cluster. While GraphScope beats Planar when using 6+ instances (each with 2× memory capacity), it requires an additional 200+s for preprocessing, which is not counted in Figure 4k. (2) None of the multi-machine system scales well. Scaling from 4 to 10 machines (using 2.5× cores), Gluon and GraphScope run 1.2× and 1.4× faster, respectively, as the communication cost gradually dominates. In contrast, the more cores are available, the better Planar works.

Cost effectiveness. For WCC on friendster, Figure 4l shows the monetary cost of Planar and multi-machine systems, calculated as the real cloud expense. (1) Although GraphScope and Gluon run faster with more resources, they do not scale cost-effectively, as also observed by [60]. (2) The cost of GraphScope is 5.45× that of Planar for a similar performance; Gluon, running much slower, spends at least 6.61× more. This further justifies the need for a single-machine system to make graph analytics accessible and affordable.

Application. Planar boosts real-world graph analytics for its cost effectiveness. Consider an in-vehicle navigation system, to find the shortest paths between two locations in a large route graph. With limited computing resources, Planar offers an ideal solution, by speeding up route planning and improving user experience.

Summary. We find the following. (1) Planar consistently outperforms the SOTA single-machine systems. It beats out-of-core MiniGraph, Blaze, GridGraph and GraphChi by up to 70.77×, 5.09×, 131.12× and 302.01×, respectively; it can handle workload over large graphs that exceed the capacity of all four baselines. It speeds up in-memory Galois and Ligra by up to 9.58× and 28.39×, respectively. (2) Over various out-of-core workloads, it reduces the I/O cost of the baselines by up to 94.5%. (3) Its partitioner beats prior ones by at least 1.87×, up to 2.12×, and its scheduling strategy consistently speeds up performance. (4) It scales well with large graphs that do not fit in memory. For in-memory computation, on average it is 3.36× faster when using 6× cores for parallelly scalable PRAM algorithms. (5) It requires less memory and is consistently faster than a 10-node Gluon cluster; it performs comparably to GraphScope with 4 machines, saving the monetary cost by 81.7%.

6 CONCLUSION

The novelty of Planar includes the following. (1) Planar is the first graph system that makes practical use of existing PRAM algorithms. (2) It proposes a parallel model that unifies in-memory and out-of-core graph computations, which goes beyond simple simulation of PRAM; it separates inter-subgraph data-partitioned parallelism from intra-subgraph SIMD parallelism to utilize multi-core parallelism, a novel combination. (3) It studies a new graph partitioning/scheduling problem, and develops a strategy that solves the unique challenges not met in multi-machine systems. Our experimental study has validated that Planar is promising in practice.

One topic for future work is to equip Planar with GPU to speed up analytics. Another topic is to fine-tune Planar for a designated task, *e.g.*, graph cleaning, for its best performance.

ACKNOWLEDGMENTS

Yang Liu and Jianxin Li are supported by the National Natural Science Foundation of China through grant No. 62225202.

REFERENCES

- [1] 2024. Friendster dataset. <https://snap.stanford.edu/data/com-Friendster.html>.
- [2] 2024. Full version, with source code availability. <https://shuhaoli.github.io/assets/papers/planar-full.pdf>.
- [3] 2024. Graph500 benchmark specifications. https://graph500.org/?page_id=12#sec-3.
- [4] 2024. Hyperlink. <http://webdatacommons.org/hyperlinkgraph/>.
- [5] Marcos K. Aguilera, Naama Ben-David, Irina Calciu, Rachid Guerraoui, Erez Petrank, and Sam Toueg. 2018. Passing Messages while Sharing Memory. In *PODC*. 51–60.
- [6] Zhiyuan Ai, Mingxing Zhang, Yongwei Wu, Xuehai Qian, Kang Chen, and Weimin Zheng. 2017. Squeezing out all the value of loaded data: An out-of-core graph processing system with reduced disk I/O. In *USENIX ATC*.
- [7] Helmut Alt, Torben Hagerup, Kurt Mehlhorn, and Franco P Preparata. 1987. Deterministic simulation of idealized parallel computers on more realistic ones. *SIAM J. Comput.* 16, 5 (1987), 808–835.
- [8] Konstantin Andreev and Harald Racke. 2006. Balanced graph partitioning. *Theory of Computing Systems* 39, 6 (2006), 929–939.
- [9] Dmitrii Avdiukhin, Sergey Pupyrev, and Grigory Yaroslavtsev. 2019. Multi-Dimensional Balanced Graph Partitioning via Projected Gradient Descent. *PVLDB* 12, 8 (2019), 906–919.
- [10] Guillaume Bagan, Angela Bonifati, Radu Ciucanu, George HL Fletcher, Aurélien Lema, and Nicky Advokaat. 2016. gMark: Schema-driven generation of graphs and queries. *TKDE* 29, 4 (2016), 856–869.
- [11] Scott Beamer. 2016. *Understanding and improving graph algorithm performance*. Ph.D. Dissertation. EECS Department, University of California, Berkeley.
- [12] Charles-Edmond Bichot and Patrick Siarry. 2013. *Graph partitioning*. John Wiley & Sons.
- [13] Florian Bourse, Marc Lelarge, and Milan Vojnovic. 2014. Balanced graph edge partition. In *SIGKDD*.
- [14] Aydin Buluç, Henning Meyerhenke, Ilya Safro, Peter Sanders, and Christian Schulz. 2016. Recent Advances in Graph Partitioning. In *Algorithm Engineering - Selected Results and Surveys*. 117–158.
- [15] Rong Chen, Jiaxin Shi, Yanzhe Chen, Binyu Zang, Haibing Guan, and Haibo Chen. 2019. Powerlyra: Differentiated graph computation and partitioning on skewed graphs. *ACM Transactions on Parallel Computing* 5, 3 (2019), 13.
- [16] Pengjie Cui, Haotian Liu, Bo Tang, and Ye Yuan. 2024. CGraph: An Ultra-fast Graph Processing System on Modern Commodity CPU-GPU Co-processor. *PVLDB* 17, 6 (2024), 1405–1417.
- [17] Dong Dai, Wei Zhang, and Yong Chen. 2017. IOGP: An Incremental Online Graph Partitioning Algorithm for Distributed Graph Databases. In *HPDC*. 219–230.
- [18] Jonas Dann, Daniel Ritter, and Holger Fröning. 2024. GraphScale: Scalable Processing on FPGAs for HBM and Large Graphs. *ACM Trans. Reconfigurable Technol. Syst.* 17, 2, Article 22 (mar 2024), 23 pages. <https://doi.org/10.1145/3616497>
- [19] Roshan Dathathri, Gurbinder Gill, Loc Hoang, Hoang-Vu Dang, Alex Brooks, Nikoli Dryden, Marc Snir, and Keshav Pingali. 2018. Gluon: A Communication-Optimizing Substrate for Distributed Heterogeneous Graph Analytics. In *PLDI*.
- [20] Roshan Dathathri, Gurbinder Gill, Loc Hoang, Vishwesh Jatala, Keshav Pingali, V. Krishna Nandivada, Hoang-Vu Dang, and Marc Snir. 2019. Gluon-Async: A Bulk-Asynchronous System for Distributed and Heterogeneous Graph Analytics. In *PACT*. 15–28.
- [21] Wenfei Fan, Tao He, Longbin Lai, Xue Li, Yong Li, Zhao Li, Zhengping Qian, Chao Tian, Lei Wang, Jingbo Xu, Youyang Yao, Qiang Yin, Wenyuan Yu, Kai Zeng, Kun Zhao, Jingren Zhou, Diwen Zhu, and Rong Zhu. 2021. GraphScope: A Unified Engine For Big Graph Processing. *PVLDB* 14, 12 (2021), 2879–2892.
- [22] Wenfei Fan, Ruochun Jin, Muyang Liu, Ping Lu, Xiaojian Luo, Ruiqi Xu, Qiang Yin, Wenyuan Yu, and Jingren Zhou. 2020. Application Driven Graph Partitioning. In *SIGMOD*.
- [23] Wenfei Fan, Jianzhong Li, Shuai Ma, Nan Tang, Yinghui Wu, and Yunpeng Wu. 2010. Graph Pattern Matching: From Intractability to Polynomial Time. *PVLDB* 3, 1-2 (2010), 264–275.
- [24] Wenfei Fan, Muyang Liu, Chao Tian, Ruiqi Xu, and Jingren Zhou. 2020. Incrementalization of graph partitioning algorithms. *PVLDB* 13, 8 (2020), 1261–1274.
- [25] Wenfei Fan, Chao Tian, Ruiqi Xu, Qiang Yin, Wenyuan Yu, and Jingren Zhou. 2021. Incrementalizing graph algorithms. In *SIGMOD*. 459–471.
- [26] Wenfei Fan, Wenyuan Yu, Jingbo Xu, Jingren Zhou, Xiaojian Luo, Qiang Yin, Ping Lu, Yang Cao, and Ruiqi Xu. 2018. Parallelizing Sequential Graph Computations. *TODS* 43, 18 (2018).
- [27] Steven Fortune and James Wyllie. 1978. Parallelism in random access machines. In *STOC*. 114–118.
- [28] Michael Garey and David Johnson. 1979. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Company.
- [29] Assefaw Hadish Gebremedhin and Fredrik Manne. 2000. Scalable parallel graph coloring algorithms. *Concurrency: Practice and Experience* 12, 12 (2000), 1131–1146.
- [30] Fady Ghanim, Uzi Vishkin, and Rajeev Barua. 2017. Easy PRAM-based high-performance parallel programming with ICE. *TPDS* 29, 2 (2017), 377–390.
- [31] Gurbinder Gill, Roshan Dathathri, Loc Hoang, Ramesh Peri, and Keshav Pingali. 2020. Single Machine Graph Analytics on Massive Datasets Using Intel Optane DC Persistent Memory. *PVLDB* 13, 8 (2020), 1304–1318.
- [32] Leslie M Goldschlager. 1978. A unified approach to models of synchronous parallel machines. In *STOC*. 89–94.
- [33] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. 2012. PowerGraph: Distributed Graph-Parallel Computation on Natural Graphs. In *USENIX OSDI*.
- [34] Raymond Greenlaw, H. James Hoover, and Walter L. Ruzzo. 1995. *Limits to Parallel Computation: P-Completeness Theory*. Oxford University Press.
- [35] Mingyang Han and Khuzaima Daudjee. 2015. Giraph Unchained: Barrierless Asynchronous Parallel Execution in Pregel-like Graph Processing Systems. *PVLDB* 8, 9 (2015), 950–961.
- [36] Alexandru Iosup, Tim Hegeman, Wing Lung Ngai, Stijn Heldens, Arnau Prat-Pérez, Thomas Manhardt, Hassan Chafi, Mihai Capota, Narayanan Sundaram, Michael J. Anderson, Ilie Gabriel Tanase, Yinglong Xia, Lifeng Nai, and Peter A. Boncz. 2016. LDBC Graphalytics: A benchmark for large-scale graph analysis on parallel and distributed platforms. *PVLDB* 9, 13 (2016), 1317–1328.
- [37] Nilesh Jain, Guangdeng Liao, and Theodore L Willke. 2013. Graphbuilder: Scalable graph ETL framework. *Graph Data Management Experiences and Systems* (2013).
- [38] Joseph JáJá. 1992. *An introduction to parallel algorithms*. Addison-Wesley.
- [39] David R Karger, Noam Nisan, and Michal Parnas. 1992. Fast connected components algorithms for the EREW PRAM. In *SPPA*. 373–381.
- [40] George Karypis and Vipin Kumar. 1995. *METIS—unstructured graph partitioning and sparse matrix ordering system, version 2.0*. Technical Report. University of Minnesota.
- [41] George Karypis and Vipin Kumar. 1998. METIS: A software package for partitioning unstructured graphs. *Partitioning Meshes, and Computing Fill-Reducing Orderings of Sparse Matrices Version 4* (1998).
- [42] George Karypis and Vipin Kumar. 1998. Multilevel k-way partitioning scheme for irregular Graphs. *JPDC* 48, 1 (1998), 96–129.
- [43] Arijit Khan. 2017. Vertex-Centric Graph Processing: Good, Bad, and the Ugly. In *EDBT*. 438–441.
- [44] Farzad Khorasani, Keval Vora, Rajiv Gupta, and Laxmi N. Bhuyan. 2014. CuSha: vertex-centric graph processing on GPUs. In *Proceedings of the 23rd International Symposium on High-Performance Parallel and Distributed Computing (Vancouver, BC, Canada) (HPDC '14)*. Association for Computing Machinery, New York, NY, USA, 239–252. <https://doi.org/10.1145/2600212.2600227>
- [45] Juno Kim and Steven Swanson. 2022. Blaze: Fast graph processing on fast SSDs. In *SC*. 1–15.
- [46] Mijung Kim and K Selçuk Candan. 2012. SBV-Cut: Vertex-cut based graph partitioning using structural balance vertices. *Data&Knowledge Engineering* 72 (2012), 285–303.
- [47] Robert Krauthgamer, Joseph Naor, and Roy Schwartz. 2009. Partitioning graphs into balanced components. In *SODA*.
- [48] Clyde P. Kruskal, Larry Rudolph, and Marc Snir. 1990. A complexity theory of efficient parallel algorithms. *Theoretical Computer Science* 71, 1 (1990), 95–132.
- [49] Aapo Kyrola, Guy E. Blelloch, and Carlos Guestrin. 2012. GraphChi: Large-scale graph computation on just a PC. In *USENIX OSDI*.
- [50] Bryant C. Lee, Uzi Vishkin, and George C. Caragea. 2007. Models for Advancing PRAM and Other Algorithms into Parallel Programs for a PRAM-On-Chip Platform. In *Handbook of Parallel Computing - Models, Algorithms and Applications*. Chapman and Hall/CRC.
- [51] Jure Leskovec and Christos Faloutsos. 2006. Sampling from large graphs. In *SIGKDD*.
- [52] Dongsheng Li, Yiming Zhang, Jinyan Wang, and Kian-Lee Tan. 2019. TopoX: Topology Refactorization for Efficient Graph Partitioning and Processing. *PVLDB* 12, 8 (2019), 891–905.
- [53] Yifan Li. 2017. *Edge partitioning of large graphs*. Ph.D. Dissertation. Université Pierre et Marie Curie-Paris VI.
- [54] Hang Liu and H Howie Huang. 2017. Graphene: Fine-Grained IO Management for Graph Computing. In *FAST*. 285–300.
- [55] Yucheng Low, Joseph Gonzalez, Aapo Kyrola, Danny Bickson, Carlos Guestrin, and Joseph M Hellerstein. 2012. Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud. *PVLDB* 5, 8 (2012), 716–727.
- [56] Lingxiao Ma, Zhi Yang, Han Chen, Jilong Xue, and Yafei Dai. 2017. Garaph: Efficient GPU-accelerated Graph Processing on a Single Machine with Balanced Replication. In *USENIX ATC*.
- [57] Steffen Maass, Changwoo Min, Sanidhya Kashyap, Woon-Hak Kang, Mohan Kumar, and Taesoo Kim. 2017. Mosaic: Processing a Trillion-Edge Graph on a Single Machine. In *EuroSys*.
- [58] Grzegorz Malewicz, Matthew H. Austern, Aart J. C. Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: A system for large-scale graph processing. In *SIGMOD*.
- [59] Daniel W. Margo and Margo I. Seltzer. 2015. A Scalable Distributed Graph Partitioner. *PVLDB* 8, 12 (2015), 1478–1489.
- [60] Frank McSherry, Michael Isard, and Derek Gordon Murray. 2015. Scalability!

- But at what COST?. In *HotOS*.
- [61] Timothy Prickett Morgan. 2018. The End of Xeon Phi - It's Xeon and Maybe GPUs From Here. <https://www.nextplatform.com/2018/07/27/end-of-the-line-for-xeon-phi-its-all-xeon-from-here/>.
- [62] Derek G Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martín Abadi. 2013. Naiad: A timely dataflow system. In *SOSP*.
- [63] Donald Nguyen, Andrew Lenharth, and Keshav Pingali. 2013. A lightweight infrastructure for graph analytics. In *SOSP*.
- [64] Gopal Pandurangan, Peter Robinson, and Michele Squizzato. 2018. The distributed minimum spanning tree problem. *Bulletin of EATCS* 2, 125 (2018).
- [65] Fabio Petroni, Leonardo Querzoni, Khuzaima Daudjee, Shahin Kamali, and Giorgio Iacononi. 2015. HDRF: Stream-Based Partitioning for Power-Law Graphs. In *CIKM*.
- [66] Alex Pothen, Horst D Simon, and Kang-Pu Liou. 1990. Partitioning sparse matrices with eigenvectors of graphs. *SIMAX* 11, 3 (1990), 430–452.
- [67] Ryan A. Rossi and Nesreen K. Ahmed. 2015. The Network Data Repository with Interactive Graph Analytics and Visualization. In *AAAI*. 4292–4293.
- [68] Amitabha Roy, Ivo Mihailovic, and Willy Zwaenepoel. 2013. X-stream: Edge-centric graph processing using streaming partitions. In *SOSP*.
- [69] Amir Hossein Nodehi Sabet, Zhijia Zhao, and Rajiv Gupta. 2020. Subway: Minimizing data transfer during out-of-GPU-memory graph processing. In *Eurosys*.
- [70] Yossi Shiloach and Uzi Vishkin. 1982. An $O(\log n)$ parallel connectivity algorithm. *Journal of Algorithms* 3, 1 (1982), 57–67.
- [71] Julian Shun and Guy E. Blelloch. 2013. Lagra: A lightweight graph processing framework for shared memory. In *SIGPLAN*. 135–146.
- [72] George M Slota, Sivasankaran Rajamanickam, and Kamesh Madduri. 2017. *PuLP/XtraPuLP: Partitioning Tools for Extreme-Scale Graphs*. Technical Report. Sandia National Laboratory.
- [73] Vasily Tarasov, Saumitra Bhanage, Erez Zadok, and Margo I Seltzer. 2011. Benchmarking File System Benchmarking: It IS Rocket Science.. In *HotOS*.
- [74] Yuanyuan Tian, Andrey Balmin, Severin Andreas Corsten, Shirish Tatikonda, and John McPherson. 2013. From "think like a vertex" to "think like a graph". *PVLDB* 7, 3 (2013), 193–204.
- [75] Leslie G. Valiant. 1990. A Bridging Model for Parallel Computation. *CACM* 33, 8 (1990), 103–111.
- [76] Leslie G. Valiant. 1990. General Purpose Parallel Architectures. In *Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity*. 943–972.
- [77] Keval Vora, Guoqing Xu, and Rajiv Gupta. 2016. Load the Edges You Need: A Generic I/O Optimization for Disk-based Graph Processing. In *USENIX ATC*.
- [78] Guozhang Wang, Wenlei Xie, Alan J. Demers, and Johannes Gehrke. 2013. Asynchronous Large-Scale Graph Processing Made Easy. In *CIDR*.
- [79] Yangzihao Wang, Andrew Davidson, Yuechao Pan, Yuduo Wu, Andy Riffel, and John D. Owens. 2016. Gunrock: a high-performance graph processing library on the GPU. In *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (Barcelona, Spain) (PPoPP '16). Association for Computing Machinery, New York, NY, USA, Article 11, 12 pages. <https://doi.org/10.1145/2851141.2851145>
- [80] Xianghao Xu, Fang Wang, Hong Jiang, Yongli Cheng, Dan Feng, and Yongxuan Zhang. 2020. A Hybrid Update Strategy for I/O-Efficient Out-of-Core Graph Processing. *TPDS* 31, 8 (2020), 1767–1782.
- [81] Da Yan, James Cheng, Yi Lu, and Wilfred Ng. 2014. Blogel: A Block-Centric Framework for Distributed Computation on Real-World Graphs. *PVLDB* 7, 14 (2014), 1981–1992.
- [82] Da Yan, James Cheng, Kai Xing, Yi Lu, Wilfred Ng, and Yingyi Bu. 2014. Pregel Algorithms for Graph Connectivity Problems with Performance Guarantees. *PVLDB* 7, 14 (2014), 1821–1832.
- [83] Pengcheng Yao, Long Zheng, Xiaofei Liao, Hai Jin, and Bingsheng He. 2018. An efficient graph accelerator with parallel data conflict management. In *Proceedings of the 27th International Conference on Parallel Architectures and Compilation Techniques* (Limassol, Cyprus) (PACT '18). Association for Computing Machinery, New York, NY, USA, Article 8, 12 pages. <https://doi.org/10.1145/3243176.3243201>
- [84] Chenzi Zhang, Fan Wei, Qin Liu, Zhihao Gavin Tang, and Zhenguo Li. 2017. Graph Edge Partitioning via Neighborhood Heuristic. In *SIGKDD*.
- [85] Kaiyuan Zhang, Rong Chen, and Haibo Chen. 2015. NUMA-aware graph-structured analytics. In *PPoPP*. 183–193.
- [86] Yunming Zhang, Mengjiao Yang, Riyadh Baghdadi, Shoaib Kamil, Julian Shun, and Saman Amarasinghe. 2018. Graphit: A high-performance graph DSL. *OOPSLA* 2 (2018), 1–30.
- [87] Da Zheng, Disa Mhembere, Randal Burns, Joshua Vogelstein, Carey E Priebe, and Alexander S Szalay. 2015. FlashGraph: Processing Billion-Node graphs on an array of commodity SSDs. In *FAST*. 45–58.
- [88] Xiangyu Zhi, Xiao Yan, Bo Tang, Ziyao Yin, Yanchao Zhu, and Minqi Zhou. 2023. CoroGraph: Bridging Cache Efficiency and Work Efficiency for Graph Algorithm Execution. *PVLDB* 17, 4 (2023), 891–903.
- [89] Shijie Zhou, Rajgopal Kannan, Viktor K. Prasanna, Guna Seetharaman, and Qing Wu. 2019. HitGraph: High-throughput Graph Processing Framework on FPGA. *IEEE Transactions on Parallel and Distributed Systems* 30, 10 (2019), 2249–2264. <https://doi.org/10.1109/TPDS.2019.2910068>
- [90] Xiaowei Zhu, Wenguang Chen, Weimin Zheng, and Xiaosong Ma. 2016. Gemini: A computation-centric distributed graph processing system. In *USENIX OSDI*.
- [91] Xiaowei Zhu, Wentao Han, and Wenguang Chen. 2015. GridGraph: Large-Scale Graph Processing on a Single Machine Using 2-Level Hierarchical Partitioning. In *USENIX ATC*.
- [92] Xiaoke Zhu, Yang Liu, Shuhao Liu, and Wenfei Fan. 2023. MiniGraph: Querying big graphs with a single machine. *PVLDB* 16, 9 (2023), 2172–2185.