



# MiniGraph: Querying Big Graphs with a Single Machine

Xiaoke Zhu

Beihang University, China  
Shenzhen Institute of  
Computing Sciences, China  
zhuxk@buaa.edu.cn

Yang Liu

Beihang University, China  
Shenzhen Institute of  
Computing Sciences, China  
liuyang@act.buaa.edu.cn

Shuhao Liu\*

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

Wenfei Fan

Shenzhen Institute of  
Computing Sciences, China  
University of Edinburgh  
United Kingdom  
Beihang University, China  
wenfei@inf.ed.ac.uk

## ABSTRACT

This paper presents MiniGraph, an out-of-core system for querying big graphs with a single machine. As opposed to previous single-machine graph systems, MiniGraph proposes a pipelined architecture to overlap I/O and CPU operations, and improves multi-core parallelism. It also introduces a hybrid model to support both vertex-centric and graph-centric parallel computations, to simplify parallel graph programming, speed up beyond-neighborhood computations, and parallelize computations within each subgraph. The model induces a two-level parallel execution model to explore both inter-subgraph and intra-subgraph parallelism. Moreover, MiniGraph develops new optimization techniques under its architecture. Using real-life graphs of different types, we show that MiniGraph is up to 76.1× faster than prior out-of-core systems, and performs better than some multi-machine systems that use up to 12 machines.

### PVLDB Reference Format:

Xiaoke Zhu, Yang Liu, Shuhao Liu, and Wenfei Fan. MiniGraph: Querying Big Graphs with a Single Machine. PVLDB, 16(9): 2172 - 2185, 2023.  
doi:10.14778/3598581.3598590

### PVLDB Artifact Availability:

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

## 1 INTRODUCTION

Big graph analytics has mostly been a privilege of big companies by employing a cluster of machines. For instance, to compute connected components of a graph with billions of vertices and trillions of edges, Yahoo! employs a 1000-node cluster with 12000 processors and 128 TB of aggregated memory [49]. To mine 3-FSM with support = 25000 on a million-edge graph, DistGraph [50] uses 128 IBM BlueGene/Q supercomputers and 32,768 GB memory. A number of graph systems have been developed to explore multi-machine parallelism, e.g., [8, 20, 25, 27, 38, 42, 51, 54, 58, 61, 63]. However, big graph analytics via multiple machines is often beyond reach of small companies, which cannot afford such enterprise clusters.

Moreover, many parallel graph systems “have either a surprisingly large COST, or simply underperform one thread” [39]. For

Table 1: WCC performance on GridGraph and MiniGraph.

| System    | friendster (mean distance: 5.1) |           | web-sk (mean distance: 13.7) |           |
|-----------|---------------------------------|-----------|------------------------------|-----------|
|           | # Supersteps                    | Disk Read | # Supersteps                 | Disk Read |
| GridGraph | 21                              | 135 GB    | 120                          | 367 GB    |
| MiniGraph | 6                               | 74 GB     | 9                            | 82 GB     |

example, single-source shortest path (SSSP) is “essentially not scalable with an increasing number of machines” [56]. This is because multi-machine systems typically adopt the shared-nothing architecture, and the more machines are used, the heavier their communication cost is incurred. In addition, machines in such a system are often under-utilized due to unbalanced workload.

To rectify the limitations of the multi-machine systems, single-machine systems have been studied to explore multi-core parallelism, notably out-of-core systems to process data that is too large to fit into the main memory of a single machine at once [9, 24, 33, 37, 46, 57, 64]. The systems are developed on a machine that has a number of CPU cores, but limited memory capacity and disk I/O bandwidth. To get over the bottleneck, these systems have mostly focused on how to optimize I/O and memory efficiency.

Despite the efforts, disk I/O remains the bottleneck of the single-machine systems. Consider Weakly Connected Components (WCC). Given a graph  $G$ , WCC is to compute the maximum subgraphs of  $G$  in which all vertices are connected to each other via a path, regardless of the direction of edges. Consider GridGraph [64], a state-of-the-art out-of-core system featuring optimal I/O. We run the out-of-box WCC implementation of HashMin [59] that comes with GridGraph as a benchmarking application. The test was conducted on a workstation powered with a 10-core CPU with hyperthreading and 16 GB DDR4 memory. Graph  $G$  is friendster or web-sk [45]; both have  $\sim 1.8$  billion edges, and amount to  $\sim 31$  GB, about twice the memory size. We measure its system I/O using `iostat`. As shown in Table 1, while the two graphs are similar in size, GridGraph incurs drastically different I/O costs. It induces 2.7× disk read on web-sk as on friendster, which is relatively denser.

An in-depth analysis reveals that the excessive I/O often stems from the vertex-centric model [25, 38, 64] adopted by GridGraph for parallel computation, referred to as VC. Under VC, an algorithm employs a user-defined function to process the immediate neighborhood of each vertex (or edge), and exchanges information between vertices via message passing. To send a message from a memory-resident vertex to a memory-absent one, it inevitably requires swapping their data in and out of the memory; as a consequence, VC generally incurs more disk I/O when  $G$  has a larger diameter.

Moreover, while VC is natural for graph algorithms such as PageRank [14] and HashMin [59] (for WCC), it is neither easy to

\*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](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 9 ISSN 2150-8097.  
doi:10.14778/3598581.3598590

write nor efficient to execute *e.g.*, an optimized WCC algorithm via Breath-First Search [28] and graph pattern matching algorithms with subgraph isomorphism or graph simulation [17] under VC.

Can we systematically reduce the I/O cost of an out-of-core graph system and improve multi-core parallelism? Given a computational problem, would other parallel models fit it better than VC?

**MiniGraph.** To answer these questions, we develop MiniGraph, an out-of-core system for graph computations with a single machine. It is the first single-machine system that extends the graph-centric model (GC) of [18, 20] from multiple machines to multiple cores. It shows that GC speeds up beyond-neighborhood computation and reduces I/O, and moreover, simplifies parallel programming.

As shown in Table 1, when computing WCC over friendster (resp. web-sk), the benefit of the beyond-neighborhood computation (GC) is evident: MiniGraph (a) takes 6 (resp. 9) supersteps to converge under the bulk asynchronous model (BSP) [52], as opposed to 21 (resp. 120) steps with GridGraph; (b) reads 74 GB (resp. 82 GB) of data in contrast to 135 GB (resp. 367 GB) of GridGraph; and (c) is less sensitive to the distribution of the input graphs.

However, it is nontrivial to migrate GC to a single-machine system. It introduces new challenges such as memory constraint and I/O cost. MiniGraph approaches the following challenges that are non-existent in prior GC systems: (a) out-of-core computations and in-memory synchronizations; (b) multi-core parallelism; and (c) resource scheduling. It has the following unique features.

*(1) A pipelined architecture.* MiniGraph proposes an architecture that pipelines access to disk for read and write, and CPU operations for query answering. The idea is to overlap I/O and CPU operations, so as to “cancel” the excessive I/O cost, and promote sequential accesses to the disk. Moreover, it employs a shared in-memory data structure for message passing and efficient synchronization. This architecture decouples computation from memory management and scheduling, giving rises to new opportunities for optimizations.

*(2) A hybrid parallel model.* MiniGraph proposes a hybrid model, supporting both the data-partitioned parallelism of GC and the operation-level parallelism of VC. Under the memory constraint, it promotes better utilization of the multi-core resources and avoids fragmentation of the input graph. Moreover, it exposes a unified interface such that the users can benefit from both and can choose one that fits their applications and graphs the best.

*(3) Two-level parallelism.* The hybrid model also enables two-level parallelism: inter-subgraph parallelism via high-level GC abstraction, and intra-subgraph parallelism for low-level VC operations. This presents new opportunities for improving multi-core parallelism. However, its relevant scheduling problem is NP-complete. This said, we develop an efficient heuristic to allocate resources, which is dynamically adapted based on resource availability.

*(4) System optimizations.* MiniGraph develops unique optimization strategies enabled by the hybrid parallel model. It employs a light-weight state machine to model the progress of cores working on different subgraphs, and tracks messages between cores. These allow it to explore shortcuts in the process to avoid redundant I/O.

**Contributions & organization.** After reviewing VC and GC parallel models (Section 2), we present MiniGraph as follows:

- its pipelined architecture and a system overview (Section 3);
- the hybrid parallel model (Section 4);
- the two-level parallel execution model, including resource scheduling (Section 5.1) and system optimizations (Section 5.2); and
- an experimental study (Section 6). Using real-life graphs, we find the following. (a) With various out-of-core workloads, MiniGraph outperforms prior out-of-core systems by 10.8× on average, up to 76.1×. (b) It is able to query a graph of size 10× of the memory capacity. (c) It outperforms the state-of-the-art multi-machine systems when they use up to 12 machines.

## 2 BACKGROUND AND MOTIVATION

In this section, we review parallel graph (programming) models.

Consider graph  $G = (V, E, L)$ , directed or undirected, where  $V$  is a finite set of vertices;  $E \subseteq V \times V$  is a set of edges; each vertex  $v$  in  $V$  (resp. edge  $e \in E$ ) carries label  $L(v)$  (resp.  $L(e)$ ) for properties.

In principle, a graph parallel model determines how users can program with the system and how the programs are executed in parallel. Two types of parallel models have been implemented in (multi-machine) graph systems, namely, VC and GC.

**The vertex-centric model (VC).** Initially proposed by Pregel [38], VC has been the *de facto* go-to model for parallel graph systems. Its principle is for users to think like a vertex: a vertex program is “pivoted” at a vertex; it may only directly access information at the current vertex and its adjacent edges [25, 38]; information is exchanged between “remote” vertices via message passing.

It is natural to program with VC for a problem when its computation can be distributed to vertices and is centered at each vertex, such as PageRank. It is, however, nontrivial to write VC algorithms for problems that are constrained by “joint” conditions on multiple vertices, *e.g.*, graph simulation [40] and subgraph isomorphism [23]. While a variety of conventional sequential algorithms have been developed for such problems, to program with VC, one has to recast the existing algorithms into the VC model. Moreover, such VC programs are often inefficient since they incur heavy message passing. The issue is more staggering for out-of-core systems since we have to repeatedly load the adjacent list of a vertex from disk to memory when the recursive/induction process backtracks.

**Example 1:** As a common benchmarking algorithm for WCC, HashMin works as follows under VC: (1) initially, each vertex is assigned a distinct numeric label; (2) in every iteration, each vertex collects the labels from its neighbors and updates its own with the minimum within the neighborhood; (3) when no more label updates can be made, it returns the total number of distinct labels in the graph.

On graph  $G$  of Figure 1a, HashMin takes 5 supersteps under BSP as shown in Figure 1b, where active vertices are in light color and inactive ones are in gray. We also count the number of message passed in each superstep, marked as “#Ops”. As the computation progresses, more vertices become inactive, and less messages are passed.

Observe the following: (1) The labels of vertices in Subgraph  $D$  are overwritten repeatedly; many initial messages and label updates are redundant. In total 92 messages are passed, while  $G$  has 16 edges only. (2) It takes 5 supersteps. Every edge must be swapped in and out of the memory for 5 times if  $G$  cannot fit in the memory entirely. It is translated to a disk read demand that is 5× of the size of  $G$ . □

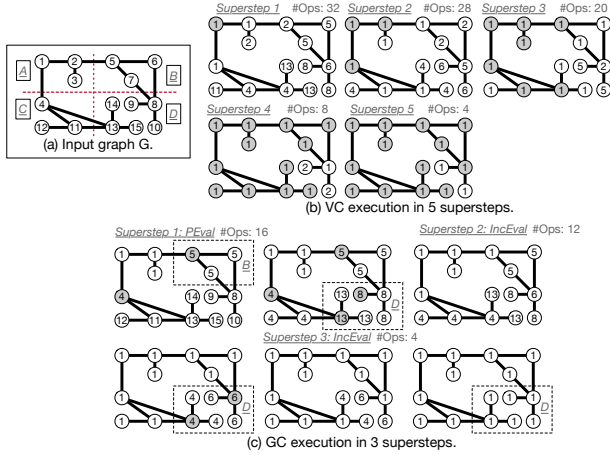


Figure 1: WCC computation on G with VC or GC.

**The graph-centric model (GC).** Proposed by GRAPE [20], GC carries out data-partitioned parallelism. Given a big graph  $G$ , it partitions  $G$  into fragments (subgraphs) using existing graph partitioners (edge-cut [10, 31], vertex-cut [13, 25, 32] or hybrid [19]), and distributes the fragments to different workers. It parallelizes sequential graph algorithms. More specifically, for a graph computational problem  $Q$ , users may provide three (existing) sequential algorithms PEval, IncEval and Assemble, referred to as a PIE program:

- PEval is a conventional algorithm for  $Q$  such that all workers execute it on their local fragments in parallel, and produce partial results (*partial evaluation*); the values of border nodes (e.g., vertices with edges to other fragments) are exchanged as messages;
- IncEval is a sequential incremental algorithm for  $Q$ ; it is repeatedly executed to refine the partial results by treating messages from other workers as updates (*incremental computation*); and
- when no more messages are exchanged, Assemble aggregates partial results from all the workers and forms the final answer.

The parallelized computation starts with PEval, and is followed by iterative IncEval until a fixpoint is reached. Under a generic condition, it guarantees to converge at correct answers as long as the sequential algorithms PEval, IncEval and Assemble are correct [21].

As opposed to VC, GC supports beyond-neighborhood computations. It simplifies parallel programming by parallelizing existing sequential algorithms; moreover, it reduces unnecessary recomputation via iterative IncEval. However, for problems such as PageRank, when graph partitions are not “well balanced” [19], the communication cost of a GC algorithm between workers may be higher than its VC counterpart. Moreover, the parallelism is explored at the subgraph level, not at the vertex/edge-level in each subgraphs.

**Example 2:** As opposed to HashMin, GC parallelizes a more efficient sequential WCC algorithm. After labeling all vertices of graph  $G$  with distinct integers, the algorithm proceeds on each fragment (subgraph)  $F_i$  of  $G$  with single-threaded worker  $W_i$ . (1) Starting with PEval,  $W_i$  initiates BFS from the vertex  $v$  that has the minimum label  $L(v)$  among all the vertices in  $F_i$ . During the BFS, it overwrites their labels with  $L(v)$ . (2) As soon as no vertex remains active,  $W_i$  generates messages comprising the updated labels of border nodes of  $F_i$ . (3) With all updates received from other workers, if  $v$ ’s update

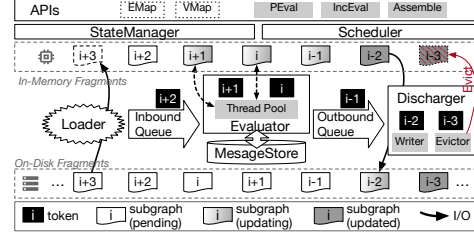


Figure 2: The pipelined architecture of MiniGraph.

carries a smaller label,  $W_i$  overwrites  $L(v)$ ; it invokes IncEval to incrementally refine the partial results until no message is exchanged among workers. (4) Finally, the algorithm calls Assemble to count WCC as the total number of distinct labels from all workers.

Suppose that graph  $G$  is partitioned into 4 subgraphs, as depicted in Figure 1a. GC completes WCC in 3 supersteps under BSP (1 round of PEval and 2 rounds of IncEval), as illustrated in Figure 1c.

Observe the following: (1) GC reduces redundant computation and needs 3 supersteps only. (2) It passes 32 messages on  $G$  (~65% from VC), and requires loading each subgraph only three times (~40% from VC). This said, the computation within each subgraph is conducted *sequentially*, and there is room for improvement. □

To the best of our knowledge, no single-machine systems supports the GC model despite its efficiency and ease of programming.

### 3 MINIGRAPH: AN OVERVIEW

In this section, we present an overview of MiniGraph.

**Challenges.** GC is designed for multi-machine systems [18, 20]. Extending it to a single-machine system introduces new challenges.

- *Out-of-core computation.* A single-machine system has to resort to secondary storage as memory extension when an input graph exceeds its memory capacity. Managing the in-memory and the on-disk parts of an input graph is crucial to performance.
- *Synchronization.* GC passes messages among workers for synchronization. With a shared memory, can a single-machine system synchronize more efficiently via shared data structures?
- *Parallelism.* GC exploits data-partitioned parallelism only. With limited memory capacity, it would result in either under-utilization of the CPU cores or excessive graph fragmentation. Systemically balancing the two factors is challenging.
- *Scheduling.* Work migration incurs low cost among cores in a single machine, due to its shared-memory architecture. Can we leverage flexible resource scheduling to improve performance?

In light of these, MiniGraph proposes a pipelined architecture, a hybrid parallel model of VC and GC, two-level parallelism, and unique optimizations on resource scheduling and I/O.

**Graph partitioning.** Previous single-machine systems organize an input graph as a large number of chunks (*a.k.a.* shards [33], blocks [64]), each consisting of a (possibly small) list of edges. This design serves the sole purpose of gathered I/O, while the system schedules directly on edges. It may not work well with the block storage.

MiniGraph explores a conceptually different approach. It partitions a large graph  $G$  such that one or more fragments can fit into the memory. A fragment is essentially a *subgraph* of  $G$ . MiniGraph adopts a two-level abstraction: (1) at a high level, it sees a subgraph

as the atomic unit for scheduling and I/O; and (2) when processing a subgraph that is in memory, it may carry out parallel computation on the edges and vertices. Our empirical study shows that the two-level abstraction improves multi-core parallelism.

MiniGraph may use any graph partitioners, *e.g.*, [10, 13, 19, 25, 31, 32, 44]. It takes as input a *partitioned graph*, and makes scheduling/execution decisions accordingly. It is nontrivial to find an optimal graph partition; the topic is deferred to the future work.

**A pipelined architecture.** MiniGraph takes as input the subgraphs (fragments)  $F_0, F_1, \dots, F_{n-1}$  of a (possibly large)  $G$ ; the subgraphs are initially stored on disk. MiniGraph iteratively processes the subgraphs in a pipelined architecture, as shown in Figure 2.

This architecture overlaps I/O and CPU operations, and conducts computations on in-memory subgraphs while loading pending ones from the disk. It mitigates the excessive cost of I/O, improving CPU utilization by reducing idle waiting. While prior out-of-core systems implement a similar notion (*e.g.*, [33, 37]), their VC inevitably generates scattered I/O; in contrast, MiniGraph promotes subgraph-based I/O of GC, and consistently accesses the disk *sequentially*.

More specifically, MiniGraph breaks down the out-of-core processing of a subgraph  $F_i$  into three consecutive stages: (1) read  $F_i$  into the memory, (2) compute and update  $F_i$ , and (3) if necessary, write the updated  $F_i$  back to the secondary storage. It has three components: Loader, Evaluator, and Discharger, for the three stages, respectively. They work asynchronously in a pipeline, and are loosely coupled via two task queues InboundQueue and OutboundQueue.

**Loader.** Working in a dedicated thread, Loader continuously selects and reads a memory-absent subgraphs from disk into the memory, as long as the memory capacity is not exhausted. Instead of many scattered reads, Loader issues a small number of bulk read requests, which can improve read throughput and reduce system interrupts.

Once a subgraph  $F_i$  is loaded into the memory, Loader pushes its token  $T_i$  into InboundQueue, where  $T_i$  includes metadata of  $F_i$  such as pointers to the data and its current state.

**Evaluator.** As the downstream component of Loader, Evaluator is responsible for efficient execution of an application. It pops token  $T_i$  from InboundQueue, locates  $F_i$  in memory, and initializes a virtual worker  $W_i$  with reserved thread allocations to execute the program on  $F_i$ . During the execution,  $W_i$  can update  $F_i$  in place, and may generate updates to other subgraphs. The update messages are cached in MessageStore, an array-like storage in memory.

To fully exploit the multi-core parallelism and cache locality, Evaluator manages a global thread pool, whose size is decided by the system hardware concurrency. Whenever threads become idle, it triggers Scheduler, which re-allocates the resources to workers (see below). Once a worker completes, Evaluator reclaims the reserved threads and pushes token  $T_i$  to OutboundQueue for discharging.

**Discharger.** As a counterpart of Loader, Discharger writes the subgraph data back to the disk. It consumes OutboundQueue continuously in a dedicated thread. For a popped token  $T_i$ , it (1) writes the updated  $F_i$  contiguously to disk, and (2) records  $F_i$ 's latest state to StateManager for maintenance. After  $F_i$  is fully discharged, Discharger evicts it from memory and reclaims space for Loader.

**Key modules.** MiniGraph implements the following unique modules for programming, execution and optimization.

**MessageStore.** MiniGraph uses MessageStore, an in-memory data structure, for synchronization among parallel workers, by caching and managing pending messages. In contrast to worker-to-worker message passing in multi-machine systems [18, 20], MessageStore works more efficiently in a shared-memory environment.

We implement it as a compact variable-size array for space efficiency (details in Appendix B of [4]). In practice, it takes only a small portion in memory. For example, under various workloads over web-sk, it consumes 1.3% of the memory at the peak.

**APIs for a hybrid parallel model.** MiniGraph exposes PIE+, a unified interface that integrates VC with GC programming (see Section 4). Users can not only parallelize sequential graph algorithms under GC to simplify parallel programming, but also further explore intra-subgraph parallelism under VC via the new interfaces.

**Scheduler.** Scheduler tracks and allocates threads in ThreadPool in which each thread corresponds to a physical CPU core. It makes decisions to assign physical threads to virtual workers to carry out (parallel) computations on fragments. It also makes active adjustments to support two-level parallelism: when a thread is available, Scheduler allocates it to either a new worker by consuming InboundQueue for speeding up inter-subgraph parallelism, or a running worker for improving intra-subgraph parallelism.

**StateManager.** As computation progresses on the subgraphs of  $G$ , StateManager maintains a state machine to model their status. It is a low-cost method for convergence detection, and it helps identify chances to skip redundant I/O or computation. It is implemented as a lightweight data structure and maintains only a few states per subgraph, taking negligible space in memory.

We will present the details of the hybrid parallel model, two-level parallelism and state manager in Sections 4, 5.1 and 5.2, respectively.

## 4 HYBRID PARALLEL MODEL

In this section we introduce the hybrid parallel model of MiniGraph. It aims to improve multi-core parallelism by enriching inter-subgraph parallelism (GC) with intra-subgraph parallelism (VC).

**Overview.** Given a graph computational problem  $Q$  and a graph  $G$ , MiniGraph processes subgraphs  $F_0, F_1, \dots, F_{n-1}$  of  $G$ . It assigns virtual worker  $W_i$  to  $F_i$  ( $i \in [0, n-1]$ ). To simplify the discussion, we adopt BSP [52] for synchronization. In each superstep, each fragment is processed exactly once. When  $F_i$  is loaded into the memory,  $W_i$  is activated and allocated CPU resources by Scheduler.

Recall from Section 2 that under GC, users provide a PIE program for  $Q$ , which consists of three sequential algorithms PEval, IncEval and Assemble. Each worker  $W_i$  executes PEval on its local subgraph  $F_i$  for partial evaluation; it then iteratively and incrementally refines the partial results via IncEval on  $F_i$  by treating messages between workers as updates, until it reaches a fixpoint, followed by Assemble to aggregate the partial results from all workers into the final answer. In principle, each  $W_i$  runs PEval and IncEval *sequentially*.

The key idea of the hybrid parallel model is two-level parallelism. At the top level, it adopts GC to benefit from data-partitioned inter-subgraph parallelism, parallelize existing sequential graph algorithms, and speed up beyond-neighborhood computations. Moreover, within each subgraph, it parallelizes the execution of PEval

---

**Algorithm 1:** A conceptual implementation of EMap

---

**Input:** active vertex set  $\Delta \in 2^V$ , update function  $f_E : (V, V) \rightarrow 2^V$ .

**Output:** updated active vertex set  $\Delta' \in 2^V$ .

```
1  $\Delta' := \emptyset$ ;
2 foreach  $u$  in  $\Delta$  do in parallel
3   foreach  $v$  in  $u$ .neighbors do in parallel
4      $\Delta' := \Delta' \cup f_E(u, v)$  /*  $F_i$  and  $\Pi_i$  may be accessed. */
5 return  $\Delta'$ ;
```

---

---

**Algorithm 2:** A conceptual implementation of VMap

---

**Input:** active vertex set  $\Delta \in 2^V$ , update function  $f_V : V \rightarrow 2^V$ .

**Output:** updated active vertex set  $\Delta' \in 2^V$ .

```
1  $\Delta' := \emptyset$ ;
2 foreach  $u$  in  $\Delta$  do in parallel
3    $\Delta' := \Delta' \cup f_V(u)$  /*  $F_i$  and  $\Pi_i$  may be accessed. */
4 return  $\Delta'$ ;
```

---

and IncEval by enforcing VC; that is, once a subgraph is in memory, it leverages the shared-memory architecture of a single machine and promotes intra-subgraph vertex/edge-level parallelism via VC.

**Example 3:** Continuing with Example 2, consider the execution of WCC under GC. On subgraph A (see Figure 1), it initiates BFS from root vertex with an initial label 1. It recurses as two independent BFSes starting from both neighbors of vertex 1, i.e., vertices 2 and 4, which then proceed as BFSes from vertices 11 (subgraph C), 13 (subgraph B), and 5 (subgraph D), so on and so forth. If we execute each recursion with 2 parallel threads, we get a  $2\times$  speedup.  $\square$

**PIE+.** In light of this, we propose a new programming interface, called PIE+. Given a problem  $Q$ , MiniGraph also takes algorithms PEval, IncEval and Assemble for  $Q$ , and follows the same workflow of GC for inter-subgraph parallelism. Moreover, it extends PIE with two additional primitives, EMap and VMap, along the same lines as the VC interface of Ligra [47]. These primitives make PEval and IncEval parallel within each subgraph at the edge/vertex level.

EMap and VMap work on a set  $\Delta$  of “active” vertices, which are to be further processed. Intuitively, EMap and VMap implement flatMap of functional programming, to execute an update function at each active vertex  $v$  in  $\Delta$  and moreover, “flat map”  $v$  to a new set of vertices that will remain active or get activated after the operation.

As shown in Algorithm 1, EMap takes as input  $\Delta$  and an edge update function  $f_E$ . For each active vertex  $u$  in  $\Delta$ , EMap makes updates along each of  $u$ ’s outgoing edge by applying  $f_E$ , executed in parallel. Here  $f_E$  is a user-defined function. It processes an edge  $e = \langle u, v \rangle$  from  $u$ , makes updates based on the data associated with  $e$ ,  $u$  and  $v$ , and returns a new (possibly empty) set of vertices that remain active after the update. Moreover,  $f_E$  may access the local subgraph  $F_i$  and intermediate data structure  $\Pi_i$  of PEval and IncEval.

Similarly, VMap takes  $\Delta$  and a vertex update function  $f_V$  as input. As shown in Algorithm 2, it aims to parallelize operations that are centered at each active vertex in  $\Delta$ . Instead of working along edges,  $f_V$  operates on each vertex  $u \in \Delta$ , processes and updates  $u$  by accessing  $F_i$  and  $\Pi_i$ , and returns an updated active vertex set.

Message passing. Similar to GC, PIE+ approaches inter-subgraph

---

**Algorithm 3:** PEval for WCC under PIE+.

---

**Input:** subgraph  $F_i = (V_i, E_i, L_i)$ .

**Output:** set  $CC_i$  of connected components in  $F_i$ .

*Message Preamble:* /\* candidate set  $C_i$  includes all border nodes. \*/

```
 $\underline{VMap}(V_i, \text{function}(v) \{v.\text{root} := v; \text{return } \emptyset; \})$ ;
1  $CC_i := \emptyset; \Pi := V_i$ ; /*  $\Pi$  is the set of unvisited vertices. */
2 while  $\Pi$  not empty do
3   find root vertex  $v_r$ , such that  $L_i(v_r) = \min_{u \in \Pi} L_i(u)$ ;
4    $CC_i := CC_i \cup \{v_r\}$ ; remove  $v_r$  from  $\Pi$ ;
5    $\Delta := \{v_r\}$ ;
6   while  $\Delta$  not empty do  $\Delta := \text{EMap}(\Delta, \text{BFSRecur})$ ;
7 return  $CC_i$ ;
```

*Message Segment:*  $M_i := \{(v, L_i(v.\text{root})) \mid v \in C_i\}$ ;  $f_{\text{aggr}} := \min$ ;

```
8 Procedure  $\text{BFSRecur}(u, v)$  :
9   if  $v$  not in  $\Pi$  then return  $\emptyset$ ;
10  remove  $v$  from  $\Pi$ ;  $v.\text{root} := v_r$ ; return  $\{v\}$ ;
```

---

synchronization via message passing. At each subgraph  $F_i$ , it declares (a) *status variables*  $\bar{x}$  for its *border nodes*, i.e., vertices that have edges to other fragments (under edge-cut); and (b) an aggregate function  $f_{\text{aggr}}$ , e.g., min and max. Intuitively, the status variables are candidates to be updated by the incremental step of IncEval, and  $f_{\text{aggr}}$  resolves conflicts when multiple workers assign different values to the same variable  $v.\bar{x}$  for a border node  $v$  of  $F_i$ .

MiniGraph adopts a *push-pull* mechanism for workers to exchange messages through MessageStore. Right after worker  $W_i$  gets initialized, it checks MessageStore to *pull* a set  $M_i$  of messages, which consists of updated status of the border nodes in  $F_i$ . After concluding computation,  $W_i$  collects the latest values  $v.\bar{x}$  for each border node  $v$ , and *pushes* them into MessageStore for aggregation. It inserts  $\bar{x}$  into  $\text{MessageStore}[v]$  if the entry is missing, and applies  $f_{\text{aggr}}(\text{MessageStore}[v], \bar{x})$  to resolve the conflict.

**Example 4:** We present a PIE+ program for WCC. It is the same as the PIE program of [21] for WCC, except that it parallelizes BFSes in PEval and IncEval within each subgraph via EMap and VMap. We consider w.l.o.g. edge-cut partitions [10, 31], in which border vertices are those that have edges to another fragment.

(1) PEval. As shown in Algorithm 3, at worker  $W_i$ , PEval computes the connected components  $CC_i$  of its local subgraph  $F_i$ , the partial result at  $F_i$ . Here  $CC_i$  is the set of root vertices that identify connected components, referred to as *CRoots*. It declares (a) a link to its CRoot, and (b) min as  $f_{\text{aggr}}$ . PEval (a) initializes an empty  $CC_i$  and a set  $\Pi$  of unvisited vertices in  $F_i$  (Line 1); and (b) conducts iterative BFS until all vertices are visited, i.e.,  $\Pi$  is emptied (Line 2–6). The BFS finds its CRoot  $v_r$  at the vertex with the lowest label (Line 3), adds  $v_r$  to  $CC_i$  to identify a unique connected component (Line 4), and sets the label of each descendant of  $v_r$  to the label  $L_i(v_r)$  of the CRoot (Line 5–6). MessageStore packs the smallest label for each border node, accessible by workers in the next superstep.

This PEval differs from its counterpart given in [21] only in the following, marked with dotted underlines. VMap is applied to parallelize the initialization of status variables (Message Preamble). EMap parallelizes its recursion (Line 6) to conduct a round of BFS. It works on a set  $\Delta$  of active vertices, which serves as the recursion



---

**Algorithm 4:** IncEval for WCC under PIE+.

**Input:** subgraph  $F_i = (V_i, E_i, L_i)$ ,  $CC_i$ , incoming messages  $M_i$ .

**Output:** refined set  $CC_i$  of connected components in  $F_i$ .

```
1  $\Delta = \{v \mid v \in M_i\}; \text{VMap}(\Delta, \text{UpdateRoot});$ 
2 return refine( $CC_i$ ); /* merge roots with the same label in  $CC_i$ . */
Message Segment:  $M_i := \{(v, L_i(v.\text{root})) \mid v \in C_i\};$ 

3 Procedure UpdateRoot( $v$ ) :
4    $v_r := v.\text{root};$ 
5   while  $v_r$  not in  $CC_i$  do  $v_r := v_r.\text{root};$ 
6   update  $v.\text{root} := v_r; L_i(v_r) := \min(L_i(v_r), M_i[v]);$ 
7   return  $\emptyset$ ;
```

---

queue for BFS. Function BFSRecur is passed to EMap as the update function  $f_E$ . For each *unvisited* neighbor  $v$  of the working vertex  $u$  (Line 9), BFSRecur overwrites  $v.\text{root}$  with  $u$ 's CRoot, and adds  $v$  to the recursion queue by returning  $\{v\}$  (Line 10). During the process, EMap makes parallel visits to the neighboring vertices.

(2) IncEval. As shown in Algorithm 4, IncEval *incrementally* refines the partial result  $CC_i$  at  $F_i$ , based on messages (updates) from other subgraphs. IncEval (a) updates the label  $L_i(v_r)$  of each CRoot  $v_r$  as the smallest of its descendants in the same connected component (Line 1); and (b) refines  $CC_i$  by merging CRoots ( $u_r, v_r \in CC_i$ ) with the same label, and sets their root to the remaining CRoot (Line 2). The messages are aggregated via min just like in PEval.

This IncEval differs from its counterpart of [21] in its VMap-based implementation of local message propagation. Here function UpdateRoot is passed to VMap as its update function  $f_V$ . For each vertex  $v$  that receives messages, UpdateRoot works in parallel to (a) find the CRoot  $v_r$  of  $v$  (Line 4–5), and (b) update the label of  $v_r$  by incorporating the incoming messages (Line 6).

(3) When no messages are generated, Assemble collects CRoots from all  $W_i$ 's and counts the number of CRoots of distinct labels.  $\square$

*Partial results.* Under GC, PEval computes partial result  $\mathcal{R}_i$  on the local subgraph  $F_i$  while IncEval refines  $\mathcal{R}_i$  incrementally based on messages. MiniGraph persists and embeds  $\mathcal{R}_i$  as the metadata of  $F_i$ , which will be discharged and loaded together with  $F_i$ .

**Writing PIE+ programs.** Besides directly parallelizing an existing PIE program (e.g., Example 4), one may cast a VC program into PIE+. Given, e.g., a GAS program of PowerGraph [25], we may convert it into a PIE+ program as follows: (1) initialize MessageStore for caching incoming messages at all vertices; (2) recast Scatter with VMap to generate the outgoing message for each vertex; (3) recast Gather with EMap to perform the message passing via adjacent edges of each vertex, and caching/aggregating them in the corresponding entries of MessageStore; (4) recast Apply with VMap to update each vertex based on its own status and its aggregated messages in MessageStore; (5) embed the implementation above in PEval and IncEval; and (6) specify Assemble based on the output.

With PIE+, users may opt to use either VC or GC that fits their applications and graphs the best. In [4], we give PIE+ programs for single-source shortest path (SSSP), PageRank (PR), breadth-first search (BFS), random walk (RW), and graph simulation (Sim).

## 5 EXECUTION MODEL

In this section, we present the execution model of MiniGraph, with resource scheduling (Section 5.1) and optimizations (Section 5.2).

### 5.1 Resource Scheduling

Given a partitioned graph, MiniGraph schedules its work as follows.

**Workflow.** Taking subgraphs  $F_0, F_1, \dots, F_{n-1}$  of graph  $G$  as input, MiniGraph executes a PIE+ program  $\mathcal{A}$  in BSP supersteps (rounds). Each round iterates over all subgraphs: the first is a PEval round, followed by IncEval rounds, on each subgraph. A new IncEval round cannot start until the last round completes, and messages generated in the current round are pulled when the next (IncEval) round starts.

In a round, on-disk subgraphs are loaded one after another, while in-memory ones are processed by concurrent workers. To process subgraph  $F_i$ , its virtual worker  $W_i$  is allocated  $p_i$  physical threads from available ones in ThreadPool by Scheduler;  $W_i$  then pulls messages generated in the last round, conducts computations on  $F_i$ , and releases resources back to Scheduler before it is deactivated.

**The scheduling problem.** Scheduler has to make two decisions at runtime: (1) when to load and process a subgraph, and (2) how to allocate resources to maximize two-level parallelism. Making optimum decisions for these is, however, highly nontrivial.

*A cost model.* Given a PIE+ program  $\mathcal{A}$ , we train a cost function  $C_{\mathcal{A}}$  to estimate the execution time of worker  $W_i$  on subgraph  $F_i$ . We adapt the cost function  $C_{\mathcal{A}_{\text{PIE}}}$  of [19] for PIE program  $\mathcal{A}_{\text{PIE}}$ :

$$C_{\mathcal{A}_{\text{PIE}}}(F_i) = \sum_{v \in F_i} h_{\mathcal{A}_{\text{PIE}}}(\bar{x}_i(v)), \quad (1)$$

where for a vertex  $v$  in  $F_i$ ,  $h_{\mathcal{A}_{\text{PIE}}}$  is a learnable polynomial regression model defined over a vector  $\bar{x}_i(v)$  of *metric variables*. Here  $\bar{x}_i(v)$  takes into account the average in/out-degree of all vertices in  $G$ , and various structural information of  $v$ , e.g.,  $v$ 's in/out-degree in  $F_i$  and  $G$ , and the number of  $v$ 's mirror across all subgraphs of  $G$ .

Model  $C_{\mathcal{A}_{\text{PIE}}}$  has proven accurate in predicting the computational cost of a PIE program, because polynomials can approximate a continuous function defined on a closed interval [55]. Extending  $C_{\mathcal{A}_{\text{PIE}}}$ , we define  $C_{\mathcal{A}}$  for a PIE+ program  $\mathcal{A}$  as

$$C_{\mathcal{A}}(F_i, p_i) = \sum_{v \in F_i} [h_{\mathcal{A}}^{\text{seq}}(\bar{x}_i(v)) + \frac{h_{\mathcal{A}}^{\text{para}}(\bar{x}_i(v))}{\min\{p_i, [d_i]\}}], \quad (2)$$

where  $p_i$  denotes the number of allocated threads for worker  $W_i$ ,  $d_i$  is the average degree of vertices in  $F_i$ , and  $h_{\mathcal{A}}^{\text{seq}}$  and  $h_{\mathcal{A}}^{\text{para}}$  are learnable models like  $h_{\mathcal{A}_{\text{PIE}}}$ . Intuitively, the cost of  $\mathcal{A}$  is broken into  $h_{\mathcal{A}}^{\text{seq}}$  for unparallelizable computations and  $h_{\mathcal{A}}^{\text{para}}$  for parallelizable operations (those defined in EMap and VMap). When more cores  $p_i$  are allocated,  $h_{\mathcal{A}}^{\text{para}}$  may get a linear speedup subject to maximum  $d_i$ .

Model  $C_{\mathcal{A}}$  generalizes  $C_{\mathcal{A}_{\text{PIE}}}$  by extending the data-partitioned parallelism of GC and adapting to the unique two-level parallel execution model of MiniGraph. It inherits two key properties of  $C_{\mathcal{A}_{\text{PIE}}}$ , which have been verified empirically in [19]: (1)  $C_{\mathcal{A}}$  depends on individual algorithm  $\mathcal{A}$ , and (2) it remains quite accurate on graphs that are of the same type (e.g., Web graphs) as training graphs, since its regression model assesses the impact of metric variables, and the variables “characterize” topological features of the type of graphs.

We train the polynomial regression models  $h_{\mathcal{A}}^{\text{seq}}$  and  $h_{\mathcal{A}}^{\text{para}}$  with training samples collected from historic running logs. One challenge

is that  $C_{\mathcal{A}}$  may not be accurate in the lack of the training data. To cope with this, we develop a training procedure to handle cold starts. We follow incremental SGD [12], initializing the model with default weights and updating it as more running logs are collected.

(1) *Initialization.* Given  $\mathcal{A}$  without historic running logs, we set

$$C_{\mathcal{A}}(F_i, p_i) = \sum_{v \in F_i} \frac{d_i^+(v)}{\min\{p_i, [d_i]\}},$$

where  $d_i^+(v)$  denotes the out-degree of vertex  $v$  in subgraph  $F_i$ . We heuristically take this as pre-trained weights for incremental SGD. Intuitively, without prior knowledge about  $\mathcal{A}$ , we assume that the cost of  $\mathcal{A}$  on  $F_i$  is proportional to the number of edges in  $F_i$ , and that operations in  $\mathcal{A}$  are fully parallelized with EMap and VMap.

(2) *Sample collection.* We collect training data for each run of  $\mathcal{A}$  over subgraph  $F_i$ . A sample includes: (a) the metric variables  $x_i(v)$  of each active vertex  $v$  in  $F_i$ , which “characterize” their individual topological structures; (b) the thread allocation  $p_i$ ; and (c) its measured computational cost  $t_i$ . After rounds of  $\mathcal{A}$  execution, we batch the new samples and feed them into the model.

(3) *Training.* Using mean square error as the loss function, incremental training computes the gradients and updates the model accordingly. It iterates through the new training data only. The model is continuously trained and updated at backend at a small cost.

*Remark.* We have experimented with  $C_{\mathcal{A}}$  and two alternative cost models: (a) a polynomial regression model with fewer metric variables and (b) a multi-layer perceptron. We find that (1) once trained for algorithm  $\mathcal{A}$ ,  $C_{\mathcal{A}}$  works well on graphs of the same type, echoing the observation of [19]; and (2) compared to model (a) and (b),  $C_{\mathcal{A}}$  strikes a better balance among accuracy, generalizability, and training cost. Due to lack of space, we present the results in [4].

*Problem formulation.* With  $C_{\mathcal{A}}$ , we model the scheduling problem as an optimization problem. Given subgraphs  $F_0, F_1, \dots, F_{n-1}$  of  $G$  and an  $m$ -core machine with memory size  $\eta$ , it is to find an optimal schedule  $\mathcal{S} = (\bar{t}, \bar{p})$ , where  $t_i \in \bar{t}$  (resp.  $p_i \in \bar{p}$ ) denotes the start time (resp. thread allocation of  $W_i$ ) for  $i \in [0, n]$ . Its objective function is

$$\arg \min_{\mathcal{S}} \max_{i \in [0, n]} \{t_i + C_{\mathcal{A}}(F_i, p_i)\}. \quad (3)$$

It is to minimize the *makespan* of  $\mathcal{S}$ , i.e., the overall completion time of all workers. At any point, a *valid* schedule is subject to three constraints: (1) the consumed memory cannot exceed  $\eta$ ; (2) the total thread allocation cannot exceed  $m$ ; and (3) the start time  $t_i$  cannot be set before  $F_i$  is fully loaded into the memory, for all  $i$  in  $[0, n]$ . Note that the final constraint takes into account the I/O cost under the memory capacity  $\eta$ . It models the behavior of Loader (see Section 3), which loads fragments continuously as long as  $\eta$  is not exhausted.

Its decision problem, denoted by DSP, is to decide, given the input and a deadline  $B$ , whether there exists a valid schedule  $\mathcal{S}$  with makespan at most  $B$ . The problem is intractable; it subsumes the NP-complete problem of *malleable parallel task scheduling* [35].

**Theorem 1:** DSP is NP-complete.  $\square$

**Proof sketch:** DSP is in NP since one can guess a schedule and check in PTIME whether it is valid and its makespan is at most  $B$ . We show that it is NP-complete by reduction from the set partition problem, which is known NP-complete [23] (see [4] for a proof).  $\square$

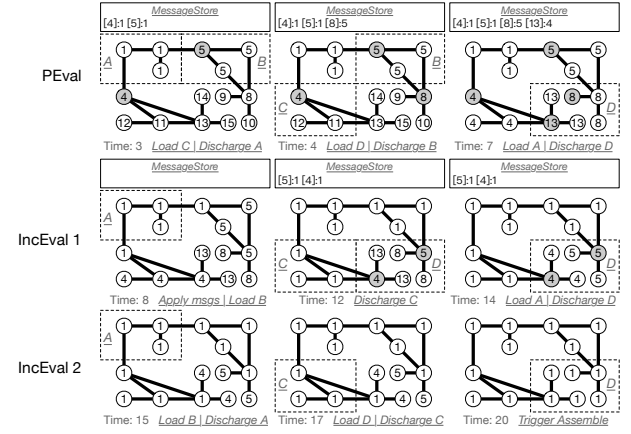


Figure 3: WCC on  $G$  with the two-level parallel execution model.

**Scheduling strategy.** Due to the intractability, Scheduler adopts a lightweight strategy. Offline for each subgraph  $F_i$ , it collects its size information and assigns  $F_i$  a number  $\hat{p}_i$ . At runtime, it ensures that each subgraph is allocated at least  $\hat{p}_i$  threads. Then whenever a set  $T$  of threads gets freed, we prioritize pending subgraphs  $F_i$  if  $|T| \geq \hat{p}_i$ . We set  $\hat{p}_i$  heuristically to estimate the optimal allocation for worker  $W_i$ , given the CPU and memory constraint. This strategy satisfies  $\hat{p}_i$  for  $W_i$  greedily and thus approximates the optimal schedule.

**Tentative resource allocation.** Scheduler allocates resources based on the subgraph size and the memory size  $\eta$ . It sets  $\hat{p}_i = \lceil \frac{s(F_i)}{\eta} m \rceil$  for worker  $W_i$ , where  $s(F_i)$  is the binary size of  $F_i$ . Intuitively,  $\hat{p}_i$  indicates the *minimum* number of threads for  $W_i$  to process  $F_i$ , to strike a balance between computing and memory resources. Following *multi-resource scheduling* [26], all threads are allocated *before* the memory gets exhausted, prioritizing CPU utilization.

**Greedy subgraph processing.** At runtime, Scheduler keeps track of a list of pending subgraphs, sorted by  $C_{\mathcal{A}}(F_i, \hat{p}_i)$ , the estimated computational cost given their tentative thread allocations. To select one for processing, it prioritizes the most computationally heavy one whose tentative allocation can be satisfied. When Scheduler is triggered by an event that  $t$  threads are freed by a worker, it is to load  $F_i$  with the highest  $C_{\mathcal{A}}(F_i, \hat{p}_i)$  as long as  $\hat{p}_i \leq t$ .

The intuition behind such ordering is to mitigate stragglers at the end of each superstep. Considering the sublinear speedup of intra-subgraph parallelism (Equation 2), prioritizing larger subgraphs can further promote inter-subgraph parallelism, since we have a better chance of overlapping them with more and shorter loads.

**Intra-subgraph parallelism.** When the unallocated threads cannot satisfy the tentative allocation of any pending subgraph, Scheduler allocates these free threads to active workers, one at a time, as it thrives to keep all CPU cores busy. It decides the receiving worker based on cost analysis. (1) With pending subgraphs, it estimates  $\Delta_i = C_{\mathcal{A}}(F_i, p_i) - C_{\mathcal{A}}(F_i, p_i + 1)$  for all active workers, and selects  $W_j$  with the highest  $\Delta_j$ , i.e.,  $W_j$  gets the most speedup from the extra thread. (2) With no subgraph pending processing for the round, it selects  $W_j$  with the maximum  $C_{\mathcal{A}}(F_j, p_j)$  to improve the straggler.

Such reallocations are temporary; whenever a new worker  $W_j$  becomes active, these available threads are retracted by Scheduler and reassigned to  $W_j$  to meet its demand of  $\hat{p}_j$  threads.

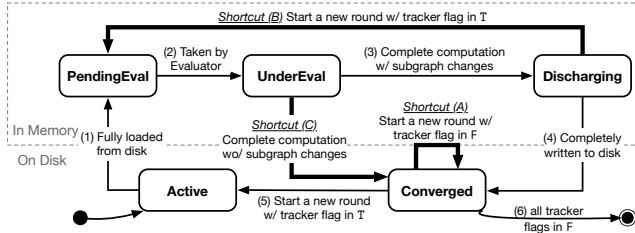


Figure 4: The state diagram of a worker.

**Example 5:** Continuing with Example 3, we compute WCC over  $G$  (Figure 1a) on a 4-core machine. Suppose that each subgraph takes 1 time unit to load/write, and each message passing takes 1 time unit. The tentative allocation assigns 2 threads to each worker. As shown in Figure 3, WCC passes all 32 message in 19 time units, leveraging both pipelined processing and two-level parallelism.  $\square$

## 5.2 System Optimizations

We now present the unique optimization strategies of MiniGraph.

Recall Example 5. In Figure 3, Subgraph  $A$  is not updated in both IncEval steps, and its worker receives no messages (updates); thus it requires no further processing, and there is no need to even load Subgraph  $A$ . To reduce the unnecessary I/O for loading/discharging Subgraph  $A$ , MiniGraph employs (a) a list  $\mathcal{M}$  of flags that helps track message exchanges among workers, and (b) a lightweight state machine for modeling the progress of each worker. Both  $\mathcal{M}$  and state machine are maintained by StateManager.

**Message tracking.** StateManager builds a list  $\mathcal{M}$  of flags, one for each worker to indicate whether the worker has any messages. It sets  $\mathcal{M}[i]$  true if worker  $W_i$  has at least one pending update to pull from MessageStore. Otherwise,  $W_i$  requires no incremental work to be done, and thus IncEval can be safely skipped in the next round.

**Example 6:** Continuing with Example 5, the step-by-step WCC execution is shown in Figure 3. Before starting Step IncEval 1, the flags are  $[F, T, T, T]$  for Subgraph  $A, B, C$  and  $D$ , respectively. Note that Subgraph  $A$  receives no messages in this round, and hence its flag is  $F$ . Similarly, the flags are  $[F, F, F, T]$  for Step IncEval 2.  $\square$

**Worker states and state transitions.** We use a finite state machine to model the progress at each worker  $W_i$ , and flag  $\mathcal{M}[i]$  to trigger state transitions of  $W_i$ . As shown in Figure 4, at any point, a worker is in one of the five states: Active, Converged, PendingEval, UnderEval, and Discharging. The first two indicate that the corresponding subgraph is on-disk, while the rest are in-memory states.

For a PEval or IncEval round, worker  $W_i$  traverses all five states. As shown in Figure 4, it starts from Active, indicating that  $W_i$  requires further computation (e.g., with unconsumed updates). Then, (1) after its subgraph data  $F_i$  is loaded into the memory by Loader,  $W_i$  becomes PendingEval, meaning that  $W_i$  is ready to run and requests thread allocation; (2)  $W_i$  turns to UnderEval, after it is taken by Evaluator and is under active computation; (3) when  $W_i$  concludes with changes made to  $F_i$ , its metadata (e.g., status variables, partial results) and messages aggregated in MessageStore,  $W_i$  is set to Discharging; and (4) once  $F_i$  is fully persisted on disk, Discharger releases its memory space and sets  $W_i$  to Converged, indicating that it is done for the round. MiniGraph concludes the round when

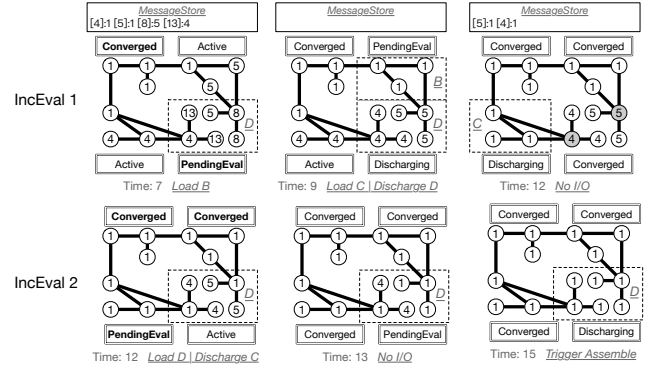


Figure 5: Optimized WCC execution on  $G$  with shortcuts.

all workers get past UnderEval. If there are messages cached in MessageStore, MiniGraph (5) resets Converged workers back to Active, and starts a new round of (IncEval) computation. Otherwise, MiniGraph (6) calls Assemble to produce the final results.

**Shortcuts in state transitions.** Under certain conditions, some states in a round of computation can be skipped without affecting the correctness. That is, MiniGraph can take some “shortcuts” in state transitions, and reduce unnecessary computation and I/O. Below are example shortcuts identified by StateManager, marked as bold arrows in Figure 4. Consider worker  $W_i$  with its flag  $\mathcal{M}[i]$ .

**Shortcut (A).** To start of a round of IncEval computation, the state of  $W_i$  is typically reset to Active following transition (5). If  $\mathcal{M}[i] = F$ , however, we may keep  $W_i$  in Converged state since it requires no further processing. That is, we can skip handling subgraph  $F_i$  in the round. This shortcut is frequently exploited when, e.g., the input graph is not well-connected and  $F_i$  is “isolated”.

**Shortcut (B).** MiniGraph starts a new round of IncEval as soon as all workers finish the current round (i.e., get past UnderEval). If  $W_i$  is still in Discharging, i.e., if its changed  $F_i$  is not yet fully persisted onto the disk,  $W_i$  is set to PendingEval directly, such that it starts the new round without going through the disk. This shortcut is exploited at the end of each round and substantially reduces I/O.

**Shortcut (C).** When  $W_i$  completes with no changes to  $F_i$  or its metadata,  $W_i$  skips Discharging and is set to Converged directly, avoiding redundant disk writes. This shortcut is effective for, e.g., SSSP.

**Example 7:** Continuing with Example 5–6, we show execution of WCC in Figure 5 by exploiting shortcuts in state transitions (in bold state labels). PEval execution remains the same as in Figure 3.

The start of Step IncEval 1 (Time 7) sees two shortcuts: (1) worker  $A$  takes Shortcut (A) to skip computation entirely; and (2) worker  $D$  skips discharging/loading by taking Shortcut (B); it enters PendingEval directly for the next round of IncEval. Because of the two shortcuts, Step IncEval 1 takes 5 time units to complete (Time 7–12), as opposed to 7 time units (Time 7–14) in Example 5.

Similarly, we optimize Step IncEval 2. It takes 3 time units (Time 12–15), a 40% reduction from 5 units (Time 14–19) in Example 5.  $\square$

## 6 EXPERIMENTAL EVALUATION

Using real-life graphs, we evaluated MiniGraph for its (1) efficiency, (2) scalability, (3) resource scheduling and (4) performance.



Table 2: Graph datasets.

| Name            | Type           | V     | E    | MaxDegree | Raw Data |
|-----------------|----------------|-------|------|-----------|----------|
| friendster [3]  | social network | 65.6M | 1.8B | 5124      | 30.14GB  |
| web-sk [45]     | Web            | 50M   | 1.9B | 8.5M      | 32GB     |
| datagen [7, 30] | synthetic      | 29M   | 2.6B | 2288      | 75GB     |
| clueWeb [45]    | Web            | 1.7B  | 7.9B | 6.4M      | 137GB    |
| hyperlink12 [6] | Web            | 7B    | 7.6B | 5.1M      | 155GB    |

**Experimental setup.** We start with the settings.

**Datasets.** We used five real-life and synthesized datasets (Table 2). Among these, datagen, clueWeb and hyperlink12 exceed the memory capacity of our testbed. Here datagen is a synthesized dataset in the LDBC [30] benchmark suite; and hyperlink12 is sampled from hyperlink [6] via BFS sampling [29]. We used edge-cut [34] as the default partitioner, and tested other partitioners.

**Baselines.** We evaluated three out-of-core systems as baselines: GridGraph [64], GraphChi [33] and XStream [46]. GridGraph is the state-of-the-art, assuming an off-the-shelf hardware platform. We did not test Mosaic [37] since it requires an Intel Xeon Phi coprocessor, which has been discontinued since 2018 [41]. We also omitted in-memory systems (e.g., [36, 43, 47, 60, 62]) since MiniGraph targets applications when the input graph cannot fit into the memory.

We tested four variants of MiniGraph: (1) MiniGraph<sub>Seq</sub>, which exposes PIE as interface without EMap and VMap (Section 4). (2) MiniGraph<sub>NoShort</sub>, which disables shortcuts (Section 5.2). (3) MiniGraph<sub>VCut</sub> and MiniGraph<sub>HCut</sub>, which use vertex-cut [2] and hybrid-cut [15], respectively, as graph partitioner.

We also compared with GraphScope [18] and Gluon [16], which use multiple machines, not a single machine as MiniGraph.

**Algorithms and queries.** We ran PIE+ programs for WCC, SSSP, PR, BFS, RW and Sim (Section 4) [4]. The first five are included in various benchmarks [5, 11]. For baselines, we used their out-of-box implementations. No baseline supports Sim out-of-box. We implemented the VC-based Sim algorithm of [22] on GraphChi. However, the algorithm cannot be implemented on GridGraph or XStream without changing their internal working, because a necessary global storage is not supported by the APIs of either system.

For each input graph for SSSP and BFS, we randomly picked 10 vertices and used them as sources. For RW, we initiated 100 walkers randomly, each taking at most 10 steps. For Sim, we randomly generated 20 query patterns, each of which has 6 vertices and 10 edges.

We deployed the single-machine systems on a workstation running Ubuntu Server 20.04 LTS, powered with an Intel Core i9-7900X CPU @3.30GHz and 64GB of DDR4-2666 memory. The CPU has a 13.75MB LLC, shared among 10 cores (20 hyperthreads).

Unless stated otherwise, we set  $n$  partitions as reported in Table 4 such that we have  $\hat{p}_i = 4$  by default (see Section 5). We set a different  $n$  for PR (s.t.  $\hat{p}_i = 10$ ), guided by our experiments with varying  $n$  (see [4] for details). Graphs are loaded from a 1TB WD Blue WDS100T2B0A SATA SSD, whose average sequential read throughput is 560MB/s. Each experiment was repeated 5 times; we report the average here. We report results for some algorithms on some graphs; the other results are consistent and are reported in [4].

**Experimental results.** We next report our findings.

**Exp-1: Efficiency.** We first evaluated the efficiency and I/O of MiniGraph versus out-of-core systems. Some experiments are con-

Table 3: Runtime statistics for SSSP, WCC and PR.

| Dataset    | Metric            | SSSP      |           | WCC       |           | PR        |           |
|------------|-------------------|-----------|-----------|-----------|-----------|-----------|-----------|
|            |                   | MiniGraph | GridGraph | MiniGraph | GridGraph | MiniGraph | GridGraph |
| friendster | # Supersteps      | 8         | 32        | 6         | 21        | 8         | 10        |
|            | Disk Read (GB)    | 78        | 115.1     | 50        | 135       | 107       | 160       |
|            | Shortcut I/O (GB) | -12       | N/A       | -22       | N/A       | -10.4     | N/A       |
|            | Avg. CPU Util.    | 33.74%    | 4.45%     | 48.2%     | 6.83%     | 68.46%    | 62.38%    |
|            | I/O-CPU Corr.     | 0.095     | -0.113    | 0.163     | -0.202    | 0.185     | -0.156    |
|            | Cache Hits        | 45.33%    | 9.59%     | 48.25%    | 12.04%    | 34.8%     | 36.2%     |
| web-sk     | # Supersteps      | 10        | 63        | 9         | 120       | 15        | 20        |
|            | Disk Read (GB)    | 112.5     | 232       | 42        | 367       | 87        | 232       |
|            | Shortcut I/O (GB) | -30.9     | N/A       | -24.6     | N/A       | -133.9    | N/A       |
|            | Avg. CPU Util.    | 15.76%    | 5.83%     | 25.04%    | 5.16%     | 42%       | 42%       |
|            | I/O-CPU Corr.     | 0.008     | 0.003     | 0.013     | 0.009     | 0.082     | -0.039    |
|            | Cache Hits        | 50.89%    | 6.37%     | 37.42%    | 11.63%    | 50.22%    | 46.04%    |

ducted under a *memory budget*, i.e., the ratio of the available memory to the size of the raw input. For instance, a 50% memory budget means that exactly *half* of the graph can be held in memory. We used Linux cgroups to enforce the memory budget, similar to [9].

**SSSP.** For SSSP over different graphs, Table 3 reports various runtime statistics compared with GridGraph, the best performing baseline. Table 4 shows the performance of all the systems.

(1) MiniGraph performs the best over all workloads. Over two real-life large clueWeb and hyperlink12, MiniGraph performs substantially better than all the baselines. With the memory capacity being only 47% (resp. 41%) of clueWeb (resp. hyperlink12), it outperforms GridGraph, by 4.6 $\times$  (resp. 3.2 $\times$ ). Moreover, GraphChi and XStream are not able to handle the two large graphs (Table 4). Over the large synthetic datagen, the improvement is relatively moderate (1.5 $\times$ ), since the graph is extremely sparse and the VC of GridGraph does not incur significant I/O. The results are consistent over small graphs, which justifies memory budgeting, a common practice for benchmarking out-of-core systems (e.g., [9]).

(2) With a 50% memory budget on friendster, MiniGraph beats GridGraph, GraphChi and XStream by 1.5 $\times$ , 2.7 $\times$  and 15.2 $\times$ , respectively. It is 2.8 $\times$ , 3.5 $\times$  and 28.9 $\times$  faster on web-sk, a better improvement, because web-sk is sparser and has a larger diameter. As remarked earlier, the VC of the baselines supports within-neighborhood operations only, and requires more supersteps to converge on graphs with larger diameters. GridGraph, e.g., takes 32 supersteps on friendster and 63 on web-sk. In contrast, MiniGraph supports beyond-neighborhood computations via its hybrid parallel model. It takes 8 supersteps on friendster and 10 on web-sk. This justifies the need for GC and two-level parallelism.

Note that over friendster, GridGraph takes 4 $\times$  the supersteps of MiniGraph, while its execution time is 2.7 $\times$ . The main reason lies in their different parallel models; a GC superstep is quite different from a VC superstep as it processes a subgraph. MiniGraph generally takes more computationally heavy supersteps than GridGraph does. This said, its superstep progresses the computation further with each superstep, and reduces the overall work required.

(3) MiniGraph also substantially reduces disk I/O (Table 3). It generates 32.3% and 51.6% less disk read than GridGraph on friendster and web-sk, respectively. More specifically, (a) shortcuts account for 32.3% and 25.8% of the I/O reduction on the two graphs, respectively. (b) The rest of I/O reduction roots in the reduced supersteps.

(4) The average CPU utilization of MiniGraph is 29.3% (resp. 9.9%) higher than GridGraph on friendster (resp. web-sk). GridGraph

Table 4: Execution time for SSSP, WCC and PR (in seconds). “/” denotes that the experiment could not finish within 4 hours.

| Data        | Memory Budget | #Partitions (PR/Others) | SSSP         |              |              |              | WCC          |              |               |              | PR             |                |              |                |
|-------------|---------------|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|---------------|--------------|----------------|----------------|--------------|----------------|
|             |               |                         | MiniGraph    | GraphChi     | GridGraph    | XStream      | MiniGraph    | GraphChi     | GridGraph     | XStream      | MiniGraph      | GraphChi       | GridGraph    | XStream        |
| friendster  | 50% (15.07GB) | 4/10                    | <b>201.8</b> | 535 (2.7×)   | 293.1 (1.5×) | 3061 (15.2×) | <b>171.8</b> | 1636 (9.5×)  | 204.7 (1.2×)  | 2037 (11.8×) | <b>190.104</b> | 450.7 (1.9×)   | 485.3 (1.9×) | 2685 (11.3×)   |
| web-sk      | 50% (16GB)    | 4/10                    | <b>326.4</b> | 1140 (3.5×)  | 917.9 (2.8×) | 9437 (28.9×) | <b>172</b>   | 620.1 (3.6×) | 704.6 (4.1×)  | 4056 (23.5×) | <b>248.3</b>   | 2288 (9.2×)    | 395 (1.6×)   | 2903 (11.7×)   |
| datagen     | 85% (64GB)    | 4/4                     | <b>85</b>    | 2483.4 (29×) | 127.1 (1.5×) | 2350 (27.6×) | <b>80.4</b>  | 2688.6 (33×) | 182.35 (2.3×) | 6124 (76.1×) | <b>306.4</b>   | 1104.03 (3.6×) | 1066 (3.5×)  | 3434.7 (11.2×) |
| c1ueWeb     | 47% (64GB)    | 4/10                    | <b>2514</b>  | /            | 11534 (4.6×) | /            | <b>2742</b>  | /            | 11665 (4.3×)  | /            | <b>2022</b>    | /              | 3803 (2.1×)  | /              |
| hyperlink12 | 41% (64GB)    | 4/4                     | <b>150.1</b> | /            | 475.6 (3.2×) | /            | <b>1940</b>  | /            | 12570 (6.5×)  | /            | <b>1553.3</b>  | /              | 5209 (3.4×)  | /              |

exhibits negative correlations ( $r < -0.1$ ) between I/O and CPU usage, indicating that CPU may starve when the system is busy at loading data. In contrast, it is  $0 < r < 0.1$  for MiniGraph; that is, its I/O does not block computations. These justify the pipelined architecture.

(5) MiniGraph has a much better CPU cache locality than GridGraph. It achieves a 45% cache hit rate on friendster and 51% on web-sk, while it is 10% for GridGraph at best (Table 3). This is because MiniGraph supports inter-subgraph parallelism, while VC systems (e.g., GridGraph) frequently access all neighbors of each vertex, which inevitably incurs random reads across the entire graph.

(6) MiniGraph also performs substantially better than its variants. We defer a detailed discussion to the case of WCC below.

**WCC.** As shown in Tables 3–4, MiniGraph beats all the baselines for WCC over all graphs. (1) It is up to 6.5×, 33.0× and 76.1× faster than GridGraph, GraphChi and XStream, respectively; note that for GraphChi and XStream, 33.0× and 76.1× do not aggregate results on c1ueWeb and hyperlink12, since the two systems could not handle the two large graphs (Table 4). (2) Over web-sk, it takes only 7.5% supersteps and 28.9% of disk read of GridGraph; moreover, it has better average CPU utilization (+41.4%) and cache locality (+36.2%). These lead to 4.1× performance improvement. (3) MiniGraph incurs 62.9% (resp. 88%) less disk read on friendster (resp. web-sk), with 16.3% (resp. 11.4%) reductions from shortcuts.

MiniGraph consistently beats its variants. As shown in Figure 6a, (a) It speeds up MiniGraph<sub>Seq</sub> by 61.0% and 78.1% on friendster and web-sk, respectively. This shows the benefit of its hybrid model that exploits intra-subgraph parallelism to improve multi-core parallelism. The hybrid model works better on denser graphs, which often allow a higher inherent intra-subgraph parallelism, to which EMap and VMap are more effective. (b) MiniGraph beats MiniGraph<sub>NoShort</sub> by 1.18× on average. The improvement is moderate, since all shortcuts apply to IncEval only, while PEval is the most costly for WCC. They reduce I/O by 44% on average, yet the effect gets diluted by other factors, e.g., stragglers. (c) Shortcuts are less effective over well-connected graphs (e.g., a 12% improvement over friendster), but are better over more “isolated” web-sk (36%).

Figure 6b shows the disk I/O and CPU utilization of MiniGraph and GridGraph over c1ueWeb. It verifies that MiniGraph overlaps I/O and CPU operations, and justifies its pipelined architecture.

**PR.** As reported in Table 4, (1) MiniGraph is 1.6–3.5× faster than GridGraph over all graphs, the best-performing baseline, although PR fits its VC model. (2) On friendster (resp. web-sk), MiniGraph incurs 33.1% (resp. 45.6%) less disk read traffic than GridGraph, with 19.6% (resp. 28.6%) reductions from shortcuts. These are higher than SSSP and WCC, since we use larger fragments in PR.

**Sim.** As shown in Figure 6c for Sim on friendster and web-sk, on average MiniGraph beats GraphChi by 2.2× and 1.6×, and reduces I/O traffic by 81% and 77.4%, respectively. One reason is that its hybrid model can support a more efficient Sim algorithm (see Appen-

dix A.1 of [4]). Its complexity is  $O((|E| + |V|)(|E_Q| + |V_Q|))$ , where  $|E|$  and  $|V|$  (resp.  $|E_Q|$  and  $|V_Q|$ ) denote the size of edges and vertices of  $G$  (resp. pattern  $Q$ ), respectively. In contrast, GraphChi can only use the Sim algorithm [22] under VC and take  $O(|E|^2(|V_Q| + |E_Q|))$  time. This further justifies the need for supporting GC.

**BFS and RW.** Figure 6d reports the speedup of MiniGraph for BFS and RW versus GridGraph. Consistent with the other algorithms, MiniGraph beats all the competitors. For BFS, MiniGraph is 1.2–1.7× faster than GridGraph over various graphs. It is at least 1.5× better over more skewed Web graphs (i.e., web-sk and c1ueWeb). Although the GC and VC algorithms for RW have similar complexity bounds, the improvement is still quite evident (52.5% on average).

**HDD vs. SSD.** We also tested the impact of using HDD as the secondary storage. When replacing SSD with a 3.6TB Seagate HDD, the average read throughput drops from 560MB/s to 125MB/s, and all the systems get slower, as expected. On friendster, e.g., it is slowed by 2.6×, while it is 3.3× and 2.9× for GridGraph and GraphChi, respectively. This is because (a) MiniGraph incurs less I/O traffic (see Table 3), and (b) it issues a small number of bulk I/O requests, and works well with HDD that has a high seek cost.

**Exp-2: Scalability.** Under 50% memory budget, we evaluated the scalability of MiniGraph with the size  $|G|$  of graphs  $G$  and the number  $m$  of CPU cores. We report the results of WCC and PR; the results of the other algorithms are consistent.

**Varying  $|G|$ .** We sampled graphs  $G$  from large c1ueWeb using Edge Sampling [29], with a scale factor  $\delta$  that controls the fraction of edges to be sampled. As shown in Figure 6f when varying  $\delta$  from 0.4 to 1.0 for WCC, (1) MiniGraph scales well with  $|G|$ . It takes 2.3× longer, while it is 2.6× for GridGraph. (2) MiniGraph also scales better than its variants MiniGraph<sub>Seq</sub> (3.5×) and MiniGraph<sub>NoShort</sub> (3.16×). These further verify the effectiveness of our hybrid parallel model and optimization strategies. (3) When  $|G|$  grows 2.5×, MiniGraph incurs 2.29× I/O traffic, by its shortcut optimizations.

**Varying  $m$ .** Varying the number  $m$  of cores from 4 to 20, we ran WCC and PR on web-sk. As shown in Figure 6g, (1) MiniGraph scales well with  $m$ . (2) It scales better than GridGraph, which does not improve much (up to +4.3%) when  $m$  varies from 8 to 20. This is because GridGraph becomes I/O-bound when  $m \geq 8$  (see also our findings in Figure 6b). (3) It also scales better than MiniGraph<sub>Seq</sub>, which barely improves when  $m \geq 4$  for the lack of intra-subgraph parallelism. These also justify the need for a hybrid parallel model.

**Exp-3: Resource scheduling.** We evaluated the scheduling strategy by testing MiniGraph in different settings of (1) thread allocation  $p_i$ , (2) the number  $n$  of partitions, and (3) graph partitioners.

**Varying  $p_i$ .** We first tested the impact of  $p_i$ , the number of threads allocated to processing a subgraph. Over c1ueWeb, Figure 6h reports how WCC and PR works with varying  $p_i$ . For each graph partitioning, we mark its tentative allocation  $\hat{p}_i$  (Section 5) with colored

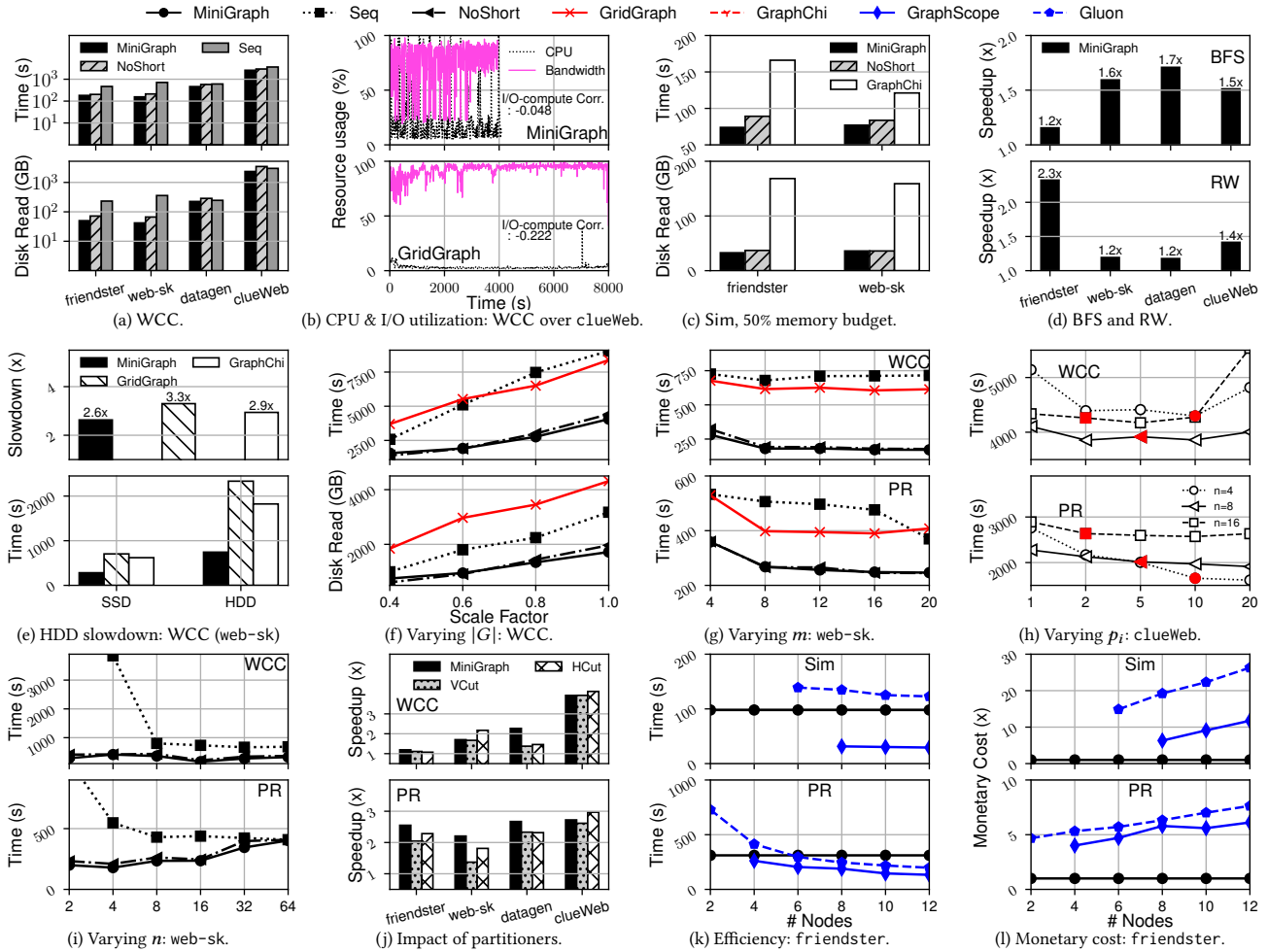


Figure 6: Efficiency, scalability, and resource scheduling of MiniGraph, and its performance vs. multi-machine systems.

points. As shown there,  $\hat{p}_i$  strikes a balance between the two levels of parallelism and improves the overall multi-core parallelism, by selecting a near-optimal value for  $p_i$  based on  $n$  and  $G$ .

**Varying  $n$ .** The number  $n$  of partitions impacts the inter-subgraph parallelism of MiniGraph. Varying  $n$  from 2 to 64 on web-sk, Figure 6i shows its impact on WCC and PR. (1) For VC-based PR, its optimal  $n$  is 4, which minimizes fragmentation while keeping the pipeline working. (2) MiniGraph<sub>Seq</sub> is very slow under a small  $n$ , as its CPU usage is limited by the low inter-subgraph parallelism. (3) For GC-based WCC, the optimal  $n$  is 16. Finding this optimum value requires a careful cost analysis (Section 8).

**Impact of partitioners.** We tested WCC and PR using different partitioners with  $n = 4$ . As reported in Figure 6j, MiniGraph beats all the baselines regardless of the partitioner used. Running PR on large clueWeb, for example, MiniGraph<sub>VCut</sub> and MiniGraph<sub>HCut</sub> outperform GridGraph by 2.71 $\times$  and 2.98 $\times$ , respectively.

**Exp-4: Single-machine vs. multi-machine.** We also evaluated the capacity of MiniGraph for graph analytics with a single machine versus multi-machine systems GraphScope and Gluon. To make a uniform testing environment, we deployed them in the cloud.

**Execution time.** Figure 6k reports the performance of the systems for Sim and PR over friendster. (1) For Sim, we used 8-vCPU 64GB-memory instances for all three systems since GraphScope and Gluon require more working memory. Gluon (resp. GraphScope) runs out-of-memory with fewer than 6 (resp. 8) nodes. MiniGraph, using a single instance, outperforms Gluon with up to 12 nodes by 29.3–46.3%. This further justifies the need for GC. While GraphScope takes less computation time, it requires an additional 200+s for preprocessing, which is not counted in Figure 6k. (2) For PR, MiniGraph runs on a single 8-vCPU 32GB-memory instance. Gluon uses a varying number of instances of the same configuration. GraphScope uses multiple 8-vCPU 64GB-memory instances, because it requires more than 32 GB in all cases of the experiment. Yet it still runs out-of-memory with 2 nodes. We find that MiniGraph performs comparably to a 4-node (resp. 6-node) deployment of GraphScope (resp. Gluon), and beats Gluon with fewer than 4 nodes.

**Cost effectiveness.** Figure 6l depicts the relative monetary cost of multi-machine systems for running Sim and PR over friendster, taking the cloud spending of MiniGraph as baseline. (1) While Gluon and GraphScope get faster with more nodes, they spend more. This echoes the observation of [39] that multi-machine par-

allelism is not always cost-effective. (2) GraphScope (resp. Gluon) pays at least 5.3 $\times$  (resp. 13.9 $\times$ ) more than MiniGraph for Sim, and costs 3.0 $\times$  (resp. 3.7 $\times$ ) more for PR. This further justifies the need for a single-machine system to save cost for small companies.

**Summary.** We find the following. (1) With the pipelined architecture and two-level parallelism, MiniGraph consistently outperforms the prior single-machine systems under all out-of-core workloads. It is up to 6.5 $\times$ , 33.0 $\times$  and 76.1 $\times$  faster than GridGraph, GraphChi and XStream, respectively; moreover, it can handle large graphs that exceed the capacity of GraphChi and XStream. (2) Under BSP, it requires only a fraction of supersteps (<29%) and disk read traffic (<53.3%) of GridGraph for SSSP and WCC. (3) It improves the CPU utilization of GridGraph, the best-performing baseline, by up to 41.4%. (4) It scales well with graphs and the number of CPU cores. It runs all the algorithms on hyperLink12 of 155GB within 32 minutes, while no baseline finishes in 1.4 hours. (5) Its shortcut optimization effectively reduces I/O cost, especially on dense graphs. (6) Its two-level execution model balances inter-subgraph and intra-subgraph parallelism. (7) MiniGraph works better than Gluon with 12 machines on Sim, and saves the monetary cost of multi-machine systems from 3.0 $\times$  to 13.9 $\times$ .

## 7 RELATED WORK

We categorize the related work as follows.

*Multi-machine systems.* A number of multi-machine systems have been developed, to support big graph analytics by scaling out, e.g., [8, 16, 20, 25, 27, 38, 42, 51, 54, 58, 63, 63]. Such systems adopt a shared-nothing architecture: they partition the input graph, and distribute the fragments to workers; all workers process their local fragments in-memory in parallel, and communicate with each other via message passing. Most of these systems adopt the VC model, except that GRAPE [20] proposes and supports GC. GraphScope [1, 18] extends GRAPE and can operate with either GC or VC.

In contrast to scaling out with multiple machines, (1) MiniGraph uses a single machine and explores multi-core parallelism. It aims to provide small businesses with a capacity of big graph analytics under limited resources. (2) It adopts a shared-memory architecture and employs the secondary storage as the memory extension. These introduce new challenges, i.e., the memory constraint and I/O cost. (3) It develops (a) a pipelined architecture to mitigate I/O cost and improve CPU utilization, and (b) an efficient synchronization scheme among workers by using a shared data structure. (4) Extending GC from multi-machine to multi-core, MiniGraph proposes a hybrid parallel model. It integrates VC and GC in a unified interface to enrich inter-subgraph GC parallelism with intra-subgraph VC parallelism, as opposed to GraphScope, which supports VC and GC as two separate sets of programming interfaces. (5) It develops a two-level execution model to support its intra- and inter-subgraph parallelism. Under the memory constraint of a single machine, this better utilizes the multi-core resources. Moreover, it introduces an efficient method for a new (intractable) resource scheduling problem. (6) MiniGraph explores new optimization methods to skip computation rounds, and hence reduce unnecessary CPU and I/O.

*Single-machine in-memory systems.* When graphs are small enough to fit into the memory of a single machine, in-memory systems [36,

43, 47, 60, 62] aim to conduct computations with high performance. Such systems focus on optimization techniques e.g., efficient task parallelization [43, 47, 62] and improved data locality [60].

In contrast, MiniGraph targets large graphs that exceed the memory capacity of a single machine. It performs iterative out-of-core computation, by actively swapping graph data between the memory and the disk. It deals with inevitable and possibly prohibitive I/O costs, a challenge that is not encountered by in-memory systems.

*Single-machine out-of-core systems.* Closer to this work are out-of-core systems [9, 24, 33, 37, 46, 53, 57, 64]. Most of the systems adopt VC and BSP. The techniques have mostly focused on reducing I/O, the major performance bottleneck. XStream [46] proposes stream reading of edges to maximize sequential disk accesses. GraphChi [33] divides the input data into small shards, each having a disjoint set of independent edges; it employs Parallel Sliding Windows to schedule shard I/O. Vora *et al.* [53] propose online adjustment of sharding, which avoids loading unneeded edges at the price of extra computation. GridGraph [64] improves the locality of edge sharding via 2-level partitioning based on source and destination IDs, and enables selective scheduling of data. Clip [9] allows asynchronous and repetitive processing of an in-memory shard; however, it must be explicitly invoked by users in their VC code, which might be error-prone. With extra hardware, Gill *et al.* [24] use persistent memory sticks as memory extension, which have higher throughput than SSDs (Intel has discontinued the Optane-only SSDs for the consumer market [48]). Mosaic [37] employs an Intel Xeon Phi coprocessor to accelerate CPU computations.

In contrast, MiniGraph (1) proposes a subgraph-based pipelined architecture to “cancel” I/O costs and promote sequential accesses to disks, to exploit the I/O bandwidth; (2) it supports both GC to simplify parallel programming and speed up beyond-neighborhood computation, and VC to parallelize edge/vertex operations in each subgraph; (3) it proposes a two-level execution model to balance inter-subgraph and intra-subgraph parallelism; and (4) it develops new optimization strategies to reduce I/O. (5) As opposed to [24, 37], it requires no dedicated hardware for storage or computation; it assumes an off-the-shelf machine that can be easily acquired.

## 8 CONCLUSION

The novelty of MiniGraph consists of (1) the first single-machine system that (a) extends GC from multiple machines to multiple cores, and (b) enriches inter-subgraph parallelism (GC) with intra-subgraph parallelism (VC) under a hybrid parallel model; (2) a pipelined architecture to overlap I/O and CPU operations; (3) a two-level execution model to schedule and balance VC and GC; (4) a unique optimization scheme to further reduce I/O. Our experimental study has verified that with a single machine, MiniGraph performs better than some 12-machine parallel systems.

One topic for future work is to extend the architecture of MiniGraph by incorporating GPU and FPGA. Another topic is to study the optimal graph partitions for a given algorithm.

## ACKNOWLEDGMENTS

We thank Jingbo Xu, Weijie Ou, and Zheng He for helpful discussions. This work was supported in part by Royal Society Wolfson Research Merit Award WRM/R1/180014 and NSFC 62225202.

## REFERENCES

- [1] 2020. *GraphScope*. Retrieved May 1, 2023 from <https://graphsco.pe.io>
- [2] 2022. *Apache Spark Partition Strategy*. Retrieved May 1, 2023 from [https://spark.apache.org/docs/1.4.0/api/java/org/apache/spark/graphx/PartitionStrategy.RandomVertexCut\\$.html](https://spark.apache.org/docs/1.4.0/api/java/org/apache/spark/graphx/PartitionStrategy.RandomVertexCut$.html)
- [3] 2022. *Friendster dataset*. Retrieved May 1, 2023 from <https://snap.stanford.edu/data/com-Friendster.html>
- [4] 2023. *Full version*. Retrieved May 1, 2023 from <https://shuao.liu.github.io/assets/papers/minigraph-full.pdf>
- [5] 2023. *Graph500 specifications*. Retrieved May 1, 2023 from <https://graph500.org>
- [6] 2023. *Hyperlink*. Retrieved May 1, 2023 from <http://webdatacommons.org/hyperlinkgraph/>
- [7] 2023. *LDBC dataset*. Retrieved May 1, 2023 from <https://graphalytics.org/datasets>
- [8] 2023. *Timely Dataflow*. Retrieved May 1, 2023 from <https://github.com/frankmcsherry/timely-dataflow>
- [9] 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*. 125–137.
- [10] Konstantin Andreev and Harald Racke. 2006. Balanced graph partitioning. *Theory of Computing Systems* 39, 6 (2006), 929–939.
- [11] Scott Beamer. 2016. *Understanding and improving graph algorithm performance*. Ph.D. Dissertation. EECS Department, University of California, Berkeley.
- [12] Dimitri P Bertsekas et al. 2011. Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. *Optimization for Machine Learning* 2010, 1–38 (2011), 3.
- [13] Florian Bourse, Marc Lelarge, and Milan Vojnovic. 2014. Balanced graph edge partition. In *SIGKDD*. 1456–1465.
- [14] Sergey Brin and Lawrence Page. 2012. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks* 56, 18 (2012), 3825–3833.
- [15] Rong Chen, Jiaxin Shi, Yanze Chen, Binyu Zang, Haibing Guan, and Haibo Chen. 2019. Powerlyra: Differentiated graph computation and partitioning on skewed graphs. *TOPC* 5, 3 (2019), 1–39.
- [16] 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*. 752–768.
- [17] Wenfei Fan. 2022. Big Graphs: Challenges and Opportunities. *PVLDB* 15, 12 (2022), 3782–3797.
- [18] 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.
- [19] 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*. 1765–1779.
- [20] Wenfei Fan, Jingbo Xu, Yinghui Wu, Wenyuan Yu, Jiaxin Jiang, Bohan Zhang, Zeyu Zheng, Yang Cao, and Chao Tian. 2017. Parallelizing Sequential Graph Computations. In *SIGMOD*. 495–510.
- [21] 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, 4, Article 18 (2018), 39 pages.
- [22] Arash Fard, M. Usman Nisar, Lakshmi Ramaswamy, John A. Miller, and Matthew Saltz. 2013. A distributed vertex-centric approach for pattern matching in massive graphs. In *IEEE BigData*. 403–411.
- [23] Michael Garey and David Johnson. 1979. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Company.
- [24] 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.
- [25] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. 2012. PowerGraph: Distributed Graph-Parallel Computation on Natural Graphs. In *OSDI*. 17–30.
- [26] Robert Grandl, Mosharaf Chowdhury, Aditya Akella, and Ganesh Ananthanarayanan. 2016. Altruistic Scheduling in Multi-Resource Clusters. In *OSDI*. 65–80.
- [27] Minyang Han and Khuzaima Daudjee. 2015. Giraph Unchained: Barrierless Asynchronous Parallel Execution in Pregel-like Graph Processing Systems. *PVLDB* 8, 9 (2015), 950–961.
- [28] John Hopcroft and Robert Tarjan. 1973. Algorithm 447: efficient algorithms for graph manipulation. *Commun. ACM* 16, 6 (1973), 372–378.
- [29] Pili Hu and Wing Cheong Lau. 2013. A survey and taxonomy of graph sampling. *arXiv preprint arXiv:1308.5865* (2013).
- [30] Alexandru Iosup, Tim Hegeman, Wing Lung Ngai, Stijn Heldens, Arnau Prat-Pérez, Thomas Manhardt, Hassan Chafio, Mihai Capotă, Narayanan Sundaram, Michael Anderson, et al. 2016. LDBC Graphalytics: A benchmark for large-scale graph analysis on parallel and distributed platforms. *PVLDB* 9, 13 (2016), 1317–1328.
- [31] George Karypis and Vipin Kumar. 1998. Multilevel-Way Partitioning Scheme for Irregular Graphs. *JPDG* 48, 1 (1998), 96–129.
- [32] 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.
- [33] Aapo Kyrola, Guy Blelloch, and Carlos Guestrin. 2012. GraphChi: Large-scale graph computation on just a PC. In *OSDI*. 31–46.
- [34] Yifan Li. 2017. *Edge partitioning of large graphs*. Ph.D. Dissertation. Université Pierre et Marie Curie-Paris VI.
- [35] Walter T Ludwig. 1995. *Algorithms for Scheduling Malleable and Nonmalleable Parallel Tasks*. Ph.D. Dissertation. Department of Computer Sciences, University of Wisconsin-Madison.
- [36] 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*. 195–207.
- [37] Steffen Maass, Changwoo Min, Sanidhya Kashyap, Woonhak Kang, Mohan Kumar, and Taesoo Kim. 2017. Mosaic: Processing a trillion-edge graph on a single machine. In *EuroSys*. 527–543.
- [38] 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*. 135–146.
- [39] Frank McSherry, Michael Isard, and Derek Gordon Murray. 2015. Scalability! But at what COST?. In *HotOS*. 14–19.
- [40] Robin Milner. 1989. *Communication and Concurrency*. Prentice Hall.
- [41] Timothy Prickett Morgan. 2018. *The End of Xeon Phi - It's Xeon and Maybe GPUs From Here*. Retrieved May 1, 2023 from <https://www.nextplatform.com/2018/07/27/end-of-the-line-for-xeon-phi-its-all-xeon-from-here/>
- [42] Derek G Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martin Abadi. 2013. Naiad: A timely dataflow system. In *SOSP*. 439–455.
- [43] Donald Nguyen, Andrew Lenharth, and Keshav Pingali. 2013. A lightweight infrastructure for graph analytics. In *SOSP*. 456–471.
- [44] Fatemeh Rahimian, Amir H. Payberah, Sarunas Girdzijauskas, and Seif Haridi. 2014. Distributed Vertex-Cut Partitioning. In *DAIS*. 186–200.
- [45] Ryan A. Rossi and Nesreen K. Ahmed. 2015. The Network Data Repository with Interactive Graph Analytics and Visualization. In *AAAI*. 4292–4293.
- [46] Amitabha Roy, Ivo Mihailovic, and Willy Zwaenepoel. 2013. X-stream: Edge-centric graph processing using streaming partitions. In *SOSP*. 472–488.
- [47] Julian Shun and Guy E. Blelloch. 2013. Ligma: A lightweight graph processing framework for shared memory. In *SIGPLAN*. 135–146.
- [48] Ryan Smith. 2022. *Intel To Wind Down Optane Memory Business - 3D XPoint Storage Tech Reaches Its End*. Retrieved May 1, 2023 from <https://www.anandtech.com/show/17515/intel-to-wind-down-optane-memory-business>
- [49] Stergios Stergiou, Dipen Rughwani, and Kostas Tsioutsoulis. 2018. Short-cutting label propagation for distributed connected components. In *WSDM*. 540–546.
- [50] Nilothpal Talukder and Mohammed J. Zaki. 2016. A distributed approach for graph mining in massive networks. *Data Mining and Knowledge Discovery* 30, 5 (2016), 1024–1052.
- [51] 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.
- [52] Leslie G. Valiant. 1990. A Bridging Model for Parallel Computation. *CACM* 33, 8 (1990), 103–111.
- [53] 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*. 507–522.
- [54] Guozhang Wang, Wenlei Xie, Alan J Demers, and Johannes Gehrke. 2013. Asynchronous Large-Scale Graph Processing Made Easy. In *CIDR*. 3–6.
- [55] Wikipedia. 2023. *Stone-Weierstrass Theorem*. Retrieved May 1, 2023 from [https://en.wikipedia.org/wiki/Stone%E2%80%93Weierstrass\\_theorem](https://en.wikipedia.org/wiki/Stone%E2%80%93Weierstrass_theorem)
- [56] Chenning Xie, Rong Chen, Haibing Guan, Binyu Zang, and Haibo Chen. 2015. SYNC or ASYNC: Time to fuse for distributed graph-parallel computation. In *PPoPP*. 194–204.
- [57] 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.
- [58] 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.
- [59] 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.
- [60] Kaiyuan Zhang, Rong Chen, and Haibo Chen. 2015. NUMA-Aware Graph-Structured Analytics. In *PPoPP*. 183–193.
- [61] Yanfeng Zhang, Qixin Gao, Lixin Gao, and Cuirong Wang. 2014. Maiter: An Asynchronous Graph Processing Framework for Delta-Based Accumulative Iterative Computation. *TPDS* 25, 8 (2014), 2091–2100.
- [62] Yunming Zhang, Mengjiao Yang, Riyadh Baghdadi, Shoaib Kamil, Julian Shun,

- and Saman Amarasinghe. 2018. Graphit: A high-performance graph DSL. *PACMPL* 2, OOPSLA (2018), 1–30.
- [63] Xiaowei Zhu, Wenguang Chen, Weimin Zheng, and Xiaosong Ma. 2016. Gemini: A computation-centric distributed graph processing system. In *USENIX OSDI*. 301–316.
- [64] 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*. 375–386.