

# Towards Communication-Efficient Out-of-Core Graph Processing on the GPU

Qiang Wang<sup>✉</sup>, Xin Ai, Yongze Yan<sup>✉</sup>, Shufeng Gong<sup>✉</sup>, Yanfeng Zhang<sup>✉</sup>, Jing Chen<sup>✉</sup>,  
and Ge Yu<sup>✉</sup>, *Senior Member, IEEE*

**Abstract**—The key performance bottleneck of large-scale graph processing on memory-limited GPUs is the host-GPU graph data transfer. Existing GPU-accelerated graph processing frameworks address this issue by managing the active subgraph transfer at runtime. Some frameworks adopt explicit transfer management approaches based on explicit memory copy with filter or compaction. In contrast, others adopt implicit transfer management approaches based on on-demand accesses with the zero-copy mechanism or unified virtual memory. Having made intensive analysis, we find that as the active vertices evolve, the performance of the two approaches varies in different workloads. Due to heavy redundant data transfers, high CPU compaction overhead, or low bandwidth utilization, adopting a single approach often results in suboptimal performance. Moreover, these methods lack effective cache management methods to address the irregular and sparse memory access pattern of graph processing. In this work, we propose a hybrid transfer management approach that takes the merits of both two transfer approaches at runtime. Moreover, we present an efficient vertex-centric graph caching framework that minimizes CPU-GPU communication by caching frequently accessed graph data at runtime. Based on these techniques, we present HytGraph, a GPU-accelerated graph processing framework, which is empowered by a set of effective task-scheduling optimizations to improve performance. Experiments on real-world and synthetic graphs show that HytGraph achieves average speedups of  $2.5\times$ ,  $5.0\times$ , and  $2.0\times$  compared to the state-of-the-art GPU-accelerated graph processing systems, Grus, Subway, and EMOGI, respectively.

**Index Terms**—GPU, graph processing, communication reduction, transfer management, out-of-core processing.

## I. INTRODUCTION

HIGH-performance graph processing is crucial for real-world graph applications, such as geo-information mining and social network analysis. Compared with CPU-based graph processing frameworks, GPU-based graph processing frameworks attract increased attention for their ability to leverage

GPU's high memory bandwidth and massive parallelism [23], [38], [44], [48], [53]. However, the limited memory capacity of GPUs presents challenges in handling large-scale real-world graphs, especially when their sizes exceed the available GPU memory.

Recently, research [13], [14], [29], [32], [39], [40], [45], [52] has been directed toward developing GPU-accelerated graph processing systems that leverage high-performance GPU processing and the substantial memory capacity of the CPU. Similar to that of out-of-core graph processing [25], [37], [43], [54] on the CPUs, GPU-accelerated graph processing suffers from low GPU utilization caused by extensive CPU-to-GPU data movement overhead. Accessing data from the CPU needs data migration over the low-bandwidth PCIe interconnect (up to 32 GB/s for PCIe 4.0), which can be an order of magnitude slower than the global memory access. Moreover, advances in PCIe interconnects have not effectively bridged the bandwidth gap, as the memory bandwidth of GPUs also increases simultaneously [35], [36]. This highlights the necessity of optimizing GPU-CPU data transfer.

Existing GPU-accelerated frameworks [13], [18], [32], [39], [40], [45], [52] mitigate data communication by tracking the evolving active vertices throughout iterative computation. In vertex-centric graph processing, computations are executed in a sequence of iterations, each processing vertices updated and marked as active in the prior iteration (i.e., active vertices), updating the out-going neighbors, and marking any neighbors with modified values as active vertices for the next iteration. In this process, it is necessary to access the edge data associated with active vertices (i.e., active subgraphs) [52]. Following existing systems [13], [18], [32], [39], [40], [45], [52], we assume that data related to vertices (such as value, neighbor index, and activity status) can reside in the GPU memory, the edge-associated data (including neighbor identities and weights) is entirely accommodated in the CPU memory, and, subgraphs comprising active edges are dynamically transferred to the GPU during iterative processing.

According to the way of reducing CPU-GPU active subgraph transfer, existing systems can be classified into two categories: **Explicit Transfer Management (ExpTM)**-based frameworks [18], [39], [40], [42], [52] and **Implicit Transfer Management (ImpTM)**-based frameworks [13], [32], [45]. In ExpTM-based frameworks, active subgraph communication is managed by the programmers. The oversized graph is split into small graph partitions, each of which can fit within the GPU

Received 11 April 2024; revised 14 January 2025; accepted 18 February 2025. Date of publication 4 March 2025; date of current version 7 April 2025. This work was supported by the National Key R&D Program of China under Grant 2023YFB4503601, in part by the National Natural Science Foundation of China under Grant U2241212, Grant 62202088, and Grant 62137001, in part by the 111 Project under Grant B16009, in part by the Fundamental Research Funds for the Central Universities under Grant N2216015 and Grant N2416011, and in part by the Distinguished Youth Foundation of Liaoning Province under Grant 2024021148-JH3/501. Recommended for acceptance by D. Li. (*Corresponding authors: Yanfeng Zhang.*)

The authors are with the School of Computer Science and Engineering, Northeastern University, Shenyang 110819, China (e-mail: wangqiang94@gmail.com; aixin0@stumail.neu.edu.cn; yongzeyan@stumail.neu.edu.cn; gongsf@mail.neu.edu.cn; zhangyf@mail.neu.edu.cn; jingchen@stumail.neu.edu.cn; yuge@mail.neu.edu.cn).

Digital Object Identifier 10.1109/TPDS.2025.3547356

TABLE I  
RUNTIME COMPARISON OF SUBWAY AND EMOGI ON VARIABLE ALGORITHMS  
AND DATASETS

	SK-2005 graph		PageRank Algorithm	
	SSSP	PageRank	sk-2005	uk-2007
Subway	14.6(s)	8.7(s)	8.7(s)	16.9(s)
EMOGI	7.5(s)	18.6(s)	18.6(s)	12.4(s)

memory. Before these subgraphs are transferred to the GPU via the explicit CUDA memory copy (`cudaMemcpy`), they must be processed by a CPU-based redundancy removal module to eliminate inactive edges. Depending on the operation mode, these approaches can be either light-weight filter-based [18], [40] or heavy-weight compaction-based [39], [52].

Recently, ImpTM-based approaches that circumvent the explicit management of data movements for active subgraphs have emerged [13], [32], [45]. ImpTM-based frameworks utilize Unified Virtual Addressing (UVA) technique to map CPU and GPU to the same memory address, allowing GPUs to directly access the required active edges in the CPU [4], [5], [13], [32]. Compared with ExpTM, ImpTM requires less engineering efforts, allowing users to directly extend a single GPU framework into an out-of-core one by managing communication through unified virtual memory (UVM) [13], [45] or zero-copy access [32]. During iterative processing, memory slices containing active edges can be transferred to the GPU implicitly in a user-transparent manner. Due to the fixed data migration mechanisms of ImpTM approaches, the transfer efficiency is highly sensitive to the graph's access pattern.

Having made extensive analysis, we find that a decision to choose one or the other approach for the best performance is determined by the memory access pattern of active edges. In GPU-accelerated graph processing frameworks based on a single approach, the performance is often suboptimal. Table I shows the performance comparison of Subway [39] (a ExpTM-compaction-based framework) and EMOGI [32] (an ImpTM-zero-copy-based framework). On sk-2005 graph [3], EMOGI outperforms the Subway on Single Source Shortest Path algorithm (SSSP), but it losses on PageRank. In contrast, for the PageRank algorithm, Subway beats EMOGI on SK dataset [3], but losses on U.K. dataset [3]. Moreover, these approaches focus on reducing redundant active subgraph transmission and lack efficient mechanisms to reuse transferred graph data iterations. Existing frameworks [13], [42] based on page-centric data caching can hardly adapt to graph processing tasks with irregular and sparse memory access patterns.

We present HyTGraph, a GPU-accelerated Graph processing system that distinguishes itself from previous frameworks by not exclusively relying on either ExpTM or ImpTM. Instead, our system employs a **Hybrid Transfer Management** method (HyTM) that combines ExpTM and ImpTM to maximize performance. HyTM splits the graph into small partitions as ExpTM does. During iterative processing, HyTM estimates ExpTM cost and ImpTM cost on-the-fly by analyzing the edge access pattern of each partition, and chooses the most cost-efficient transfer method. Building upon HyTM method, HyTGraph provides an effective vertex-centric graph caching method incorporating

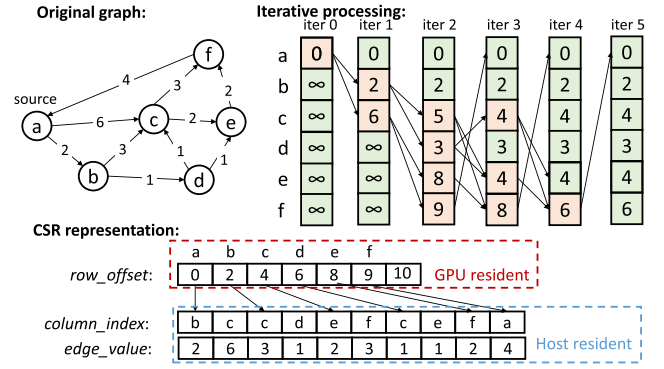


Fig. 1. An example of SSSP computation. The orange box indicates the active vertices and the green box indicates the inactive one. The input graph is organized into a CSR format.

fine-grained cache management and GPU parallel refreshing for efficient graph transfer reusing across iterations. Furthermore, HyTGraph provides a contribution-driven asynchronous scheduling to accelerate convergence. Experimental results on real-world and synthetic graphs demonstrate that HyTGraph achieves an average speedup of 5.0X over Subway [39], 2.0X over Grus [45], and 2.5X over EMOGI [32].

This work extends the conference version [46], enhancing the communication efficiency through fine-grained GPU graph caching: (1) We present a vertex-centric graph caching framework that reduces CPU-GPU communication through fine-grained GPU data caching. The framework further incorporates a decoupled cache refreshing mechanism along with GPU parallel processing to achieve low-overhead cache management (Section VII). (2) We conduct experimental studies to verify the effectiveness of the caching method in Section IX. These include: (a) Updates to the overall performance evaluation and discussion of HyTGraph in Sections IV and IX-G, incorporating the caching optimizations (b) Supplement data caching components in Sections IX-E and IX-F to assess the data transfer reduction and consequent performance improvement; (c) New experimental evaluations in Section IX-F to assess the cost and benefit of GPU data caching.

## II. BACKGROUND

### A. Vertex-Centric Graph Processing

Vertex-centric programming [16], [30] has been widely adopted in graph processing frameworks for its simplicity and powerful expression ability. It uses a generic function to define the behavior of a vertex and its neighbors. Considering the message passing direction, the function can be either pull- or push-based [44]. During computation, the function is iteratively evaluated on all vertices until the algorithm terminates. Fig. 1 illustrates an example of SSSP, an algorithm to find the shortest paths from a given source vertex to all the other vertices. It starts from a source vertex *a*. In each iteration, the accessed vertices broadcast their shortest distances to the outgoing neighbors and update the shortest distance if the new path is shorter than the old one. The algorithm converges when no more vertices

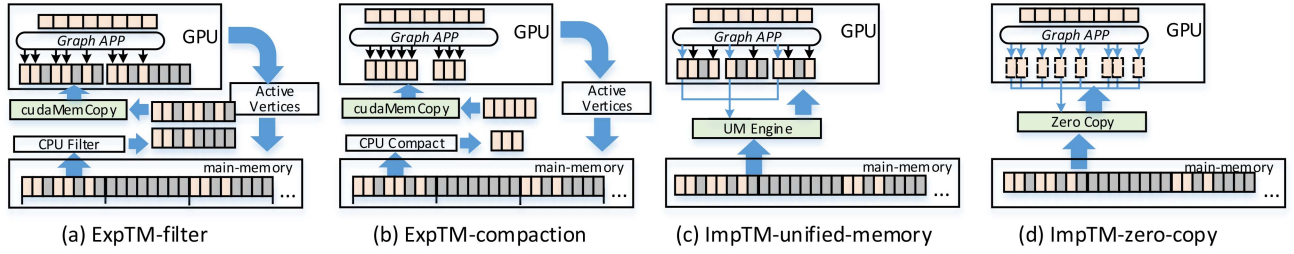


Fig. 2. An example of the four approaches. The thin blue arrow, thin black arrow, and thick blue arrow represent the remote memory access, local memory access, and host-GPU data transfer, respectively. ExpTM-based approaches need to transfer the active vertices back to the host side for compaction or filtering.

are updated. During iterative computation, only the vertices updated by the previous iteration (active vertices) need to be processed.

**GPU graph processing:** Recent research explores the massive parallelism of GPUs [23], [38], [44], [48], [53] to accelerate graph processing. Despite achieving promising results, the processing capability of these studies is limited by the GPU memory. For example, a common 16 GB GPU can accommodate only about 600 million edges (assuming each edge occupies 8 bytes). While multi-GPU processing is an intuitive approach for scaling to large-scale graphs, the cost of expanding GPU memory is prohibitively high, with the price of 1 GB of GPU memory often tens or even hundreds of times that of 1 GB of CPU memory [34], [35]. Fortunately, we observe that vertex data, which requires frequent random accesses, is often exponentially smaller than edge data. Although edge data consumes significant memory, it is read-only and typically accessed sequentially. Therefore, placing vertex data on GPUs and offloading large-scale edge data to the CPU provides a cost-effective solution that leverages the high GPU computational power and the affordable host memory resources [32], [39]. A heterogeneous computing platform with a single 16 GB GPU and 320 GB of CPU memory can handle large-scale graphs with hundreds of millions of vertices and tens of billions of edges. In contrast, achieving this with multi-GPU in-memory processing could require dozens of such GPUs, along with the complexity of optimizing inter-GPU communication. If the input graph scales further, combining out-of-memory processing with multi-GPUs could provide a promising solution. This would involve partitioning vertex data across multiple GPUs and exploring new interconnect technologies to extend host memory, which is discussed as future work in Section X. In this work, we focus on optimizing the data communication in CPU-GPU heterogeneous computing, especially exploring efficient transfer management modules and caching mechanisms to minimize unnecessary edge communications.

## B. ExpTM Approaches

**ExpTM-filter:** GraphReduce [40], GTS [24], and Graphie [18] adopt a filter-based method to reduce the inactive subgraph transfer. They monitor the active edges of the partitioned subgraphs and transfer only those containing active edges. Fig. 2(a) shows an illustrative example. This method filters out partitions containing no active edges without additional processing. Therefore, active partitions will be entirely transferred to the GPU, even if

only one edge is active. When the proportion of active edges is low, the volume of unnecessary data transfer will be large.

**ExpTM-compaction:** In contrast, some other frameworks [39], [42], [52] introduce CPU-assisted compaction to reduce communication. Before transferring a partition containing active edges to the GPU, these frameworks use CPUs to remove the inactive subgraph and compact the remaining data into a continuous memory space to facilitate explicit memory copy. Fig. 2(b) shows an illustrative example of Subway [39], a typical ExpTM-compaction-based system. Compared with the filter-based frameworks [18], [40], compaction-based frameworks can minimize the data transfers by removing all inactive edges. But at the cost, it involves additional CPU and memory manipulation overhead.

## C. ImpTM Approaches

**ImpTM-unified-memory:** Unified-virtual-memory (UVM) defines a managed memory space where both GPU and CPU share a single address space, maintaining a coherent memory image [13], [45]. Fig. 2(c) shows an illustrative example. During computation, memory pages (4 KB in default) containing requested data are automatically migrated to GPUs. UVM provides page-centric GPU data caching, enabling subsequent accesses to the same memory page to be served directly from the GPU's global memory, thus avoiding additional data transfers. However, the "automated migration" cost is not free. Migrating new pages to the GPU memory triggers page-fault processing, which requires not only data transfer but also significant costs in Translation Lookaside Buffer (TLB) invalidation and page fault updating [32].

**ImpTM-zero-copy:** Zero-copy memory access offers a more lightweight approach. The method maps pinned CPU memory to GPU address spaces, enabling GPUs to directly access CPU memory through the Transaction Layer Packet (TLP) of PCIe [32]. Compared to UVM-based method, zero-copy access provides a finer granularity of CPU-GPU data access. As per the PCIe 3.0 specification, up to 256 outstanding memory requests can be processed concurrently by each TLP, with each request accommodating data payloads of 32, 64, 96, or 128 bytes [32], tailored to the size of the various adjacency lists. This capability allows zero-copy access to facilitate simultaneous, fine-grained access to the edge data. Moreover, zero-copy access incurs lower transfer overhead than UVM-based approaches as it eliminates the need for page-fault handling. As a sacrifice, zero-copy access



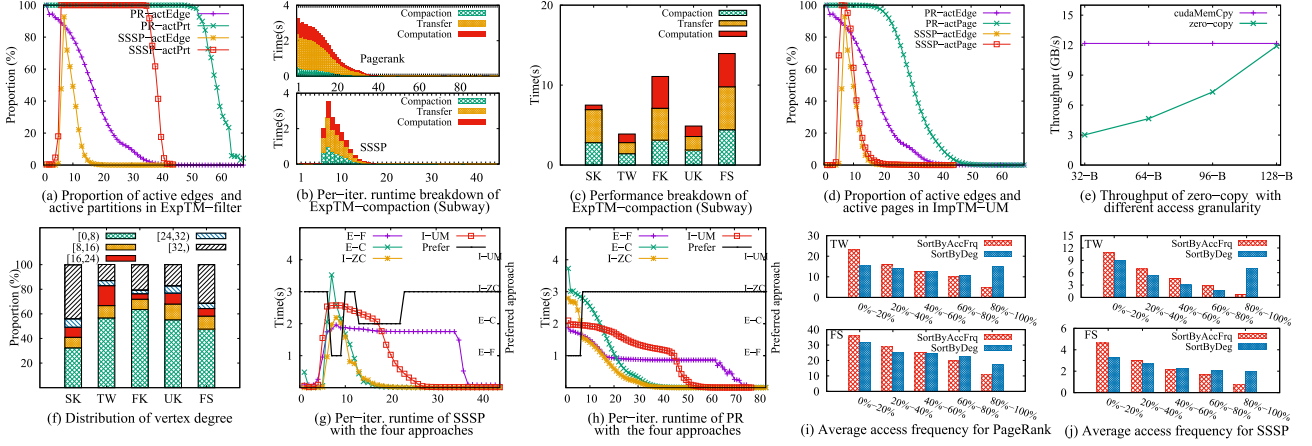


Fig. 3. Performance analysis of the four engines using the two algorithms.

cannot provide data-reusing functions. Repeat accesses to the same data will cause multiple separate CPU-GPU data transfers.

### III. ANALYSIS OF EXISTING APPROACHES: A MOTIVATING STUDY

In this section, we conduct an experimental analysis of the existing approaches using two graph algorithms with diverse memory access patterns: SSSP and PageRank.

#### A. Analysis of ExpTM

*ExpTM-filter (ExpTM-F)*: As mentioned above, filter-based ExpTM has a large amount of redundant transfers even if the proportion of active edge is low. We run PageRank and SSSP on friendster-konect [2] graph to explore the redundant data transfer problem with the partition number set to 256. Fig. 3(a) shows the proportion curves of active edges and partitions containing active edges (active partitions). We can observe that the proportion of active partitions does not decrease immediately with the proportion of active edges. For SSSP and PageRank algorithms, the active edges account for only 28.3% and 12.3% of the total transfer volume. Although ExpTM-filter method shows ineffectiveness with a small number of active edges, it exhibits advantages when dealing with subgraphs containing a large proportion of active edges. This is because it can maximize the utilization of PCIe bandwidth through the `cudaMemcpy()` function.

*ExpTM-compaction (ExpTM-C)*: The compaction-based ExpTM significantly reduces data communication by transferring compacted active data. However, this approach involves substantial overhead due to CPU-based compaction processing, particularly when a large proportion of active edges are involved. As highlighted by Subway [39], in scenarios where the proportion of active edges is high (e.g., 80%), the overhead associated with compaction can even outweigh the benefit of transfer reduction [39]. Fig. 3(b) depicts the per-iteration runtime breakdown of Subway and showcases instances where costs exceed benefits. Additionally, Fig. 3(c) shows the overall performance breakdown of SSSP algorithm on Subway, with the

preprocessing overhead excluded from the execution time. We can observe that across all five datasets, the compaction stage accounts for 34.5% of the overall runtime.

#### B. Analysis of ImpTM

*ImpTM-unified-memory (ImpTM-UM)*: Unified-virtual-memory provides on-demand memory migration. However, this method falls short of efficiency in graph processing applications. On the one hand, recent studies show that the peak bandwidth of unified-memory can only reach 73.9% of that of explicit memory copy due to the high “automated migration” overhead [32]. On the other hand, the migration granularity of 4 K bytes often leads to the inclusion of substantial inactive data [13], [32] in each memory page. This issue is caused by the inherent irregularity and poor locality of graph algorithms [32], [42]. Fig. 3(d) shows the proportion of active edges and active memory pages of each iteration. For SSSP and PageRank algorithms, the active edges account for only 54.5% and 65.0% of the total transfer volume, respectively. Therefore, ImpTM-unified-memory exhibits suboptimal communication performance on large graphs, regardless of whether the ratio of active edges is high or low. Additionally, the page-centric data caching also contributes little to the performance improvement due to the low space utilization of active subgraphs in the memory pages. The UVM-based approach only shows efficiency when the graph size is sufficiently small to be fully accommodated within the GPU memory.

*ImpTM-zero-copy (ImpTM-ZC)*: Maximizing PCIe bandwidth utilization is crucial for improving zero-copy access performance. As indicated by EMOGI [32], saturating most of the 256 memory requests in each TLP with 128-byte data is essential for maximizing the PCIe bandwidth. The underlying reason is as follows: Each TLP not only carries the payloads of memory requests but also includes a header field containing control information. Smaller memory requests result in a higher number of TLPs for the same data volume, thus allocating a significant portion of the bandwidth to header field transfers. Fig. 3(e) shows the throughput of zero-copy access across various

memory request sizes (ranging from 32 bytes to 128 bytes). We can observe that when the memory request size is 128 bytes, the zero-copy access can achieve almost the same performance as `cudaMemcpy` (the maximum PCIe utilization). However, when the request size is set to 32 bytes, the throughput decreases significantly. To optimize bandwidth usage, EMOGI [32] proposes merged and aligned optimization with which each warp of threads access consecutive neighbors of one vertex in a 128-byte cache line size. Nonetheless, guaranteeing that most memory requests reach the 128-byte mark is challenging, especially when considering that each neighbor typically requires 4 bytes, and real-world graphs, adhering to a power-law distribution, often feature vertices with fewer than 32 neighbors. As shown in Fig. 3(f), on average, 74.7% of vertices have fewer than 32 neighbors, with 51.1% having fewer than 8. This leads to inconsistent performance of zero-copy access across real-world graphs.

### C. Performance Comparison of the Four Approaches

We report the per-iteration runtime of ExpTM-filter, ExpTM-compaction, ImpTM-unified-memory, and ImpTM-zero-copy on friendster-konect [2] using two typical graph algorithms: the graph traversal algorithm SSSP and the iterative algorithm PageRank [44]), illustrated in Fig. 3(g) and (h). The implementations of ExpTM-filter, ImpTM-unified-memory, and ImpTM-zero-copy leverage the processing kernel of SEP-Graph [44]. The `cudaMemAdviseSetReadMostly` optimization is enabled for ImpTM-unified-memory (the evicted memory pages will be discarded directly instead of written back to the CPU memory). For the ExpTM-compaction implementation, We use Subway [39] due to its highly-optimized compaction engine and efficient GPU kernel from Tigr [38]. All methods are set to execute synchronously to maintain a consistent number of active vertices across iterations. We employ a “Prefer” curve to indicate the best-performing approach in each iteration. By examining the proportion curves of active edges for SSSP and PageRank in Fig. 3(a), we observe that ExpTM-filter excels when the proportion of active edges is high, attributable to its efficient utilization of PCIe bandwidth.

When the proportion of active edge is low, ImpTM-zero-copy outperforms other approaches in most iterations, thanks to its fine-grained data migration mechanism. However, for the SSSP algorithm, ExpTM-compaction occasionally outperforms ImpTM-zero-copy during certain iterations. This variability can be attributed to the unstable performance of zero-copy under different vertex degrees. As mentioned above, zero-copy’s effectiveness is affected not only by the active edge proportion but also by the average degree. Given a fixed number of active edges, a higher count of active vertices leads to a higher number of fragmented edge data accesses and an increase in partially filled TLP requests, leading to reduced bandwidth utilization efficiency. ImpTM-unified-memory consistently underperforms across all scenarios, even with GPU data caching. Its coarse-grained page migration mechanism falls short of recognizing the irregular and sparse data access patterns of graph processing, thereby diminishing both cache and PCIe bandwidth efficiency.

### D. Cross-Iteration Data Reusing

Existing transfer management approaches primarily focus on intra-iteration communication optimization. However, duplicated vertex accesses across iterations also result in certain edge data being transferred multiple times. As shown in Fig. 3(i-j), we evaluate the SSSP and PageRank algorithms on the TW and FS graphs. Vertices are sorted and grouped into five categories based on their access frequency (red bars) and vertex degree (blue bars), and the average access frequency for each group is presented. The results reveal that the top 20% most frequently accessed vertices exhibit significantly higher access frequencies than the average (black line) and other groups, offering considerable opportunities for communication reuse. Existing data caching mechanisms, however, exhibit certain limitations. The ImpTM-unified-memory approach, which employs page-centric dynamic data migration and caching, fails to adapt effectively to the irregular access patterns of graph data. Recent studies have explored the significance of high-degree vertices (i.e., hub vertices) in graph processing, highlighting their frequent access during computation [28]. Consequently, extracting hub vertices for data caching becomes a natural and intuitive solution. However, we observe that high-degree vertices in different graph processing tasks do not always align with high-frequency accessed vertices, as shown by the blue columns in Fig. 3(i-j). On the TW dataset, the bottom 20% low-degree vertices exhibit even higher access frequencies than the middle groups. Moreover, the overlap between the top 20% of vertices ranked by degree and those ranked by access frequency is limited, ranging from 34% to 45%. This inspires us to design a more targeted, access frequency-guided caching mechanism for efficient data reuse.

### E. Summary of Existing Approaches

Having made intensive analysis, we observe that the data transfer overhead is affected by three primary factors: data transfer volume, PCIe bandwidth Utilization, and CPU pruning cost. As depicted in Table II, existing approaches fail to optimize these factors cooperatively. Adopting a single transfer management method often limits the system’s effectiveness to only one or several specific scenarios, as shown in the last column in Table II. Additionally, existing GPU data caching mechanisms, which rely on fixed-size memory page migration or degree-based subgraph pre-selection, are not suited for handling the irregular data access patterns inherent in graph processing. In addition to systems [13], [32], [39], [40] mentioned above, Scaph [52] adopt ExpTM-compaction. Different from Subway, Scaph [52] and Ascetic perform compaction on the partitioned graph. It distinguishes the partitions with a small proportion of active edges and compacts them for subsequent GPU processing. In contrast, partitions with a large proportion of active edges are entirely loaded to the GPU. Grus [45] is an ImpTM-based framework that manages the edge-associated data in main memory with priorities, prefetching high-priority data to the GPU through unified-memory and accessing low-priority data through zero-copy. In addition, some frameworks [14], [29] also use CPU-GPU co-processing to accelerate graph processing. We review them in Section XI.

TABLE II  
SUMMARY OF EXISTING APPROACHES AND REPRESENTATIVE SYSTEMS

Approach	System	Transfer Volume	Bandwidth Utilization	Prune Cost	GPU Data Cache	Preferred Scenario
ExpTM-F	GraphReduce [40] Graphie [18] GTS [24]	High	High	Low	N/A	• Subgraph with a large proportion of active edges
ExpTM-C	Subway [39] Scaph [52] Ascetic [42]	Low	High	High	N/A	• Subgraph with a small proportion of active edges and small average degree
ImpTM-UM	HALO [13] Grus [45]	Medium	Low	NA	Page-centric	• Small graph that can fit into GPU memory
ImpTM-ZC	EMOGI [32]	Low	Unstable	No	N/A	• Subgraph with a small proportion of active edges and high average degree
HyTM	Our approach	Low	high	Low	Vertex-centric	• Adapt to subgraph with various active degree

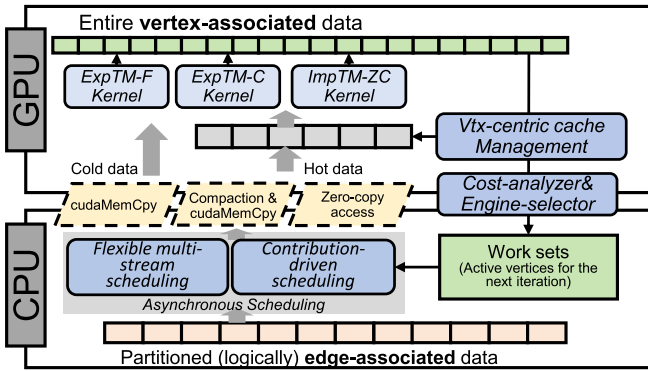


Fig. 4. Overview of HyTGraph.

#### IV. HYTMGRAPH OVERVIEW

We present HyTGraph, a GPU-accelerated graph processing system based on hybrid transfer management (HyTM) to maximize performance. Fig. 4 shows an overview. HyTGraph organizes the graph into a CSR structure. Following [18], [52], HyTGraph logically splits the edge data into  $N$  edge-balanced partitions  $\{P_0, P_1, \dots, P_{N-1}\}$ . During computation, partitions containing active vertices are processed with their most cost-efficient engines. HyTGraph provides three functions to achieve efficient HyTM.

*Cost-based engine selection:* HyTGraph uses a GPU-resides cost analyzer and engine selector to compute the data transfer overhead for different approaches and selects the most cost-efficient one for each partition. Based on the analysis in Section III, we choose ExpTM-F, ExpTM-C, and ImpTM-ZC as engine candidates. In addition, HyTGraph provides a task combining optimization to merge small subgraph pieces into larger tasks to minimize scheduling overhead.

*Contribution-driven asynchronous scheduling:* HyTGraph introduces asynchrony to improve computation efficiency. Rather than simply recomputing the loaded subgraph multiple times [39], [52], HyTGraph adopts a contribution-driven priority scheduling method to prioritize those partitions that contribute more to convergence. To further improve resource utilization, HyTGraph uses multiple CUDA streams to overlap the computation kernel, data transfer, and CPU-based active subgraph compaction.

*Vertex-centric graph cache:* HyTGraph provides a vertex-centric graph caching method, distinguishing itself from the page-centric methods in existing GPU graph processing systems [13], [42]. HyTGraph finely caches the frequently accessed data and compactly stores them in the GPU to maximize the cache utilization. Moreover, instead of performing the heavy-weight individual-vertex data replacement, HyTGraph uses a batched cache refreshing method that minimizes data replacement overhead through GPU parallel processing.

#### V. COST-BASED ENGINE SELECTION

##### A. Cost Analysis and Engine Selection

Most existing activeness-tracking-based frameworks use the proportion of active edges as the evaluation metric [18], [29], [40], [52] to select appropriate processing engines, which provides an intuitive and lightweight distinguishing method. However, it is not suitable in HyTM as the proportion of active edges cannot reflect the cost of different transfer approaches (as detailed in III-C). In this work, we directly model the overhead of different transfer management methods and accurately select the most cost-efficient execution engine.

*Cost of ExpTM-filter:* The ExpTM-filter-based approach entirely transfers partitions containing active edges to the GPU using `cudaMemcpy`, thus incurring only data transfer overhead. This overhead can be estimated by the saturated TLPs (as discussed in Section III, Fig. 3(e)). Given a partition  $P_i$ , the number of memory requests can be calculated with  $\sum_{v \in P_i} D_o(v) * d_1/m$ , where  $\sum_{v \in P_i} D_o(v)$  is the edge number of partition  $i$ ,  $D_o(v)$  represents the number of neighbors of  $v$ ,  $d_1$  represents the memory occupation of a vertex, and  $m$  represents the maximum capacity of an outstanding memory request (128-bytes in PCIe-3.0 specification). Denote  $MR$  as the maximum number of outstanding memory requests in TLP ( $MR = 256$  in PCIe 3.0) and  $\lceil \cdot \rceil$  as the round-up operation, we formalize the transfer overhead of each  $P_i$  as follow:

$$Tef_i = \left\lceil \left( \sum_{v \in P_i} D_o(v) \right) * d_1/m/MR \right\rceil * RTT, \quad (1)$$

where  $\lceil (\sum_{v \in P_i} D_o(v)) * d_1/m/MR \rceil$  is the number of required TLPs, and  $RTT$  represents the round trip time for PCIe to process a saturated TLP request.



*Cost of ExpTM-compaction:* ExpTM-compaction involves additional CPU-based compaction, so its cost consists of two parts: the data transfer overhead and the CPU compaction overhead. Since the compaction needs to reorganize the active edges and change their positions, we also need to generate a vertex index array and transfer it to GPU for addressing the compacted neighbors. Then the transfer volume can be formalized as  $\sum_{v \in A_i} D_o(v) * d_1 + |A_i| * d_2$ , where  $A_i$  represents the active vertex subset of  $P_i$  and  $d_2$  represents the memory occupation of each vertex index. The CPU-based compaction is related to both data transfer volume and CPU compaction throughput, which can be computed using  $\sum_{v \in A_i} D_o(v) * d_1 + |A_i| * d_2 / Thpt_{cpt}$ , where  $Thpt_{cpt}$  represents the CPU compaction throughput. Then, the cost of ExpTM-compaction can be formalized as:

$$T_{ec_i} = \left\lceil \left( \sum_{v \in A_i} D_o(v) * d_1 + |A_i| * d_2 \right) / m / MR \right\rceil * RTT + \sum_{v \in A_i} D_o(v) * d_1 + |A_i| * d_2 / Thpt_{cpt} \quad (2)$$

*Cost of ImpTM-zero-copy:* The ImpTM-zero-copy approach provides vertex-centric on-demand access in the cacheline size. Therefore, each active vertex  $v$  takes one or several independent memory requests depending on its neighbor size. Generally, the number of required memory requests of vertex  $v$  can be formalized as  $\lceil D_o(v) * d_1 / m \rceil$ . Considering that we can hardly guarantee the neighbors of all vertices start from a cacheline-aligned position, some vertices may have the misaligned neighbor array and thus require one additional memory transaction [32]. We introduce a function  $am()$ , which returns 1 for unaligned vertices and 0 for the others.<sup>1</sup> Therefore, the transfer overhead of ImpTM-zero-copy can be formalized as:

$$T_{iz_i} = \left\lceil \left( \sum_{v \in A_i} (\lceil D_o(v) * d_1 / m \rceil + am(v)) \right) / MR \right\rceil * RTT_{zc} \quad (3)$$

where  $(\sum_{v \in P_i(V) \cap A_i} \lceil D_o(v) * d_1 / m \rceil + am(v))$  is the required memory transactions of  $P_i$ . It should be noted that the TLP round trip time of zero-copy ( $RTT_{zc}$ ) is not the same as that in ExpTM ( $RTT$ ) because the payload of each TLP in zero-copy may be unsaturated. This makes  $RTT_{zc}$  always less than the  $RTT$ s in ExpTM-filter and ExpTM-compaction. In this paper, we use a dumping factor  $\gamma$  to compute  $RTT_{zc}$  for each partition as follows:  $RTT_{zc} = \gamma * RTT + (1 - \gamma) * (\sum_{v \in A_i} D_o(v) / \sum_{v \in P_i} D_o(v)) * RTT$ , where  $(\sum_{v \in A_i} D_o(v) / \sum_{v \in P_i} D_o(v))$  is the proportion of active edge.  $\gamma * RTT$  represents the fixed time to process a TLP, and  $(1 - \gamma) * (\sum_{v \in A_i} D_o(v) / \sum_{v \in P_i} D_o(v)) * RTT$  represents the time related to the size of payload. By referring to [32], we set  $\gamma$  to 0.625.

*Transfer engine selection:* HyTGraph compares  $T_{ef_i}$ ,  $T_{ec_i}$ , and  $T_{iz_i}$  to choose the most cost-efficient execution engine.

<sup>1</sup>In the implementation, the number of memory requests of each active vertex  $\lceil D_o(v) * d_1 / m \rceil + am(v)$  can be directly computed by using the length and physical position of the neighbor data.

---

**Algorithm 1:** Cost-Based Engine Selection.

---

**Input:** active vertex set  $\{A_0, \dots, A_{N-1}\}$  of  $N$  partitions,  
**Output:** tasks prefer ExpTM-filter  $\{Vf_0 \dots Vf_{M-1}\}$   
 ( $M < N$ ), task prefer ExpTM-compaction  $Vc$ , and task prefer ImpTM-zero-copy  $Vz$ .

1: initialize a selection array  $\{p_0, \dots, p_{N-1}\}$  on GPU.

**Cost analysis and engine selection:**

2: **for** each  $A_i$  in  $\{A_0, \dots, A_{N-1}\}$  **do** in parallel

3:   Compute  $T_{ef_i}$ ,  $T_{ec_i}$ , and  $T_{iz_i}$  according to Formula (1,2,3)

4:   **if**  $T_{ec_i} < \alpha * T_{ef_i}$  and  $T_{ec_i} < \beta * T_{iz_i}$  **then**

5:      $p_i = \text{'ExpTM-C'}$ ;

6:     insert  $A_i$  to  $Vc$ ; //pre-combine on GPU

7:   **else if**  $T_{ef_i} < T_{iz_i}$  **then**

8:      $p_i = \text{'ExpTM-F'}$ ;

9:   **else**

10:      $p_i = \text{'ImpTM-ZC'}$ ;

11:     insert  $A_i$  to  $Vz$ ; //pre-combine on GPU

12:   **end if**

13: **end for**

14: Copy  $Vc$ ,  $\{p_0, \dots, p_{N-1}\}$  and  $\{A_0, \dots, A_{N-1}\}$  to host. **Task Combination:**

15:  $i = 0, j = 0, length = 0$ ;

16: **while**  $i < N$  **do**

17:   **if**  $p_i = \text{'ExpTM-F'}$  and  $length < k$  **then**

18:     insert  $A_i$  to  $Vf_j$ ;

19:      $length = length + 1$ ;

20:   **else**

21:      $length = 0, j = j + 1$ ;

22:   **end if**

23:    $i = i + 1$ ;

24: **end while**

---

However, theoretically modeling the CPU compaction throughput  $Thpt_{cpt}$  is a challenging task because ExpTM-compaction introduces parallel and random writes on the host memory. This makes  $Thpt_{cpt}$  vary with active edges nonlinearly. In practice, we compute  $T_{ec_i}$  by considering only the transfer overhead and compare it with  $T_{ef_i}$  and  $T_{iz_i}$ . If  $T_{ec_i}$  is less than  $\alpha * T_{ef_i}$  and  $T_{ec_i}$  is less than  $\beta * T_{iz_i}$ , we choose ExpTM-compaction. The first condition comes from Subway's observation [39], where  $\alpha$  is set to 80%. The second condition is based on the observation from Section III: when a partitioned subgraph has few active edges but many active vertices, and zero-copy requires many unsaturated memory requests to transfer the data, leading to ExpTM-compaction is a better choice than ImpTM-zero-copy. In the implementation,  $\beta$  is set to 40%. If these conditions are not met, we compare  $T_{iz_i}$  with  $T_{ef_i}$ . If  $T_{iz_i}$  is less than  $T_{ef_i}$ , we choose ImpTM-zero-copy. Otherwise, we choose ExpTM-filter. The value of  $RTT$  can be arbitrarily specified, as it can be eliminated in subsequent comparisons. Since the cost computation between partitions is independent, HyTGraph performs the engine selection on the GPU for high performance. Algorithm 1 line (2-13) shows the overall execution flow.

### B. Partition Granularity and Task Combination

A key to implementing high-performance hybrid transfer management is to determine the optimal granularity for task scheduling. The existing frameworks [18], [29], [40], [52] directly use the partitioned subgraphs as scheduling unit. This method is simple and straightforward, however, can lead to low task scheduling performance. If the partition size is too large, the coarse-grained cost computation may lead to inappropriate engine selection. If the partition size is small, it may lead to higher kernel scheduling overhead and fragmented data transfers for more partitions.

To achieve fine-grained engine selection and low overhead task scheduling simultaneously, HyTGraph decouples the graph partitioning and task scheduling and optimizes them separately. HyTGraph partitions the graph into small partitions (32 MB each partition) to provide fine-grained cost analysis. While in the computation, HyTGraph packages partitions choosing the same engine into large chunks for processing. Specifically, for partitions using ExpTM-filter, HyTGraph merges  $k$  consecutive partitions into a large one ( $k = 4$  in HyTGraph) to reduce scheduling overhead (Line 15–24 in Algorithm 1). For partitions using ExpTM-compaction, HyTGraph merges all their active vertices and writes their neighbors to a consecutive memory chunk to leverage explicit memory copy (line 6 in Algorithm 1). For partitions using ImpTM-zero-copy, HyTGraph merges all active vertices (line 11 in Algorithm 1) and processes them with one CUDA kernel to leverage the implicit computation-communication overlapping feature of zero-copy access.

## VI. CONTRIBUTION-DRIVEN ASYNCHRONOUS TASK SCHEDULING

Asynchronous computation allows newly updated results to be used immediately in subsequent computation, which has been proven to be effective in GPU-based graph processing [6], [44]. However, simply processing the transferred subgraph multiple times may cause the local updates to be abolished by the subsequent computation from other partitions, leading to increased computation and data transfer. This is known as the stale computation problem [11], [47]. Our experiments reveal that existing framework with multi-round processing can even cause higher transfer volume (See Section IX-E for details) than synchronous scheduling in some cases. To address this issue, HyTGraph adopts a contribution-driven priority scheduling method.

*Hub-vertex-driven priority scheduling:* Due to the power-law property of real-world graphs, some important vertices with high incoming and outgoing degrees often become the hubs in the computation path. These vertices become critical upstream dependencies of a large number of vertices because of the large outgoing degree. On the other hand, these vertices also have a high probability of being activated in the iterative computation due to large incoming degrees. If these vertices do not accumulate sufficient updates before being scheduled, the downstream computation results based on the current value are likely to be abolished by subsequent new updates. Based on this observation,

we propose a hub-vertex-driven priority scheduling approach. By ensuring that the hub vertices accumulate enough contributions before being scheduled, HyTGraph can reduce the possible stale computations on the downstream vertices. Implementing hub-vertex-driven scheduling in GPU-accelerated platforms is challenging, as the hub vertices distribute randomly among the entire graph, which makes hub vertices hard to gather and transfer. To solve this problem, HyTGraph adopts a hub vertex sorting method [50], that groups the top 8% important vertices at the beginning of the CSR structure. The importance score of each vertex is measured by the following equation:

$$H(v) = \frac{D_o(v) * D_i(v)}{D_{o\max} * D_{i\max}} \quad (4)$$

$D_i(v)$ ,  $D_o(v)$ ,  $D_{i\max}$ , and  $D_{o\max}$  represent the incoming-, outgoing-, maximum incoming-, and maximum outgoing- degree, respectively. In this way, the hub vertices are grouped together, and the non-hub-vertices remain in their natural order. HyTGraph recomputes the loaded subgraph only once because most updates can only pass two hops effectively [43]. Another benefit of this method is that the vertices having a high probability of being activated (with large in-degree) are stored together, which can help gather high-active subgraphs for the cost-based engine selection. It is worth mentioning that the hub sorting does not need to be performed every time. The hub vertex can be sorted only once in the preparation stage, and all subsequent executions (of different algorithms) can benefit from it.

*$\Delta$ -driven priority scheduling:* For iterative graph algorithms based on arithmetic accumulation, e.g.,  $\Delta$ -based PageRank and PHP [49], the contribution of vertices is directly reflected in their delta values (the messages to be accumulated). Prioritizing vertices with larger  $\Delta$  value can help downstream vertices accumulate more valid updates [49]. Following the original vertex-centric  $\Delta$ -driven priority scheduling, HyTGraph provides a partition-centric  $\Delta$ -driven scheduling method that computes an overall  $\Delta$  value for each partition and prioritizes those with larger contributions.

## VII. VERTEX-CENTRIC GRAPH CACHING

The proposed hybrid transfer management focuses on transfer reduction within each iteration. However, certain iterative graph algorithms (e.g., Pagerank) involve repeated data access among iterations [42], which benefit less from HyTM, and can hardly optimized with exiting caching mechanisms due to coarse-grained data caching and challenging cache vertices determination. In this section, we introduce an access frequency-guided, vertex-centric graph caching (VCGC) method. During computation, HyTGraph tracks the access frequency of all vertices in real time, sorts them in by access frequency to determine the cache candidates, and performs vertex-centric cache replacement using the hybrid communication engine. This approach leverages a more directed access frequency metric and finely transfers and caches the graph in a vertex-centric manner.

Enable high-performance VCGC is a non-trivial task as frequent eviction and loading for variable-length adjacency lists



**Algorithm 2:** Data Caching During Iterative Processing.

---

```

1:  $V_{curr} = \text{init}(\text{"ALG"})$  // initial active vertices
2:  $V_{cache} = \emptyset$  // cached vertices
3:  $N_{batch} = 0$  // counter of processed vertices
4: while  $V_{curr} \neq \emptyset$  do
5:   if  $N_{batch} > \gamma|V|$  then
6:      $V_{cand} = \text{hot\_candidates\_VCDC}(C_{cap})$ 
7:     if  $\text{overlap}(V_{curr}, V_{cache})$  then
8:        $\text{cache\_update\_VCDC}(V_{curr}, V_{cache})$ 
9:       Evict cold data and alloc space for new
         candidates
10:       $c\_tag = 1$ 
11:    end if
12:     $\text{reset\_hotness\_VCDC}()$ 
13:     $N_{batch} = 0$ 
14:  end if
15:   $\{V_c, V_z, V_f\} = \text{engine\_selection\_HyTM}(V_{curr})$ 
16:   $V_{next} = \text{proc\_HyTM\_VCGC}(\{V_c, V_z, V_f\}, V_{cand} \setminus$ 
     $V_{cache})$ 
17:    Load accessed candidates into the cache
18:    Record the number of accesses
19:   $V_{curr} = V_{next}$ 
20:   $N_{batch} += |V_{curr}|$ 
21:  if  $c\_tag == 1$  then
22:     $V_{cache} = V_{cand}$ 
23:  end if
24: end while

```

---

on the GPU incur substantial memory manipulation overhead. To address this issue, HyTGraph provides several key features, including 1) a periodic cache refreshing approach that decouples the eviction and loading of frequently accessed data; 2) CSR-based compacted data organization to maximize memory utilization, and uses GPU parallel processing to accelerate the cache determination. We first introduce the overall execution flow and then detail the design and implementation of each component.

*Execution flow:* The cache management flow during iterative processing is outlined in Algorithm 2, with functions associated with data caching indicated by the **VCDC** suffix. To begin, HyTGraph initializes the active vertices (Line 1) and the necessary data for cache maintenance (Line 2 and 3) and then performs the iterative computation (Line 4-24). Before engine selection and computation, HyTGraph first checks whether the cache condition is triggered based on the hotness (Line 5-7). If the condition is met, HyTGraph refreshes the vertex-centric cache by evicting cold data and reloading the new data in batches through GPU parallel processing. To avoid additional host-GPU data transfer, HyTGraph separates the cache filling and eviction, first evicting cold data and allocating space for the new candidates (in Line 8) and then loading the data to the GPU during HyTM processing. This design ensures low cache management overhead.

*Hotness-based cache candidate selection:* HyTGraph determines the cache candidates by sorting vertices by the hotness and obtaining the top-K vertices as the candidates. The "hotness"

of a vertex is defined by the access volume over a certain period, which is maintained using a  $|V|$ -length hotness array during iterative processing. In the candidate determination stage (Line 8), vertices are sorted by the hotness, with the topK hot vertices whose aggregated neighbor sizes do not exceed a given capacity ( $C_{cap}$ ) as the candidates, i.e.,  $V_{cand}$ . Following the selection, the hotness array is reset for the next step. This process, encompassing hotness evaluation and top-K computation, incurs less overhead compared to iterative processing. HyTGraph also performs the selection process on the GPU to achieve high performance.

*Cache refresh condition checking:* Determining appropriate replacement timing is crucial to efficient communication and low replacement overhead. Length replacement cause increased cache miss, while highly frequent replacements can result in the overhead outweighing the benefits. To address this issue, HyTGraph employs a lazy condition-checking approach to balance the cost and benefit. Firstly, instead of using iterations as the metric for timing checks, HyTGraph calculates the cache candidates  $V_{cand}$  after processing a given number of vertices, denoted by  $\gamma|V|$  (Line 5). This is to adapt to the asynchronous and  $\Delta$ -based priority processing in HyTGraph. Secondly, HyTGraph evaluates whether the data requiring replacement exceeds  $\beta$  of the new candidates (Line 7) and performs cache refreshing (Line 8) only if the condition is met. Otherwise, iterative computation continues with the old cache. The 30% threshold is configured based on the insight that it ensures sufficient changes in access frequency are accumulated while avoiding frequent cache refreshing caused by minor fluctuations in vertex hotness. In graph computation, the access frequency distribution typically stabilizes after the first few iterations. Therefore, this 30% setting allows HyTGraph to cache most of the correct vertices after a few refresh cycles. Once stabilized, the system avoids unnecessary adjustments triggered by small changes, which would otherwise require reorganizing the entire CSR structure.

*Vertex-centric cache refreshing:* Cache replacement involves the removal of old vertices and the loading of new ones. HyTGraph decouples them into two phases and processes them separately to reduce replacement overhead. In the cache updating stage (Line 8), HyTGraph performs GPU parallel cold data deletion and compacts the remaining data to make room for the new candidates (Line 8). Initially, a data deletion kernel is launched to identify all outdated vertices ( $V_{cache} \setminus V_{cand}$ ) and mark them as invalid. Subsequently, the remaining valid data ( $V_{cand} \cap V_{cache}$ ) is compacted in the cache. Finally, new spaces are allocated for the new vertices ( $V_{cand} \setminus V_{cache}$ ) in the tail of the cache, and the corresponding write indices are generated. To avoid additional data communication for loading data into the cache, HyTGraph integrates the cache loading stage into the computation kernel (Line 16). During computation, the GPU computation kernel will store the accessed data of new cache candidates ( $V_{cand} \setminus V_{cache}$ ) in the allocated cache space according to the indices generated in the cache updating stage. Such an implementation ensures cache refreshing efficiency by fully utilizing GPU parallelism and avoiding additional host-GPU data communications.

## VIII. IMPLEMENTATION AND OPTIMIZATIONS

*Flexible multi-stream scheduling:* The processing engines of ExpTM-F, ExpTM-C, and ImpTM-zero-copy-ZC require different resources, including CPUs for active edge compaction, GPU for the computation kernel, and PCIe for the host-GPU data transfer. To overlap the resource utilization and improve the parallelism, HyTGraph uses multiple CUDA streams to process the tasks concurrently. During the iterative processing, the task scheduler monitors the available streams and assigns them to tasks that have not been scheduled. The operating system will automatically overlap data transfer and kernel computation of different streams. HyTGraph first schedules the ExpTM-Filter tasks with specific priority (as discussed in Section VI) to leverage the contribution-driven priority scheduling. Then the ImpTM-zero-copy and ExpTM-compaction tasks are scheduled. The CPU-based active edge compaction can be overlapped with the kernel computation and data transfer of ImpTM-zero-copy and ExpTM-filter.

*Hotness computation:* HyTGraph uses a byte per vertex to record the access frequency during computation, optimizing GPU memory usage. The involved sorting, TopK, and compaction operations are implemented using the CUB library [8].

*CPU data compaction:* HyTGraph provides a simple yet efficient compaction engine on the CPU following the design of Subway [39]. To track locations of compacted edge data, HyTGraph re-generate a new compressed neighbor index for fast edge location.

*Computation kernel:* HyTGraph uses SEP-Graph's processing kernel for its mature optimizations and enables neighbor shifting [44] for it to support ExpTM-F and ExpTM-C engines.

## IX. EXPERIMENTAL EVALUATION

### A. Experimental Setup

*Environments:* Our main test platform is equipped with one Intel Silver 4210 2.20 Ghz 10-core CPU, 128 GB DRAM, and an NVIDIA GTX 2080Ti GPU with 34SMX clusters, 4352 cores, and 11 GB GDDR6 memory. The GPU is enabled with CUDA 10.1 runtime and 418.67 driver. The host side runs Ubuntu 18.04 with Linux kernel version 4.13.0. All source codes are compiled with -O3 optimization.

*Graph algorithms and datasets:* We evaluate HyTGraph with four algorithms. Besides SSSP and PageRank, the other two algorithms are Breadth-First Search (BFS) and Connect Component (CC) [44]. We use both real-world graphs and synthesized graphs in our evaluation. The major parameters of graph datasets that are used in our experiments are presented in Table III. The synthetic graphs are generated by PaRMAT [22] with the input parameters  $a = 0.5$ ,  $b = 0.2$ , and  $c = 0.2$  to maintain a power-law distribution.

*Systems for comparison:* We compare HyTGraph with three representative GPU-accelerated graph processing systems Subway [39], EMOGI [32], and Grus [45]; a in-memory GPU graph processing system cuGraph [9]; and a CPU-based graph processing system Galois [33]. Additionally, we also implement ExpTM-filter and ImpTM-unified-memory in HyTGraph's

TABLE III  
DATASET DESCRIPTION

Dataset	V	E	E / V	Size
Road-USA [2] (RU)	23.9M	57.7M	2.4	461M
Orkut [2] (RU)	3.1M	117M	39	968M
sk-2005 [3] (SK)	50.6M	1.93B	38	28GB
Twitter [2] (TW)	52.5M	1.96B	37	32GB
Friendster-konec [2] (FK)	68.3M	2.59B	37	42GB
uk-2007 [3] (UK)	105.1M	3.31B	31	55GB
Friendster-snap [2] (FS)	65.6M	3.61B	55	58GB
RMAT [22]	1-100M	0.1-6.4B	-	-

TABLE IV  
COMPARISON WITH OTHER SYSTEMS

		Overall runtime (s)						
Alg.	System	RU	OK	SK	TW	FK	UK	FS
PR	Galois	4.68	4.66	21.3	66.3	293.6	28.5	342.4
	cuGraph	0.31	0.5	-	-	-	-	-
	ExpTM-F	6.41	3.54	37.7	34.8	60.7	34.3	162.8
	ImpTM-UM	0.76	0.22	6.89	16.5	75.4	22.4	102.7
	Grus	0.60	<b>0.17</b>	<b>1.72</b>	12.2	52.2	14.8	79.8
	Subway	<b>0.17</b>	0.26	8.68	38.1	73.7	16.9	108.4
	EMOGI	3.11	0.52	18.6	21.4	51.1	12.4	68.3
	HyTGraph	1.67	0.34	1.84	<b>8.76</b>	<b>21.6</b>	<b>4.04</b>	<b>29.2</b>
SSSP	Galois	23.92	1.71	26.7	12.9	51.5	15.2	33.1
	cuGraph	8.05	0.89	-	-	-	-	-
	ExpTM-F	68.5	2.14	60.9	15.1	50.4	60.9	70.1
	ImpTM-UM	3.57	<b>0.18</b>	12.7	10.1	37.2	18.6	34.9
	Grus	14.98	0.20	25.2	11.2	70.8	5.32	16.9
	Subway	<b>2.65</b>	0.21	14.6	10.9	20.8	18.4	27.7
	EMOGI	33.75	0.38	7.46	4.09	14.9	4.71	11.8
	HyTGraph	5.24	0.24	<b>5.02</b>	<b>1.81</b>	<b>7.08</b>	<b>2.56</b>	<b>5.86</b>
CC	Galois	20.8	0.41	23.9	15.7	35.9	55.1	39.4
	cuGraph	<b>0.31</b>	0.22	-	-	-	-	-
	ExpTM-F	36.21	0.15	21.9	5.47	10.9	41.6	11.8
	ImpTM-UM	3.25	0.15	<b>1.43</b>	1.49	3.27	7.88	4.16
	Grus	12.37	<b>0.13</b>	2.09	1.36	3.21	5.17	4.69
	Subway	2.99	0.19	11.67	6.52	8.61	14.7	14.1
	EMOGI	14.95	0.21	4.01	1.96	2.71	4.54	3.76
	HyTGraph	5.25	0.18	3.23	<b>1.19</b>	<b>2.07</b>	<b>3.23</b>	<b>2.56</b>
BFS	Galois	<b>1.28</b>	0.88	16.2	7.55	12.5	15.2	14.7
	cuGraph	1.66	0.68	-	-	-	-	-
	ExpTM-F	26.87	1.02	20.3	3.86	8.87	25.1	9.54
	ImpTM-UM	2.45	0.45	1.13	1.29	1.97	2.33	6.25
	Grus	9.90	<b>0.13</b>	<b>0.83</b>	1.11	1.85	2.37	3.35
	Subway	1.63	0.28	7.39	5.79	6.85	9.04	13.49
	EMOGI	9.93	0.24	1.06	1.04	<b>1.44</b>	1.26	<b>1.97</b>
	HyTGraph	3.75	0.22	0.93	<b>0.85</b>	1.82	<b>0.88</b>	2.54

codebase for a fair comparison. We use the default configuration of these systems. The number of compaction threads in Subway and HyTGraph is set to  $2 \times$  CPU cores. The cache capacity of HyTGraph is set to 2 GB, which is the available memory when processing the largest graph and aligned for all datasets. All reported results are measured by averaging the number of 5 runs. HyTGraph is open-sourced at <https://github.com/iDC-NEU/HyTGraph>.

### B. Overall Performance on Small Graphs

We compare HyTGraph with in-memory processing frameworks on small graphs, using both a power-law graph (Orkut)

and a uniform road network (Road-USA), to analyze the performance trade-off between out-of-memory and in-memory processing, as shown in Table IV. In addition to cuGraph, UVM-based Grus and ImpTM-UM degrade to in-memory processing as most data can be cached in GPUs. Subway also switches to in-memory processing by copying the data entirely to GPU. In contrast, we disable the graph caching function in HyTGraph to evaluate its performance as a fully out-of-memory framework. Our findings reveal that Grus, Subway, ImpTM-UM, and cuGraph achieve the best results across various configurations on these small graphs. On the Road-USA graph, the disadvantages of out-of-memory processing become more evident due to its disproportionately long absolute runtime relative to the graph size. This is primarily attributed to the small average degree and large diameter of the road network, which result in a higher number of iterations. Algorithms such as CC and SSSP require 5,000–6,000 iterations to converge, each involving substantial fragmented edge data accesses (each access involves up to 9 neighborhoods per). EMOGI with zero-copy access exhibits inferior performance in this scenario, mainly due to PCIe underutilization caused by transferring sliced data. In addition, the performance on the road network is also affected by load balancing. We observe that in-memory processing frameworks with fine-grained workload balancing, such as cuGraph and Subway (which leverages Tigr’s degree-optimized load balancing [38]), outperform other frameworks in different cases. Notably, Galois outperforms GPU baselines for the BFS algorithm on the Road-USA. This demonstrates that CPU-based frameworks still hold advantages in certain scenarios.

### C. Overall Performance on Large Graphs

*Comparison with ExpTM-F, Subway, and EMOGI:* Table IV shows the overall results. Due to the heavy redundant transfer, ExpTM-F shows inferior performance. The speedup of HyTGraph over ExpTM-F ranges from 2.81X (for PageRank on FK) to 28.52X (for BFS on U.K.) with an average of 9.89X. Neither Subway nor EMOGI is always better than the other. The speedup of HyTGraph over Subway ranges from 2.91X (for SSSP on SK) to 10.27X (for BFS on U.K.) with an average of 5.02X. Subway’s critical performance bottleneck lies in its heavy CPU-based compaction and preprocessing (For SSSP algorithm, the preprocessing and compaction overhead account for 46.9%–74.9% of the total runtime). On CC, SSSP, and PageRank, HyTGraph is faster than EMOGI by 2.01X on average, with its speedups ranging from 1.24X to 7.91X. With the help of zero-copy access, EMOGI achieves significant performance improvement on low-activeness subgraphs. While for the high-activeness subgraphs, especially those with dense and small degree vertices, EMOGI usually has low host-GPU utilization due to unsaturated memory requests. In contrast, HyTGraph achieves efficient data transfer on both high-activeness and low-activeness partitions by adopting hybrid transfer management. On BFS, HyTGraph outperforms Subway and EMOGI on SK, TW, and U.K.. On FK and FS, EMOGI shows better performance because most of the accesses on these two graphs are sparse. Moreover, compared

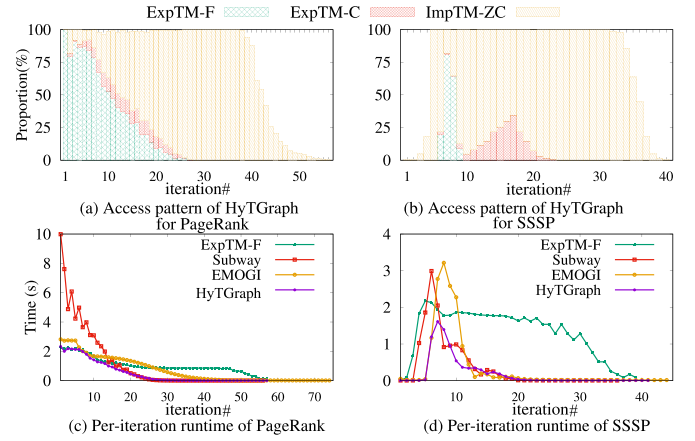


Fig. 5. Execution path of HyTGraph and per-iteration runtime comparison with ExpTM-filter, EMOGI and Subway (on FK).

with HyTGraph, EMOGI avoids the cost analysis, engine selection, and task combination. Since BFS traverses each vertex only once and thus has no cross-layer communication, we disable the vertex-centric data caching for it.

*Comparison with Unified-Memory (UM)-based systems: ImpTM-UM and Grus:* On the SK graph, the UM-based frameworks demonstrate superior performance for PageRank, CC, and BFS algorithms because the accessed edge-associated data can be entirely cached in the GPU memory. UM-based approaches only transfer the data once. However, when processing large graphs, the performance of ImpTM-UM degrades significantly because the implicit data transfer requires expensive page replacement and data transfer overhead. The experimental results show that on the four large graphs, HyTGraph achieves on average 3.01X and 2.55X speedups over ImpTM-UM and Grus, respectively.

*Comparison with CPU-based Approach:* From Table IV, we can observe that the GPU-accelerated graph processing frameworks show significant performance improvement over CPU-based Galois. Specifically, HyTGraph shows an average of average 10.34x speedup over Galois.

### D. Execution Path Analysis of HyTM

To demonstrate the performance improvement of hybrid processing, we record the execution path of HyTGraph on PageRank and SSSP to show the proportion of partitions using ExpTM-filter, ExpTM-compaction, and ImpTM-zero-copy in each iteration. Fig. 5(a) shows the result on PageRank. The proportion of active partitions is high in the early iterations, HyTGraph prefers ExpTM-filter. As the algorithm converges and many vertices become inactive, the proportion of ImpTM-zero-copy increases. For SSSP in Fig. 5(b), there are few active vertices in the early and last few iterations, HyTGraph prefers ImpTM-zero-copy. When most vertices are activated in the middle iterations, HyTGraph prefers ExpTM-filter to improve the transfer efficiency. As the number of active vertex decreases, ExpTM-compaction is also used in some partitions. Fig. 5(c) and (d) show the



TABLE V  
TRANSFER REDUCTION ANALYSIS

Transfer volume / Edge volume						
Alg.	System	SK	TW	FK	UK	FS
PR	ExpTM-F	57.6X	52.4X	58.3X	30.9X	121.6X
	Subway	2.46X	<b>5.48X</b>	10.74X	1.79X	12.44X
	EMOGI	3.31X	20.6X	24.6X	3.81X	25.23X
	HytGraph <sup>-</sup>	2.17X	10.9X	12.01X	1.68X	12.62X
	HyTGraph	<b>1.36X</b>	8.15X	<b>8.82X</b>	<b>1.38X</b>	<b>10.8X</b>
SSSP	ExpTM-F	44.3X	11.2X	28.1X	24.3X	24.1X
	Subway	4.23X	2.07X	<b>3.32X</b>	1.78X	3.19X
	EMOGI	3.29X	1.74X	4.81X	1.11X	2.69X
	HytGraph <sup>-</sup>	3.25X	1.25X	4.60X	1.13X	2.52X
	HyTGraph	<b>2.53X</b>	<b>1.10X</b>	3.87X	<b>1.06X</b>	<b>2.12X</b>

per-iteration runtime of ExpTM-F, Subway, EMOGI, and HyTGraph. As these systems adopt different asynchronous processing strategies, the active vertex number of different systems in each iteration is not the same. HyTGraph cannot consistently outperform the others in each iteration. However, through the hybrid transfer management, HyTGraph achieves the minimum overall runtime.

#### E. Transfer Reduction Analysis

We analyze the effectiveness of HyTGraph's transfer reduction by comparing it with ExpTM-filter, Subway (ExpTM-compaction), and EMOGI (ImpTM-zero-copy), using PageRank and SSSP algorithms across five real-world graphs. Transfer volumes are normalized against the volume of edges. As shown in Table V, ExpTM-filter exhibits the highest transfer volume. With the help of fine-grained zero-copy access, EMOGI achieves considerable transfer reduction. However, due to the lack of effective asynchronous scheduling, the transfer volume is still large. Subway, aimed at minimizing data transfer through CPU data compaction, shows diverse performance across different algorithms due to the unstable naive asynchronous processing. Specifically, Subway excels in the PageRank algorithm, where additional computation markedly enhances convergence. However, for the value-replacement-based SSSP algorithm, Subway loses its edge due to the stale computation problem [11]. Since processing the transferred subgraph only twice each scheduling, HyTGraph's variant without data caching (HytGraph<sup>-</sup>) shows marginal improvement over Subway for the PageRank algorithm. On small graphs with few partitions, e.g., the TW graph, HytGraph<sup>-</sup> even demands 2X data transfer compared to Subway. Conversely, for the SSSP algorithm, HytGraph<sup>-</sup> achieves significant transfer reduction in most scenarios through hybrid transfer management and contribution-driven priority scheduling. After incorporating vertex-centric graph caching, HyTGraph further diminishes data transmission by 7% to 37% over HytGraph<sup>-</sup>. These enhancements significantly improve HyTGraph's communication efficiency.

#### F. Performance Improvement Analysis

To access the performance gain offered by different optimizations, our methodology initiates from a fundamental baseline

utilizing ImpTM-zero-copy (EMOGI) and sequentially incorporates hybrid transfer management, vertex-centric graph caching, task combining, and contribution-driven scheduling. We selected ImpTM-zero-copy as the initial point of comparison due to its superior performance across all baselines. Fig. 6 shows the normalized runtime, reflecting the performance gains attributed to each optimization.

*Effectiveness of HyTM:* The hybrid transfer management method brings speedups of averaging 2.87X, 1.66X, 1.34X, and 1.07X for PageRank, SSSP, CC, and BFS algorithms, respectively. For algorithms with significant variations in active vertices, such as SSSP, CC, and PageRank, HyTM demonstrates notable performance improvements. Conversely, for algorithms with a few active vertices throughout the execution, such as BFS, ImpTM-zero-copy still holds advantages. This is because it benefits from the efficient zero-copy access and avoids the overhead associated with cost analysis and task management in HyTM. On FK and FS graph, ImpTM-zero-copy even demonstrates better performance than HyTM. However, in most algorithms with diverse access patterns, HyTGraph demonstrates considerable effectiveness.

*Effectiveness of task combining and contribution-driven scheduling:* The task combining (TC) can bring speedups of averaging 1.28X, 1.37X, 1.19X, and 1.05X for PageRank, SSSP, CC, and BFS algorithms, respectively. The contribution-driven scheduling (CDS) further brings speedups of averaging 2.18X, 1.21X, 1.25X, and 1.06X over the hybrid processing with TC. Overall, integrating the two optimizations brings speedups of averaging 2.78X, 1.67X, 1.47X, and 1.16X over the naive hybrid transfer management method on the four algorithms. PageRank algorithm benefits most because the proposed asynchronous processing can effectively accelerate convergence by prioritizing the vertices with large contributions, i.e., the vertex value. In contrast, BFS benefits little from the two designs because each vertex is activated only once during the iterative processing, leading to a small overall transfer overhead.

*Effectiveness of vertex-centric graph caching (VCGC):* Since the BFS algorithm traverses each vertex only once, it can not benefit from the VCGC optimization. Therefore, we do not enable VCGC for it in the evaluation. On the other three algorithms, VCGC provides speedups ranging from 0.98X to 1.38X by reducing cross-layer communication. We observe that the effectiveness of data caching differs across algorithms. For PageRank and SSSP, VCGC exhibits consistent improvements ranging from 1.13X to 1.38X. In contrast, for the CC algorithm, VCGC demonstrates slightly weaker performance on TW and FK graphs. This discrepancy arises because the runtime of these cases has short runtimes and low volumes of cross-layer communication, rendering the caching improvements insufficient to outweigh the costs associated with cache replacement. Nevertheless, as an enhancement to the hybrid transfer management method, VCGC effectively reduces transmission volumes in most scenarios. As a recommendation, VCGC is particularly advantageous for algorithms with extensive cross-layer accesses, such as PageRank and SSSP and their variants.

*Comparison with hub-vertex caching:* We evaluate the performance of the proposed VCGC and the hub-vertex-based

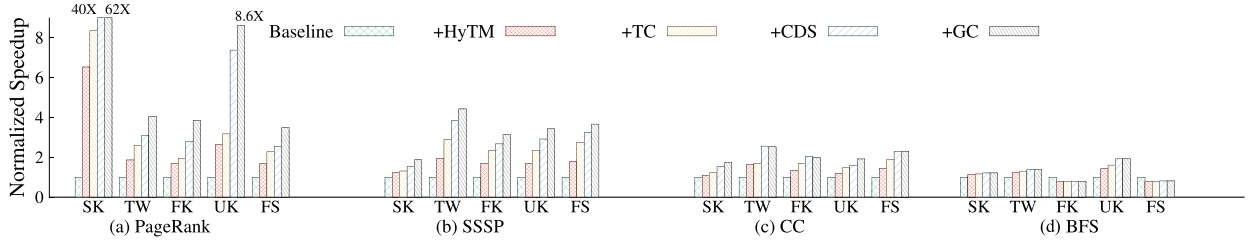


Fig. 6. Performance gain analysis of hybrid transfer management (HyTM), Task Combining (TC), Contribution-Driven Scheduling (CDS), and vertex-centric Graph Caching (GC).

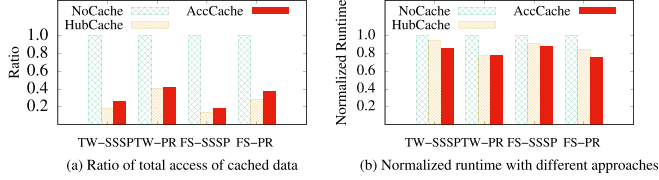


Fig. 7. Comparison with hub-vertex caching.

TABLE VI  
RUNTIME BREAKDOWN OF HYTGRAPH WITH VERTEX-CENTRIC GRAPH CACHING ON PAGERANK

Engine	Component	time of different components (s)				
		SK	TW	FK	UK	FS
W/o caching	overall	2.85	11.5	30.2	4.71	40.9
	Computation	1.84	8.76	21.6	4.04	29.2
W/ caching	Cand. Selection	+0.07	+0.08	+0.1	+0.07	+0.09
	Cache Refreshing	+0.19	+0.48	+0.68	+0.44	+0.52
Benefit	-	-1.26	-3.29	-9.26	-1.17	-12.2

caching approach by comparing: 1) the proportion of reduced edge accesses relative to the total accesses, and 2) the normalized runtime compared to HyTGraph without data caching, as shown in Fig. 7. Both caching mechanisms effectively reduce the volume of edge accesses. On average, VCGC achieves a 24.2% higher reduction in edge accesses compared to hub-vertex sorting. However, in specific scenarios such as TW-PR, hub-vertex sorting performs comparably to VCGC. Specifically, VCGC and hub-vertex caching reduce edge accesses by 14.1%–40.9% (avg. 25.1%) and 17.3%–42.4% (avg. 31.0%), respectively, compared to HyTGraph without data caching. In terms of runtime improvement, the difference between the two approaches is less pronounced. VCGC achieves an average improvement of approximately 7.9% over hub-vertex caching. This is due to HyTGraph’s contribution-driven priority scheduling, which performs additional iterations for processing hub vertices, inherently reducing remote data accesses. This optimization overlaps with the caching mechanism, preventing the access reductions from fully translating into runtime gains. Looking forward, VCGC’s access-frequency-based caching mechanism can be extended to a wider range of graph processing tasks, such as real-time streaming graph analysis [15], graph sampling [21], and multi-hop subgraph analysis [26]. It is expected to achieve high runtime efficiency in these applications.

*Overhead of graph caching:* Table VI presents the overhead associated with graph caching, including the hotness-based candidate selection and Vertex-centric cache refreshing. We also

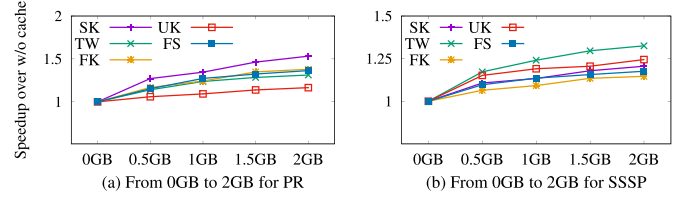


Fig. 8. Performance with varying cache sizes.

compare the results with the performance gains brought by graph caching on PageRank. We observe that the time dedicated to candidate selection and cache refreshing is small, accounting for approximately 3% to 12% of total runtime. Compared to the benefit of optimizing cross-layer accesses, the total overhead of cache management is also minor, ranging from 5% to 42% of performance gain. With an appropriate cache refresh configuration, HyTGraph typically caches the most correct vertices after a single (or two for TW and FK graph) cache refreshing cycle. In summary, for graph iterative algorithms with extensive cross-iteration repeat vertex accesses (e.g., PageRank), GPU graph caching can significantly enhance performance, with the incurred management overhead being effectively outweighed by the benefits.

### G. Sensitivity Analysis

*Varying cache sizes:* To evaluate the impact of varying cache sizes, we run SSSP and PR on the five large graphs, starting from no hot data caching (0 GB) and linearly increasing the cache size to 2 GB. As shown in Fig. 8. We observe that the first 0.5B of data caching yields a significant performance improvement, approximately 42% to 62% compared to the improvement of caching 2 GB of data. This is because, due to the power-law distribution, frequently accessed vertices represent a small proportion of the total but account for a large number of edge accesses. As the cache size increases from 0.5 GB to 2 GB, the per-GB benefit of caching graph data decreases. Nevertheless, the overall performance improvement remains non-negligible.

*Varying graph sizes:* We compare HyTGraph with Grus, Subway, and EMOGI under variable graph sizes and report the results in Fig. 9. Subway encounters an integer overflow issue and fails to process graphs with 6.4 billion edges. Grus exhibits superior performance on small graphs because the data only needs to be loaded once. However, as graph sizes increase, the performance declines because the overhead of data

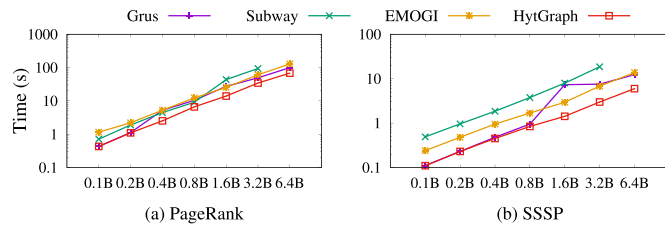


Fig. 9. Performance comparison with increasing graph size, the graphs are generated by RMAT with sizes from 0.1 Billion to 6.4 Billion (64X).

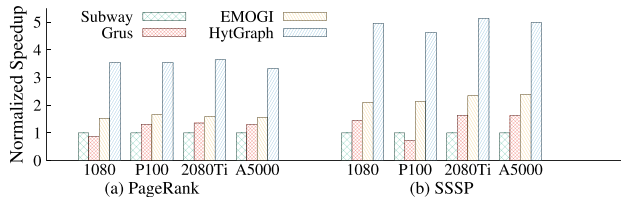


Fig. 10. Performance comparison on different GPUs (FS).

migration of unified memory increases. We observe that most system demonstrates linear scaling with increasing graph sizes. Benefiting from the vertex-centric GPU data caching, HyTGraph achieves comparable performance with Grus and outperforms other systems on small graphs. On large graphs, Grus underperforms on the SSSP algorithm due to the increased overhead of frequent memory page swapping of low active data. In contrast, HyTGraph can efficiently process large graphs through a combination of hybrid transfer management and vertex-centric data caching. As the graph size increases from 0.1B to 6.4B (64X), the runtime of Grus, EMOGI, and HyTGraph for PageRank increases by 231.2X, 111.6X, and 157.86X, respectively. For the SSSP algorithm, the runtime of Grus, EMOGI, and HyTGraph increases by 111.8X, 57.08X, and 54.22X, respectively. Due to limited cache capacity, the performance benefits of HyTGraph's caching do not scale proportionally as the graph size increases, which makes it appear to have weaker scalability compared to EMOGI without caching. However, HyTGraph still demonstrates a significant advantage in absolute runtime. Furthermore, compared to the UVM-based page caching approach used in Grus, HyTGraph's caching method offers superior scalability.

*Varying GPUs:* We assess HyTGraph's performance across various platforms equipped with different GPUs, including a GTX 1080 connected to one i7-8700 k CPU through PCIe3.0x8, a NVIDIA P100 connected to dual Xeon-silver 4210 CPUs through PCIe3.0x16, a GTX 2080Ti connected to dual Xeon Silver 4210 CPUs through PCIe3.0x16, and an NVIDIA A5000 connected to dual Xeon silver 4316 CPUs through PCIe4.0x16. The FS graph serves as the input data. We normalize the runtime of all systems to Subway and present the results in Fig. 10. We can observe that HyTGraph outperforms all three competitors across all platforms. For PageRank (and SSSP), HyTGraph achieves speedups of 3.4X-3.6X (4.6X-5.1X), 2.6X-4.0X (3.1X-6.3X), and 2.1-2.2X (2.1X-2.3X) over Subway, Grus, and EMOGI, respectively.

## X. LIMITATIONS AND FUTURE WORK

*Extending HyTGraph to multiple GPU acceleration:* Currently, HyTGraph assumes the vertex data can fit into a single GPU. When processing graphs with vertex data that exceed a single GPU's memory capacity, the data can be partitioned across multiple GPUs. This approach introduces additional challenges, including managing host-GPU communication bandwidth, inter-GPU communication overhead, and achieving a balance among computation, communication, and cache utilization. We take developing multi-GPU-based HyTGraph as future work.

*Adapting to GPU platforms with advanced interconnects:* Recently, the hardware manufacturers have introduced advanced interconnects, such as NVIDIA NVlink [35], Intel CXL [10], and AMD Infinity Fabric [1]). These technologies facilitate the construction of larger host memory pools for scaling large graphs and enable devices and host memory to connect through faster and more heterogeneous communication links. While these advancements offer new opportunities for processing large-scale graphs, they also present challenges in optimizing irregular data transfers across complex communication links. Extending HyTGraph to support these emerging interconnect technologies is part of our future work.

## XI. RELATED WORK

*In-GPU-memory graph processing:* The high parallelism of GPU has attracted great attention [12], [19], [20], [31], [48], [51], [53] in graph processing community. Cusha [23] uses two novel data structures, named GShards and CW, to avoid non-coalesced memory access. Gunrock [48] performs computation on the frontier with data-centric abstraction. Tigr [38] proposes a virtual transformation to transform skewed graphs into virtual vertices for load-balancing. SEP-Graph [44] optimize execution paths by adaptively switching Sync/Async, Push/Pull, and data-driven/topology-driven modes.

*Out-of-core GPU graph processing:* GPU-accelerated graph processing has attracted extensive attention. Besides the systems mentioned above [13], [18], [32], [39], [40], [42], [45], [52], recent studies also propose CPU-GPU co-processing to accelerate large graphs computation [14], [29]. However, the CPU-based low-activeness subgraph processing may become a new bottleneck. Besides graph processing, researchers have also investigated GPU-accelerated pattern matching on large graphs [7], [17] that optimize communication by sharing execution or combining ZC and UVM.

*Cross-iteration graph reusing in heterogeneous memory systems:* Halo, Ascetic, and Liberator optimize cross-iteration graph reuse on GPUs. Halo [13] leverages UVM for automated data migration and caching, introducing Harmonic Locality Ordering to reorganize the graph and improve cache efficiency. Ascetic [42] and Liberator [27] split data into static and dynamic cache regions to enhance memory utilization. GraphMP [41] employs a compressed edge cache mechanism to maximize cache efficiency. However, all these systems suffer from the access amplification issue, due to coarse-grained page-level memory transfers.



## XII. CONCLUSION

We present HyTGraph, a highly efficient GPU-accelerated graph processing framework by adaptively switching transfer management methods involving explicit transfer management and implicit transfer management. This hybrid approach maximizes the host-GPU bandwidth and is necessary to achieve the shortest overall execution time. Moreover, HyTGraph provides a vertex-centric graph caching method that further reduces communication through data transfer reusing. Our intensive experiments show the high effectiveness of HyTGraph.

## REFERENCES

- [1] AMD infinity fabric link, 2024. Accessed: Oct. 11, 2024. [Online]. Available: <https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/other/56978.pdf>
- [2] "The Koblenz network collection," 2021. Accessed: Sep. 1, 2021. <http://konect.uni-koblenz.de/>
- [3] "Laboratory for web algorithmics," 2021. Accessed: Sep. 2021. [Online]. Available: <http://law.di.unimi.it/>
- [4] N. Agarwal, D. W. Nellans, M. Stephenson, M. O'Connor, and S. W. Keckler, "Page placement strategies for GPUs within heterogeneous memory systems," in *Proc. 20th Int. Conf. Architectural Support Program. Lang. Operating Syst.*, Istanbul, Turkey, 2015, pp. 607–618.
- [5] R. Ausavarungnirun et al., "Mosaic: A GPU memory manager with application-transparent support for multiple page sizes," in *Proc. 50th Annu. IEEE/ACM Int. Symp. Microarchitecture*, Cambridge, MA, USA, 2017, pp. 136–150.
- [6] T. Ben-Nun, M. Sutton, S. Pai, and K. Pingali, "Groute: An asynchronous multi-GPU programming model for irregular computations," in *Proc. 22nd ACM SIGPLAN Symp. Princ. Pract. Parallel Program.*, Austin, TX, USA, 2017, pp. 235–248.
- [7] J. Chen, Q. Wang, Y. Gu, C. Li, and G. Yu, "Unified-memory-based hybrid processing for partition-oriented subgraph matching on GPU," *World Wide Web*, vol. 25, no. 3, pp. 1377–1402, 2022.
- [8] "NVIDIA cub," 2023. [Online]. Available: <https://github.com/NVIDIA/cub>
- [9] "Rapids cudagraph," 2023. [Online]. Available: <https://docs.rapids.ai/api/cudagraph/stable/>
- [10] CXL, "Compute express link specification revision 1.1," 2022. [Online]. Available: <https://www.computeexpresslink.org/>
- [11] W. Fan et al., "Adaptive asynchronous parallelization of graph algorithms," in *Proc. Int. Conf. Manage. Data*, 2018, Houston, TX, USA, 2018, pp. 1141–1156.
- [12] A. Gaihare, Z. Wu, F. Yao, and H. Liu, "XBFS: Exploring runtime optimizations for breadth-first search on GPUs," in *Proc. 28th Int. Symp. High-Perform. Parallel Distrib. Comput.*, Phoenix, AZ, USA, 2019, pp. 121–131.
- [13] P. Gera, H. Kim, P. Sao, H. Kim, and D. A. Bader, "Traversing large graphs on GPUs with unified memory," *Proc. VLDB Endow.*, vol. 13, no. 7, pp. 1119–1133, 2020.
- [14] A. Gharaibeh, L. B. Costa, E. Santos-Neto, and M. Ripeanu, "A yoke of oxen and a thousand chickens for heavy lifting graph processing," in *Proc. Int. Conf. Parallel Architectures Compilation Techn.*, Minneapolis, MN, USA, 2012, pp. 345–354.
- [15] S. Gong et al., "Automating incremental graph processing with flexible memoization," *Proc. VLDB Endow.*, vol. 14, no. 9, pp. 1613–1625, 2021.
- [16] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin, "Powergraph: Distributed graph-parallel computation on natural graphs," in *Proc. 10th USENIX Symp. Operating Syst. Des. Implementation*, Hollywood, CA, USA, 2012, pp. 17–30.
- [17] W. Guo, Y. Li, M. Sha, B. He, X. Xiao, and K. Tan, "GPU-accelerated subgraph enumeration on partitioned graphs," in *Proc. 2020 Int. Conf. Manage. Data*, Portland, OR, USA 2020, pp. 1067–1082.
- [18] W. Han, D. Mawhirter, B. Wu, and M. Buland, "Graphic: Large-scale asynchronous graph traversals on just a GPU," in *Proc. 2019 USENIX Annu. Tech. Conf.*, Renton, WA, USA, 2019, pp. 429–442.
- [19] P. Harish and P. J. Narayanan, "Accelerating large graph algorithms on the GPU using CUDA," in *Proc. 14th Int. Conf. High Perform. Comput.*, Goa, India, 2007, pp. 197–208.
- [20] S. Hong, S. K. Kim, T. Oguntebi, and K. Olukotun, "Accelerating CUDA graph algorithms at maximum warp," in *Proc. 16th ACM SIGPLAN Symp. Princ. Pract. Parallel Program.*, San Antonio, TX, USA, 2011, pp. 267–276.
- [21] H. Hu, F. Liu, Q. Pei, Y. Yuan, Z. Xu, and L. Wang, "lambda-grapher: A resource-efficient serverless system for GNN serving through graph sharing," in *Proc. ACM Web Conf.*, Singapore 2024, pp. 2826–2835.
- [22] F. Khorasani, R. Gupta, and L. N. Bhuyan, "Scalable SIMD-efficient graph processing on GPUs," in *Proc. 24th Int. Conf. Parallel Architectures Compilation Techn.*, 2015, pp. 39–50.
- [23] F. Khorasani, K. Vora, R. Gupta, and L. N. Bhuyan, "CuSha: Vertex-centric graph processing on GPUs," in *Proc. 23rd Int. Symp. High-Perform. Parallel Distrib. Comput.*, 2014, pp. 239–252.
- [24] M. Kim, K. An, H. Park, H. Seo, and J. Kim, "GTS: A fast and scalable graph processing method based on streaming topology to gpus," in *Proc. Int. Conf. Manage. Data*, San Francisco, CA, USA, 2016, pp. 447–461.
- [25] A. Kyrola, G. E. Blelloch, and C. Guestrin, "Graphchi: Large-scale graph computation on just a PC," in *Proc. 10th USENIX Symp. Operating Syst. Des. Implementation*, Hollywood, CA, USA, 2012, pp. 31–46.
- [26] C. Li et al., "Bytograph: A high-performance distributed graph database in bytedance," *Proc. VLDB Endowment*, vol. 15, no. 12, pp. 3306–3318, 2022.
- [27] S. Li et al., "Liberator: A data reuse framework for out-of-memory graph computing on GPUs," *IEEE Trans. Parallel Distrib. Syst.*, vol. 34, no. 6, pp. 1954–1967, Jun. 2023.
- [28] H. Liu and H. H. Huang, "Enterprise: Breadth-first graph traversal on GPUs," in *Proc. Int. Conf. High Perform. Comput. Netw. Storage Anal.*, Austin, TX, USA, pp. 68:1–68:12, 2015.
- [29] L. Ma, Z. Yang, H. Chen, J. Xue, and Y. Dai, "Garaph: Efficient GPU-accelerated graph processing on a single machine with balanced replication," in *Proc. 2017 USENIX Annu. Tech. Conf.*, Santa Clara, CA, USA, 2017, pp. 195–207.
- [30] G. Malewicz et al., "Pregel: A system for large-scale graph processing," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, Indianapolis, Indiana, USA, 2010, pp. 135–146.
- [31] D. Merrill, M. Garland, and A. S. Grimshaw, "High-performance and scalable GPU graph traversal," *ACM Trans. Parallel Comput.*, vol. 1, no. 2, pp. 14:1–14:30, 2015.
- [32] S. Min, V. S. Maitlody, Z. Qureshi, J. Xiong, E. Ebrahimi, and W. Hwu, "EMOGI: Efficient memory-access for out-of-memory graph-traversal in gpus," *Proc. VLDB Endow.*, vol. 14, no. 2, pp. 114–127, 2020.
- [33] D. Nguyen, A. Lenharth, and K. Pingali, "A lightweight infrastructure for graph analytics," in *Proc. 24th Symp. Operating Syst. Princ.*, Farmington, PA, USA, 2013, pp. 456–471.
- [34] NVIDIA, "NVIDIA A100 tensor core GPU," 2022. [Online]. Available: <https://www.nvidia.com/en-us/data-center/a100/>
- [35] NVIDIA, "NVIDIA H100 tensor core GPU," 2022. [Online]. Available: <https://www.nvidia.com/en-us/data-center/h100/>
- [36] NVIDIA, "NVIDIA Tesla P100," 2022. [Online]. Available: <https://www.nvidia.com/en-us/data-center/tesla-p100/>
- [37] A. Roy, I. Mihailovic, and W. Zwaenepoel, "X-stream: Edge-centric graph processing using streaming partitions," in *Proc. ACM SIGOPS 24th Symp. Operating Syst. Princ.*, Farmington, PA, USA, 2013, pp. 472–488.
- [38] A. H. N. Sabet, J. Qiu, and Z. Zhao, "Tigr: Transforming irregular graphs for GPU-friendly graph processing," in *Proc. 23rd Int. Conf. Architectural Support Program. Lang. Operating Syst.*, Williamsburg, VA, USA 2018, pp. 622–636.
- [39] A. H. N. Sabet, Z. Zhao, and R. Gupta, "Subway: Minimizing data transfer during out-of-GPU-memory graph processing," in *Proc. 15th EuroSys Conf.*, Heraklion, Greece, 2020, pp. 12:1–12:16.
- [40] D. Sengupta, S. L. Song, K. Agarwal, and K. Schwan, "Graphreduce: Processing large-scale graphs on accelerator-based systems," in *Proc. Int. Conf. High Perform. Comput. Netw. Storage Anal.*, Austin, TX, USA, 2015, pp. 28:1–28:12, 2015.
- [41] P. Sun, Y. Wen, T. N. B. Duong, and X. Xiao, "GraphMP: I/O-efficient big graph analytics on a single commodity machine," *IEEE Trans. Big Data*, vol. 6, no. 4, pp. 816–829, Dec. 2020.
- [42] R. Tang et al., "Ascetic: Enhancing cross-iterations data efficiency in out-of-memory graph processing on GPUs," in *Proc. 50th Int. Conf. Parallel Process.*, Lemont, IL, USA, 2021, pp. 41:1–41:10.
- [43] K. Vora, "LUMOS: Dependency-driven disk-based graph processing," in *Proc. USENIX Conf. Usenix Annu. Tech. Conf.*, 2019, pp. 429–442.

- [44] H. Wang, L. Geng, R. Lee, K. Hou, Y. Zhang, and X. Zhang, "Sep-graph: Finding shortest execution paths for graph processing under a hybrid framework on GPU," in *Proc. 24th ACM SIGPLAN Symp. Princ. Pract. Parallel Program.*, Washington, DC, USA, 2019, pp. 38–52.
- [45] P. Wang, J. Wang, C. Li, J. Wang, H. Zhu, and M. Guo, "Grus: Toward unified-memory-efficient high-performance graph processing on GPU," *ACM Trans. Archit. Code Optim.*, vol. 18, no. 2, pp. 22:1–22: 25, 2021.
- [46] Q. Wang, X. Ai, Y. Zhang, J. Chen, and G. Yu, "Hytgraph: GPU-accelerated graph processing with hybrid transfer management," in *Proc. 39th IEEE Int. Conf. Data Eng.*, Anaheim, CA, USA, 2023, pp. 558–571.
- [47] Q. Wang et al., "Automating incremental and asynchronous evaluation for recursive aggregate data processing," in *Proc. 2020 Int. Conf. Manage. Data*, SIGMOD 2020, Portland, OR, USA, June 14–19, 2020, pp. 2439–2454.
- [48] Y. Wang, A. A. Davidson, Y. Pan, Y. Wu, A. Riffel, and J. D. Owens, "Gunrock: A high-performance graph processing library on the GPU," in *Proc. 23rd Int. Symp. High-Perform. Parallel Distrib. Comput.*, Vancouver, BC, Canada, 2014, pp. 239–252.
- [49] Y. Zhang, Q. Gao, L. Gao, and C. Wang, "Maiter: An asynchronous graph processing framework for delta-based accumulative iterative computation," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 8, pp. 2091–2100, Aug. 2014.
- [50] Y. Zhang, V. Kiriansky, C. Mendis, S. P. Amarasinghe, and M. Zaharia, "Making caches work for graph analytics," in *Proc. IEEE Int. Conf. Big Data*, Boston, MA, USA, 2017, pp. 293–302.
- [51] Y. Zhang, X. Liao, H. Jin, B. He, H. Liu, and L. Gu, "Digraph: An efficient path-based iterative directed graph processing system on multiple gpus," in *Proc. 24th Int. Conf. Architectural Support Program. Lang. Operating Syst.*, Providence, RI, USA 2019, pp. 601–614.
- [52] L. Zheng et al., "Scaph: Scalable GPU-accelerated graph processing with value-driven differential scheduling," in *Proc. 2020 USENIX Annu. Tech. Conf.*, 2020, pp. 573–588.
- [53] J. Zhong and B. He, "Medusa: Simplified graph processing on GPUs," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 6, pp. 1543–1552, Aug. 2014.
- [54] X. Zhu, W. Han, and W. Chen, "Gridgraph: Large-scale graph processing on a single machine using 2-level hierarchical partitioning," in *Proc. 2015 USENIX Annu. Tech. Conf.*, Santa Clara, 2015, pp. 375–386.



**Qiange Wang** received the PhD degree in computer science from Northeastern University, China, in 2022. He is currently a postdoctoral research fellow with the National University of Singapore. His research interests include distributed graph processing, learning, and management systems.



**Xin Ai** is currently working toward the PhD degree in computer science with Northeastern University. His research interests include parallel and distributed graph computing and learning system.



**Yongze Yan** received the master's degree in computer science from the Northeastern University of China, Shenyang, in 2023. He is currently working toward the PhD degree. His major research interests include GPU data processing and database on emerging hardware.



**Shufeng Gong** received the PhD degree in computer science from Northeastern University, China, in 2021. He is currently a lecturer with Northeastern University, China. His research interests include cloud computing, distributed graph processing, and data mining.



**Yanfeng Zhang** received the PhD degree in computer science from Northeastern University, China, in 2012. He is currently a professor with Northeastern University, China. His research consists of distributed systems and Big Data processing. He has published many papers in the above areas. His paper in SoCC 2011 was honored with "Paper of Distinction".



**Jing Chen** received the master's degree in computer science from Northeastern University, China, in 2022. Her research interest includes GPU graph processing systems.



**Ge Yu** (Senior Member, IEEE) received the PhD degree in computer science from the Kyushu University of Japan, in 1996. He is now a professor with Northeastern University, China. His current research interests include distributed and parallel systems, cloud computing, Big Data management, and blockchain techniques and systems. He has published more than 200 papers in refereed journals and conferences. He is the CCF fellow and the ACM senior member.