# Performance Optimization of GPU Breadth-First Search Algorithm

July 21, 2025

**Abstract**

In recent times, the use of large-scale graphs involving millions of vertices are quite common in various practical applications like social networks, search engines, maps, etc., and often take up high computational resources when processing such graphs for traversal or finding shortest path between nodes. Although practical sequential implementations are possible using high-end computers and have been reported, such resources are not easily accessible. Moreover, the performance of such algorithms degrades drastically with increasing graph size. To this end, possiblity of parallelisation of graph algorithms using Graphical Processing Units (GPU) have gained significant importance, considering their high computational power and affordability, though their restrictive programming model is often difficult to get past. In this paper, we present the parallelisation of three fundamental class of graph problem - Breadth First Search, Single Source Shortest Path, and All-Pair Shortest Path, using CUDA based implementation on GPUs. We have also profiled the performance of the algorithms on a diverse set of generated and procured large graphs and have observed significant speedups as compared to the results of a sequential algorithm.

## 1 Introduction

Graph algorithms are a very common requirement in several problem domains including several scientific and engineering applications, which involves processing large graphs that have millions, if not billions, of vertices. Fundamental graph operations like breadthfirst search (BFS), depth-first search (DFS), shortest path, etc., occur frequently in these domains. While fast sequential implementations of the famous algorithms for these problems exists [1, 2], they are of the order of number of vertices and edges. On very large graphs, these algorithms become practically impractical. But when we switch to using parallel algorithms, we can achieve much more reasonable and practical implementations, in terms of both time and hardware cost on these basic graph problems [3].

The potential of GPU computing for large-scale graph processing was first explored by Harish and Narayanan [4], who implemented BFS and SSSP using CUDA on graphs with up to 10 million vertices. Among them, Breadth First Search (BFS) is one of the most fundamental graph traversal algorithms, and its efficient implementation on GPUs

has a significant impact on the performance of the entire graph computing system. This experiment is based on the NVIDIA RTX4070 GPU platform, aiming to compare and analyze various typical GPU BFS implementation methods, evaluate their performance differences, resource utilization, and potential bottlenecks in actual operation, explore the impact of thread organization, memory access patterns, and frontier queue management on algorithm efficiency, and provide data support and practical guidance for the optimization and system design of subsequent graph algorithms.

# 2 BREADTH FIRST SEARCH

In the BFS problem you are given an undirected, unweighted graph G (V , E) and a source vertex S, and we need to find the minimum number of edges needed to reach every vertex in G from source vertex S. An asymptotically optimal sequential solution for the problem takes O (V + E) time. BFS has been used in state space searching, graph partitioning, automatic theorem proving, etc., and is one of the most used graph orientation in various practical graph algorithms. We will now discuss various CUDA based approaches to tackle the problem of BFS.

## 2.1 First Approach: Simple BFS

Luo et al. [5] introduced one of the earliest GPU BFS implementations that outperformed CPU counterparts, using a hierarchical queue and layered kernel design. Unlike prior work, their method offered the same computational complexity as sequential BFS, while exploiting hardware concurrency more effectively. This project uses a similar approach. In this approach, we will use level synchronization to solve BFS. The graph is traversed level by level and once visited it is not visited again. The BFS distance stores the level at which the nodes were processed. The level variable keeps track of the level being processed. Maintaining a queue for each vertex would incur additional overheads and thereby slow down the speed of execution. Instead in this implementation we give one thread to each vertex. In each iteration, if a vertex is at the same level as the level being processes, it fetches its cost from the cost array and updates its neighbours if more than 'its cost + 1' i.e. the neighbour was univisted. Here, a global flag keeps track of the changes made in each iteration. If no change is made in an iteration it implies that all vertices have been visited and BFS completed. The above approach has been presented in Algorithm 1 and the corresponding kernel in Algorithm 2.

---
**Algorithm 1:** simpleBFS_Host
---
**Input:** Vertex array $V_a$, Edge array $E_a$, Source vertex $S$
**Output:** Distance array $D_{ista}$, Parent array $P_a$

**1** Create distance array $D_{ista}$ and parent array $P_a$ of size $|V|$;
**2** Initialize all elements of $D_{ista}$, $P_a$ to $\infty$;
**3** $D_{ista}[S] \leftarrow 0$;
**4** $level \leftarrow 0$;
**5** $flag \leftarrow$ True;
**6** **while** $flag$ **do**
**7** $\quad$ $flag \leftarrow$ False;
**8** $\quad$ Invoke $\texttt{simpleBFS\_kernel}(level, V_a, E_a, D_{ista}, flag)$;
**9** $\quad$ $level \leftarrow level + 1$;
**10** **end**

---
**Algorithm 2:** simpleBFS_kernel
---
**Input:** Current level $level$, Vertex array $V_a$, Edge array $E_a$, Distance array
$\quad\quad\quad$ $D_{ista}$, Flag $flag$

**1** $tid \leftarrow \texttt{getThreadID()}$;
**2** $f \leftarrow$ False;
**3** **if** $tid < |V|$ **and** $D_{ista}[tid] = level$ **then**
**4** $\quad$ $u \leftarrow tid$;
**5** $\quad$ **foreach** $v \in neighbours\ of\ u$ **do**
**6** $\quad\quad$ **if** $level + 1 < D_{ista}[v]$ **then**
**7** $\quad\quad\quad$ $D_{ista}[v] \leftarrow level + 1$;
**8** $\quad\quad\quad$ $f \leftarrow$ True;
**9** $\quad\quad$ **end**
**10** $\quad$ **end**
**11** **end**
**12** **if** $f = True$ **then**
**13** $\quad$ $flag \leftarrow$ True;
**14** **end**

---

## 2.2 Second Approach: Queue BFS

In this approach as well, we will use level synchronization to solve BFS. Like we have seen earlier, the graph is traversed level by level and once visited it is not visited again. But here, we maintain a queue cQ of the vertices to be visited. A cost array, Ca, stores the minimal number of edges of each vertex from source S. In each iteration, we traverse through the queue, and for each vertex u, if we find a unvisited neighbour v, we update the distance value for v, and add it to the queue of vertices to be expanded next, nQ. We use atomicMin and atomicAdd operations to avoid race conditions. The above approach has been presented in Algorithm 3 and the corresponding kernel in Algorithm 4.

---

**Algorithm 3:** queueBFS_Host

---

**Input:** $V_a, E_a, S$ ;                          // The graph $G(V, E)$ and source $S$
**Output:** Distance array $D_{ista}$, Parent array $P_a$

**1** Create cost array $D_{ista}$ and parent array $P_a$ of size $|V|$ and initialise all values
  to $\infty$;

**2** Create two arrays $cQ$ and $nQ$, initialise $cQ \leftarrow \{S\}$, $nQ \leftarrow \emptyset$;

**3** $D_{ista}[S] \leftarrow 0$;

**4** $P_a[S] \leftarrow -1$;

**5** $l \leftarrow 0$ ;                          // Start with the source vertex

**6 while** $|cQ| > 0$ **do**

**7**  | Invoke queueBFS_kernel$(l, V_a, E_a, D_{ista}, P_a, cQ, nQ)$;

**8**  | swap$(cQ, nQ)$;

**9**  | $nQ \leftarrow \emptyset$;

**10** | $l \leftarrow l + 1$;

**11 end**

---

**Algorithm 4:** queueBFS_kernel

---

**Input:** $l, V_a, E_a, D_{ista}, P_a, cQ, nQ$

**1** $tid \leftarrow$ getThreadID();

**2 if** $tid < |cQ|$ **then**

**3**  | $u \leftarrow cQ[tid]$;

**4**  | **foreach** $v$ *neighbor of* $u$ **do**

**5**  |  | **if** $D_{ista}[v] = \infty$ **and** $atomicMin(D_{ista}[v], l + 1) = \infty$ **then**

**6**  |  |  | $P_a[v] \leftarrow u$;

**7**  |  |  | $pos \leftarrow$ atomicAdd$(nQ\_size, 1)$;

**8**  |  |  | $nQ[pos] \leftarrow v$;

**9**  |  | **end**

**10** | **end**

**11 end**

---

## 2.3   Third Approach: Scan BFS

A work-efficient parallel BFS algorithm should perform O (V + E) work. To achieve
this, each iteration should examine only the edges and vertices in that iteration's logical
edge and vertex-frontiers, respectively. Edge frontiers are vertices that needs to be
examined in the next step of the BFS. In this algorithm, frontier is managed outof-core
and is fully produced in off-chip memory for consumption by the next BFS iteration
after a global synchronization step.

The intuition behind the linear optimisation of BFS is quite same as the standard
BFS. We maintain a queue of unvisited vertex, levelsynchronised, and we terminate
when the queue is empty. The only challenge is to populate this queue efficiently. Here,
we employ prefix scan to find out the position of vertices in queue for the next iteration.
For each vertex in the queue, we assign a new thread to compute the frontiers for the
next iteration. The above approach has been presented in Algorithm 5, and the kernels
in the subsequent Algorithms 6, 7, 8 and 9.

---

**Algorithm 5:** scanBFS_Host

---

**Input:** $V_a, E_a, S$ ;                          // The graph $G(V, E)$ and source $S$
**Output:** Distance array $D_{ista}$, Parent array $P_a$

1  Create arrays $Dega, PreDega$ of size $|V|$ and initialise all values to 0;
2  Create distance array $D_{ista}$ of size $|V|$ and initialise all values to $\infty$;
3  Create mask array $P_a$ of size $|V|$ and initialise all values to $-1$;
4  Create queues $cQ \leftarrow \{S\}$, $nQ \leftarrow \emptyset$;
5  $D_{ista}[S] \leftarrow 0$, $P_a[S] \leftarrow -1$, $l \leftarrow 0$;
6  **while** $|cQ| > 0$ **do**
7  $\quad$ Invoke $\texttt{nextLayer}(l, V_a, E_a, P_a, D_{ista}, cQ)$;
8  $\quad$ Invoke $\texttt{countDegrees}(V_a, E_a, P_a, cQ, Dega)$;
9  $\quad$ Invoke $\texttt{scanDegrees}(|cQ|, Dega, PreDega)$;
10 $\quad$ Perform prefix sum on $Dega$, store result in $PreDega$;
11 $\quad$ $nQ \leftarrow PreDega[|cQ|/\texttt{NUM\_THREADS}]$;
12 $\quad$ Invoke $\texttt{populateNextQueue}(V_a, E_a, P_a, cQ, nQ, Dega, PreDega)$;
13 $\quad$ $cQ \leftarrow nQ$;
14 $\quad$ $l \leftarrow l + 1$;
15 **end**

---

---

**Algorithm 6:** nextLayer

---

**Input:** $l, V_a, E_a, P_a, D_{ista}, cQ$

1  $tid \leftarrow \texttt{getThreadId()}$;
2  **if** $tid < |cQ|$ **then**
3  $\quad$ $u \leftarrow cQ[tid]$;
4  $\quad$ **foreach** $v$ *neighbor of* $u$ **do**
5  $\quad\quad$ **if** $D_{ista}[v] > l + 1$ **then**
6  $\quad\quad\quad$ $D_{ista}[v] \leftarrow l + 1$;
7  $\quad\quad\quad$ $P_a[v] \leftarrow u$;
8  $\quad\quad$ **end**
9  $\quad$ **end**
10 **end**

---

---

**Algorithm 7:** countDegrees

---

**Input:** $V_a, E_a, P_a, cQ, Dega$

1  $tid \leftarrow \texttt{getThreadId()}$;
2  **if** $tid < |cQ|$ **then**
3  $\quad$ $u \leftarrow cQ[tid]$,    $d \leftarrow 0$;
4  $\quad$ **foreach** $v$ *neighbor of* $u$ **do**
5  $\quad\quad$ **if** $P_a[v] = E_a.\boldsymbol{index}(v)$ ***and*** $v \neq u$ **then**
6  $\quad\quad\quad$ $d \leftarrow d + 1$;
7  $\quad\quad$ **end**
8  $\quad$ **end**
9  $\quad$ $Dega[tid] \leftarrow d$;
10 **end**

---

---

**Algorithm 8:** scanDegrees

    **Input:** $cQ_{\text{size}}, Dega, PreDega$

    /* Get the Id of the thread                                                       */

**1**  $tid \leftarrow$ getThreadId()

**2**  **if** $tid < cQ_{size}$ **then**

**3**      Create a shared array $preSum$ of size NUM_THREADS

**4**      $m \leftarrow$ threadId.x

**5**      $preSum[m] \leftarrow Dega[tid]$

**6**      sync_threads

**7**      $n \leftarrow 2$

**8**      **while** $n \leq$ *NUM_THREADS* **do**

**9**           **if** $\textbf{\textit{bitwiseAnd}}(m, n-1) = 0$ *and* $tid + (2 \cdot n) < cQ_{size}$ **then**

**10**               $preSum[m] \leftarrow preSum[m] + preSum[tid + (2 \cdot n)]$

**11**           **end**

**12**           sync_threads

**13**           $n \leftarrow 2 \cdot n$

**14**      **end**

**15**      **if** $m = 0$ **then**

**16**           $PreDega[\frac{tid}{\text{NUM\_THREADS}} + 1] \leftarrow preSum[m]$

**17**      **end**

**18**      $n \leftarrow$ NUM_THREADS

**19**      **while** $n > 1$ **do**

**20**           **if** $\textbf{\textit{bitwiseAnd}}(m, n-1) = 0$ *and* $tid + (n/2) < cQ_{size}$ **then**

**21**               $temp \leftarrow preSum[m]$

**22**               $preSum[m] \leftarrow preSum[m] + preSum[tid + (n/2)]$

**23**               $preSum[tid + (n/2)] \leftarrow temp$

**24**           **end**

**25**           sync_threads

**26**           $n \leftarrow n/2$

**27**      **end**

**28**      $Dega[tid] \leftarrow preSum[m]$

**29**  **end**

---

**Algorithm 9:** assignVerticesNextQueue

**Input:** $V_a, E_a, P_a, cQ, nQ, Dega, PreDega$

```
/* Get the Id of the thread                                          */
```
1   $tid \leftarrow$ getThreadId();
2   **if** $tid < cQ_{size}$ **then**
3     Initialise a shared variable $i$;
4     **if** $threadIdx.x = 0$ **then**
5       $i \leftarrow$ PreDega[NUM_THREADS];
6     **end**
7     sync_threads;
8     $s \leftarrow 0$;
9     **if** $threadIdx.x \neq 0$ **then**
10      $s \leftarrow$ Dega$[tid - 1]$;
11    **end**
12    $u \leftarrow$ cQ$[tid]$;
13    $c \leftarrow 0$;
14    **foreach** $v$ **in** *neighbours of* $u$ **do**
15      **if** $P_a[v] = E_a.\textbf{index}(v)$ **and** $v \neq u$ **then**
16       $nQ[i + s + c] \leftarrow v$;
17       $c \leftarrow c + 1$;
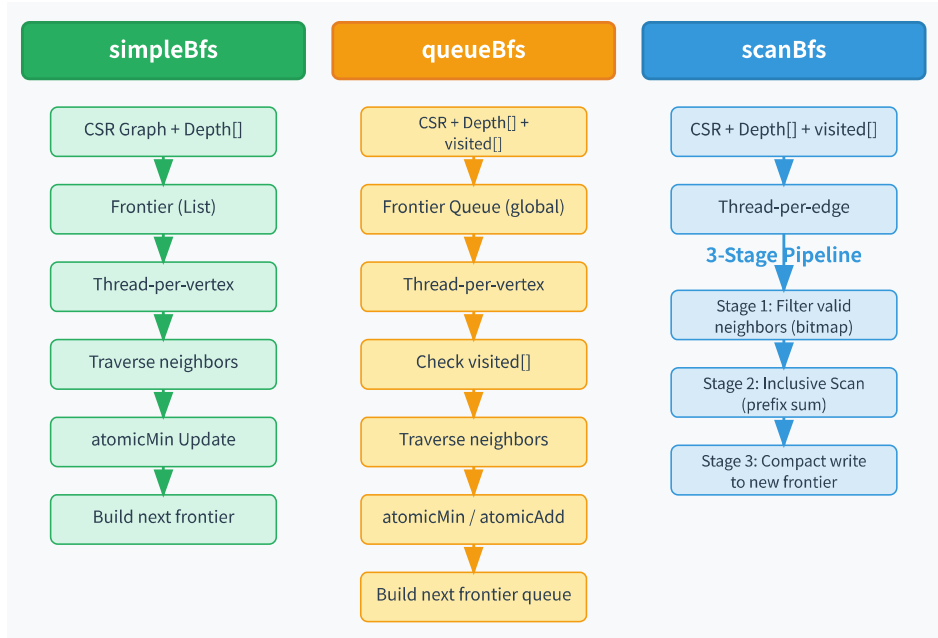18      **end**
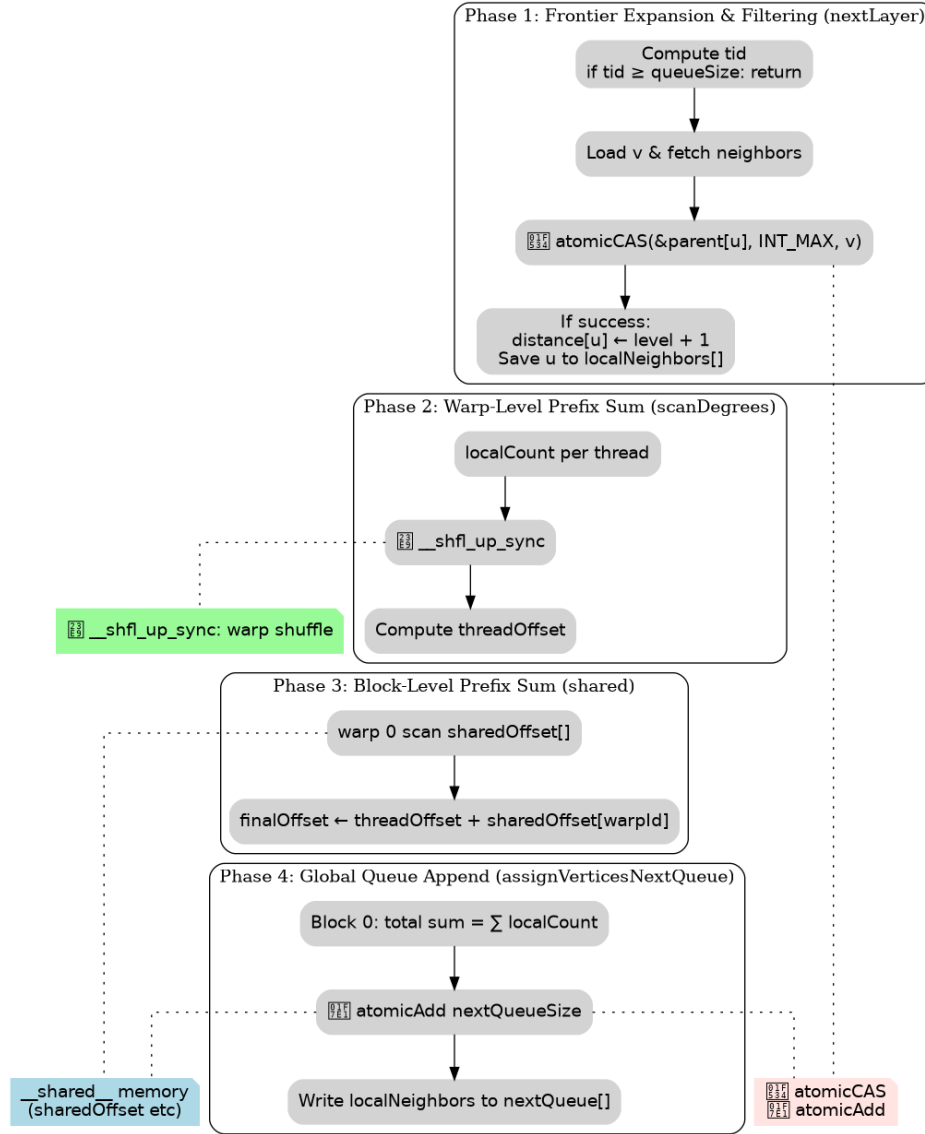19    **end**
20 **end**



Figure 1: bfs execution models

Figure 2: scanBfsFusedKernel

## 2.4  Fourth Approach: Fused Scan BFS

To reduce kernel launch overhead and improve memory locality, we design a fused implementation of the scan-based BFS approach, termed Fused Scan BFS. This approach integrates the four phases—nextLayer, countDegrees, scanDegrees, and assignVerticesNextQueue—into a single CUDA kernel. The kernel not only performs frontier expansion and filtering, but also handles intra-block parallel prefix sum and global output position computation, thereby minimizing inter-kernel synchronization.

The kernel execution is launched in a grid of thread blocks, each containing a fixed number of threads (e.g., 256). Each thread first processes a vertex v from the current frontier queue currentQueue[tid]. The thread traverses the adjacency list of v, filters out visited neighbors via atomic CAS on the parent array, and records valid successors into a local buffer localNeighbors[32]. During this process, atomic operations are used to update the distance and parent arrays.

To determine the global output position in the nextQueue, the kernel performs a warp-level prefix sum on the per-thread output count using CUDA's ˍshflˍupˍsync intrinsic. Each warp's total output is stored in shared memory (sharedOffset). Subsequently, a block-level scan is performed on the warp-level totals by warp 0 to compute the prefix sums across warps. The final global offset is computed by atomically adding the block's total output count to nextQueueSize.

This design avoids repeated memory writes and kernel switches by combining all stages into one pass. Shared memory is heavily utilized to store intermediate offsets, per-thread output counts, and warp sums, ensuring minimal global memory access.

Key advantages of this fused approach include: reduced kernel launch latency due to single-kernel structure; improved memory locality since each thread handles vertex expansion and result placement in situ; efficient load balancing via warp-centric output accumulation and prefix sum; scalable global queue writing using a global atomic offset obtained once per block.

The launch configuration dynamically determines grid and block sizes based on the frontier size. The shared memory size is also computed accordingly to accommodate all intermediate buffers.

In summary, the Fused Scan BFS design represents a highly optimized traversal strategy tailored for GPU architectures, leveraging warp-level and block-level coordination to accelerate irregular graph expansion workloads.

## 2.5 Fifth Approach: Warp-Merge Load-Balanced Fused BFS

To further improve the performance of BFS traversal on irregular graphs, we implement a warp-merge load-balanced fused kernel, referred to as scanBfsFusedKernel_WM. Unlike the previous fused kernel which assigns one thread to each frontier vertex, this implementation utilizes warp-level collaboration to handle the large variance in vertex degrees more efficiently. The idea is to group threads into warps, assign multiple vertices to each warp, and then allow threads within the warp to collectively process the outgoing edges of these vertices.

**Specific implementation** Each warp begins by cooperatively loading up to 32 vertices from the frontier queue into shared memory. For each vertex, its adjacency offset and degree are also cached. Then, a warp-level prefix sum is performed to compute the relative position of each vertex's edges in the warp-wide edge list. This ensures that threads can access the correct range of edges during edge traversal. To achieve this, we use warp shuffle instructions to implement an efficient intra-warp prefix sum, avoiding global synchronization.

Once the total number of edges to be processed by the warp is known, threads iterate over this edge list in a strided manner, such that each thread processes every 32nd edge. For each edge, the thread checks whether the neighbor has been visited using an atomic compare-and-swap operation. If the neighbor is discovered for the first time, its distance and parent are updated, and it is inserted into the next frontier queue using an atomic add.

This warp-centric strategy helps distribute high-degree vertices across multiple threads, alleviating the problem of load imbalance caused by assigning large workloads to individual threads. Furthermore, it leverages the CUDA warp execution model to achieve

**Stage 1: Input Loading Phase**

**Current Frontier Queue**
v0, v1, v2, ..., v31
(32 vertices per warp)

v0
start: 10
degree: 5

v1
start: 15
degree: 3

...

load →

**Shared Memory**
vList[32]: v0, v1, ..., v31
vStartOffsets[32]
vDegrees[32]
prefixDegrees[32]

Each warp loads 32 vertices and their degree information using laneId

**Stage 2: Intra-Warp Processing Phase**
Warp (32 threads)

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |

__shfl_up_sync (prefix sum)

**Prefix Sum on Degrees**
prefixDegrees[i] =
sum(degrees[0..i])

**Strided Edge Processing**
🔒 atomicCAS
Update distance[u], parent[u]
atomicAdd to nextQueue counter

**Adjacency List**
edge0, edge1
edge2, edge3
...

**Stage 3: Global Write-back and Queue Update Phase**

**Thread Counts**
threadCounts[0] = 2
threadCounts[1] = 1

aggregate →

**Shared Global Offset**
sharedOffset[32]
sharedGlobalOffset

write →

**Next Frontier Queue**
nextQueue[pos] = u
u0, u1, u2, ..., uN
(newly discovered vertices)

**Parallel Prefix Sum (Optional)**
Block-level scan for sharedOffset calculation
Warp-level scan for sharedGlobalOffset
Efficient position calculation for queue writes

**Legend:** → Data Flow --→ Thread Synchronization → Atomic Operations ▢ Shared Memory
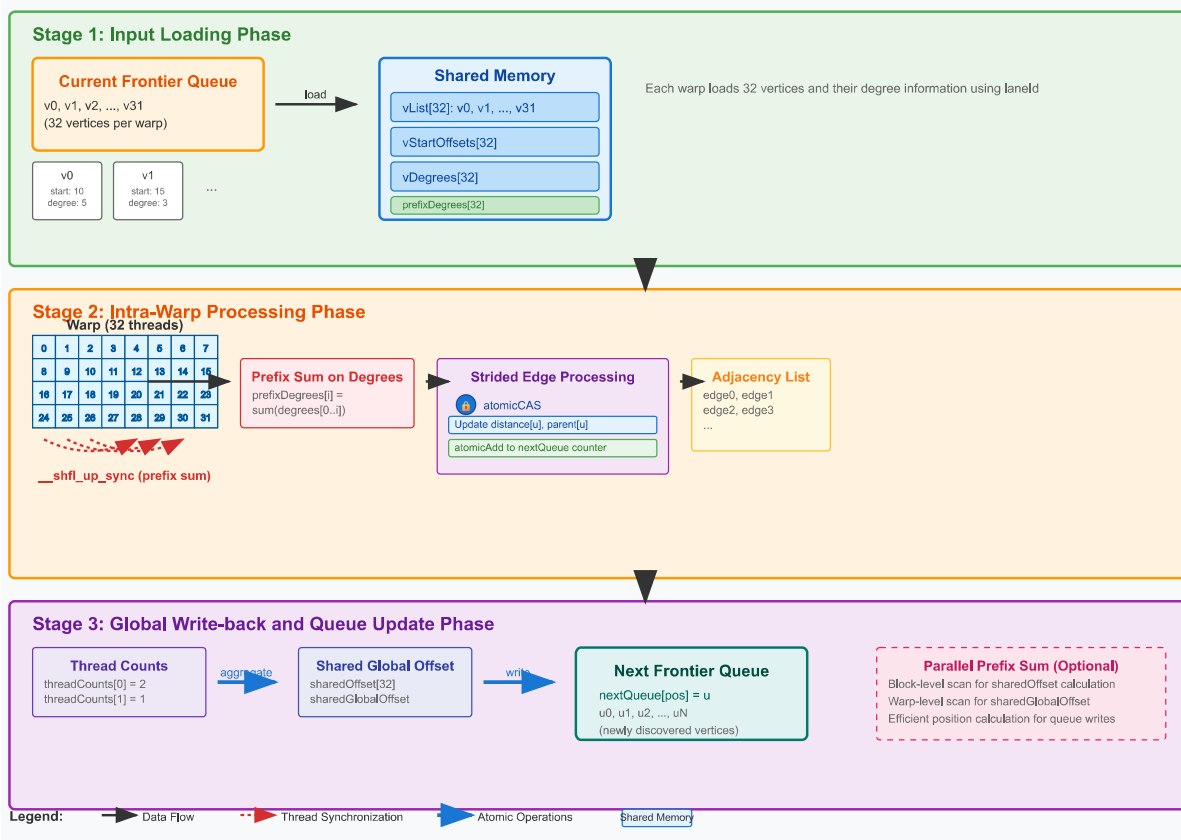
Figure 3: Warp-Merge Load-Balanced BFS Kernel Execution

low-overhead synchronization and maximizes memory coalescing by accessing adjacency lists in a structured manner. Since all warps share the same fused kernel logic, it also avoids redundant kernel launches and host-device synchronization.

The kernel launch configuration dynamically adjusts the grid size based on the number of warps required, and shared memory allocation is carefully computed to hold intermediate metadata such as vertex lists, offsets, degrees, prefix sums, and per-thread output counters. At the end of each iteration, newly discovered vertices are stored into the global queue, and the algorithm proceeds to the next BFS level. This implementation demonstrates improved performance especially for scale-free graphs and other graphs with highly skewed degree distributions.

## 2.6 Sixth Approach: Coalesced Mapping Load-Balanced Fused BFS

In this section, we present a fused BFS implementation based on the Coalesced Mapping (CM) load balancing strategy, namely scanBfsFusedKernel_CM. This approach adheres to the fused design philosophy by integrating vertex expansion, degree computation, prefix sum, and neighbor assignment into a single kernel. Unlike the warp-level WM strategy, CM emphasizes block-level coordination, where each block is responsible for handling a batch of vertices. These vertices' adjacency edges are compressed into a contiguous segment, and threads cooperatively process them in a balanced manner.
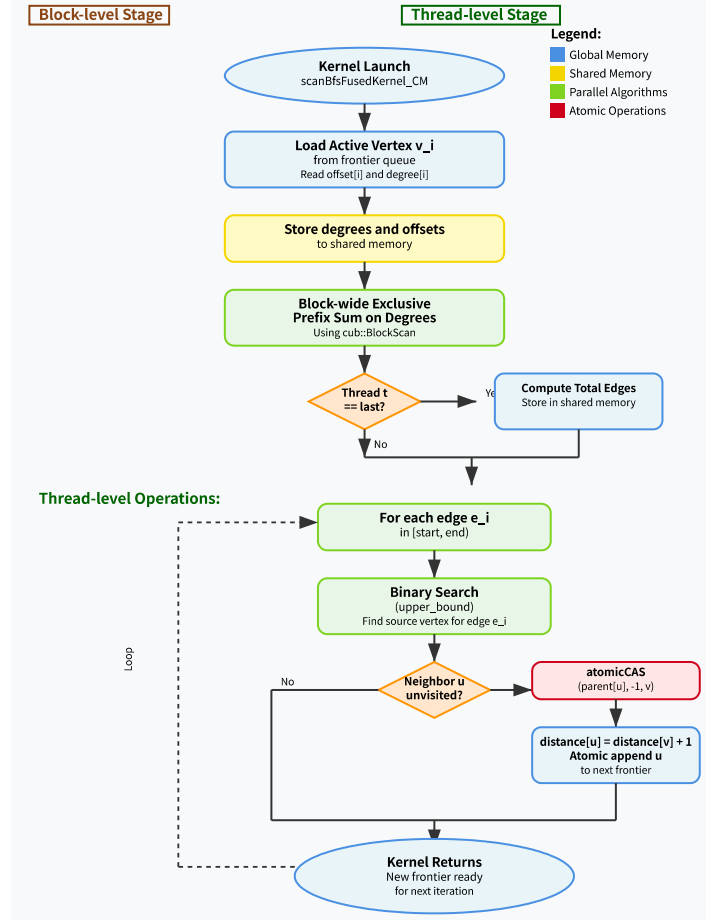
10

Figure 4: scanBfsFusedKernel_CM Flow

**Specific process**   At the beginning of kernel execution, each thread loads its assigned vertex from the current frontier queue along with its corresponding edge offset and degree. These values are stored in shared memory for fast access. A parallel exclusive prefix sum is then computed on the degrees using CUB's BlockScan::ExclusiveSum, allowing the block to determine a compacted index range for all outgoing edges. The final thread in the block also computes the total number of edges to be processed.

Threads then proceed to cooperatively iterate over this edge range. A binary search (upper_bound-like logic) is employed to determine the source vertex for each edge, allowing each thread to identify which vertex its assigned edge belongs to. This enables a uniform distribution of edge processing across all threads, helping avoid load imbalance caused by skewed vertex degrees.

When a thread encounters a neighbor that has not yet been visited, it attempts to atomically update the parent pointer using atomicCAS. If successful, the thread updates the neighbor's distance and appends it to the next-level frontier queue using an atomic increment. Compared to vertex-centric approaches, this edge-centric fused strategy allows for finer-grained work assignment, better memory coalescing, and improved parallel efficiency.

The host-side routine runCudaScanfusedBfs_CM manages the overall BFS process, repeatedly launching the fused kernel and updating the frontier queues until no vertices remain. Overall, the CM-style BFS demonstrates high efficiency, particularly for graphs

with significant degree variance, where its balanced thread scheduling and compacted memory accesses prove highly effective.

# 3   Experimental Environment and Evaluation Methodology

**Platform**   All experiments in this study were conducted on a workstation equipped with an NVIDIA RTX 4070 GPU, featuring 5888 CUDA cores, a base clock of 1.92 GHz, and 12 GB of GDDR6 memory with a theoretical memory bandwidth of 504 GB/s. The software environment included CUDA Toolkit 12.3, NVIDIA driver version 550.xx, and Nsight Compute 2024.1 and Nsight Systems 2024.1 for profiling analysis. The operating system was Ubuntu 22.04 LTS with Linux kernel version 6.2.

Table 1: Representative graphs used in benchmarking

| Graph | Vertices | Edges | Domain |
|---|---|---|---|
| soc-orkut | 3M | 212.7M | SN |
| web-uk-2005 | 129K | 23M | WG |
| kron_g500-logn21 | 2.1M | 182.1M | GG |
| roadNet-CA | 1.9M | 5.5M | RN |
| sc-msdoor | 415K | 19.8M | SC |
| as-caida2007 | 26.5K | 106.8K | CA |
| com-youtube | 1.13M | 2.99M | YT |
| p2p-Gnutella08 | 6.3K | 20.8K | GN |
| wiki-talk | 2.39M | 5.02M | WT |
| wiki-talk-temporal | 1.14M | 3.31M | WT |

**Datasets**   To comprehensively evaluate the performance of BFS implementation, we used 10 graph datasets from SuiteSpearse Matrix Collection and SNAP, with representative datasets such as: 'roadNet-CA' (a sparse road network graph with high diameter and low degree), 'soc-orkut' (a social graph with small-world properties and power-law degree distribution), 'kron_g500-logn21' (a synthetic graph generated using Kronecker models, known for high scalability), 'msdoor' (a moderately-sized mesh-like graph), and 'web-uk-2005' (a large-scale web crawl graph with high heterogeneity). These datasets cover diverse graph topologies, from highly structured road networks to irregular and large-scale internet graphs, thereby enabling a broad analysis of algorithmic behavior under different structural constraints.

**Evaluation method**   Each BFS variant— simpleBfs, queueBfs, scanBfs, scanfusedbfs, scanfusedBfs_WM and scanfusedBfs_CM compiled using CMake 3.26 and GCC 11.4. The binaries were generated in release mode, and the kernel launch configurations (grid and block sizes) were explicitly controlled to align with each algorithm's internal design. The 'bfs_exec' program accepted '.mtx' input files and command-line options for specifying source vertex and verbosity levels. To ensure consistency, the BFS traversal

was always initiated from vertex 0, and all experiments were executed on GPU device 0 without concurrent background processes.

The evaluation process involved three key stages. First, end-to-end performance including graph loading, memory allocation, kernel execution, and result retrieval—was measured using Python scripts that invoked the BFS binary and logged timing information across runs. Each algorithm-dataset pair was executed five times, and the average values were reported. Second, detailed GPU kernel performance metrics were collected using Nsight Compute (ncu). The profiling focused on key indicators such as kernel execution time, compute and memory throughput, SM occupancy, instruction-level parallelism, and memory access patterns. These metrics enabled us to isolate kernel-level efficiency and identify bottlenecks such as warp divergence or bandwidth saturation.

The implementation correctness of each BFS version was verified by comparing their output parent arrays against the reference implementation, ensuring semantic equivalence of traversal results. To measure relative performance across implementations, we used the kernel execution time of simpleBfs as the baseline and computed the speedup factor for each method on each dataset. In addition, Python scripts using 'matplotlib' and 'pandas' were employed to generate visualization figures, including stacked bar charts of kernel and non-kernel time and relative speedup plots.

This experimental framework allowed us to perform a holistic evaluation of BFS implementations across algorithmic, architectural, and system layers, revealing not only micro-level efficiency (e.g., per-kernel throughput) but also macro-level end-to-end trends (e.g., overhead from synchronization and data movement). Through the integration of large-scale real and synthetic graphs, modern profiling tools, and structured comparative analysis, we aim to uncover both the advantages and trade-offs of each BFS strategy under realistic deployment conditions.

**Step1** :This experiment uses the camke tool to generate executable files. It can automatically execute programs, collect data, and draw images through script files.

```
mkdir build
cd build
cmake ..
make
bash ./run_all.sh
```

**Step2** :Analyze Kernel startup and system bottlenecks using nsys

```
nsys profile −o report_xxx ./xxxx
```

Obtain the.nsys-rep file, which can be opened with the nsight sys GUI.

**Step2** : Analyze the internal execution efficiency of the Kernel using ncu

```
ncu −−set full −o *** ./xxx
```

After completion, a ***.ncu-rep file will be generated on the server and can be opened locally using Nsight Compute.
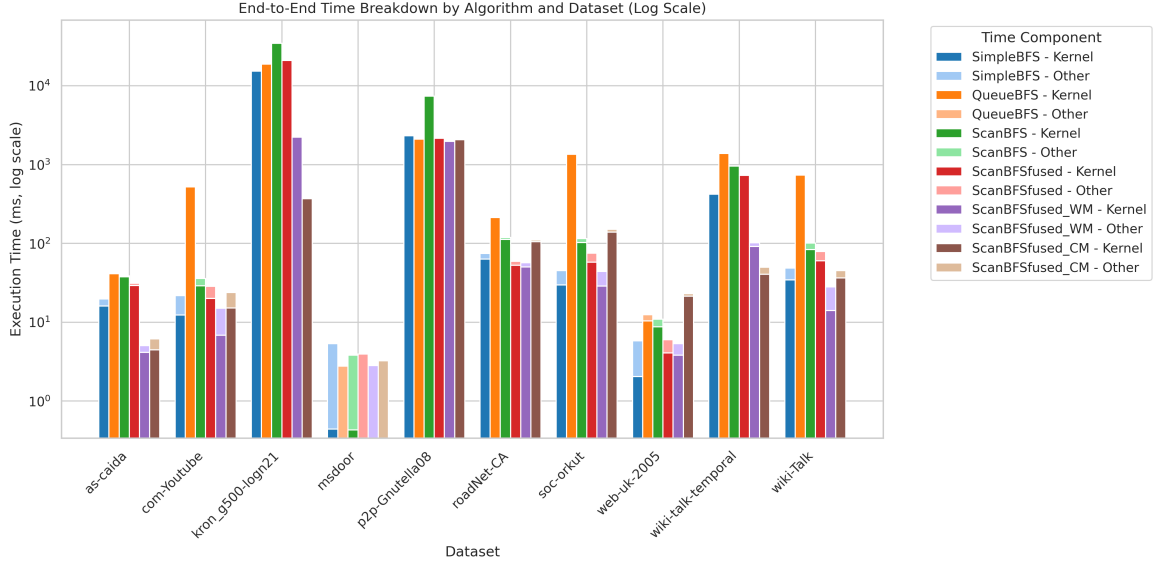
# 4 Comprehensive Analysis and Conclusion



Figure 5: end to end breakdown stacked log



Figure 6: relative speedup

**Overall Performance**

**SimpleBFS Analysis**   The simpleBfs algorithm performs well on structured graphs with shallow frontiers and low-degree nodes, where its straightforward design and low overhead enable efficient execution. It is particularly effective on datasets like msdoor and roadNet-CA, outperforming more complex algorithms in both kernel and end-to-end time. However, its scalability is limited. On larger, high-degree graphs such as

Table 2: Kernel Execution Time (ms) of Different BFS Algorithms on Various Datasets. Lower is better.

| Dataset | SimpleBFS | QueueBFS | ScanBFS | Scanfused | Scan_WM | Scan_CM |
|---|---|---|---|---|---|---|
| kron_g500-logn21 | 17144.2 | 20746.3 | 38571.3 | 21316.8 | 2249.69 | 367.089 |
| msdoor | 0.37488 | 0.119808 | 0.347392 | 0.248224 | 0.457952 | 0.367552 |
| roadNet-CA | 70.3312 | 219.964 | 136.855 | 55.5148 | 65.5848 | 97.3145 |
| soc-orkut | 29.8281 | 1354.49 | 102.302 | 57.7611 | 28.7526 | 138.723 |
| web-uk-2005 | 2.05197 | 10.3867 | 8.77875 | 4.09763 | 3.81114 | 21.3515 |
| as-caida | 15.2383 | 40.6508 | 34.1332 | 28.5826 | 3.81856 | 2.86198 |
| com-Youtube | 12.1147 | 499.67 | 27.0602 | 19.3098 | 5.75101 | 13.3751 |
| p2p-Gnutella08 | 0.982016 | 1.79917 | 1.86675 | 0.843616 | 0.691552 | 1.5968 |
| wiki-Talk | 35.4796 | 743.981 | 84.3735 | 60.8229 | 15.4551 | 43.7454 |
| wiki-talk-temporal | 448.514 | 1437.08 | 1007.98 | 765.706 | 94.5894 | 36.1944 |

soc-orkut and kron_g500-logn21, performance degrades significantly due to full-graph scans, inefficient memory access, and lack of load balancing. Overall, while suitable for low-parallelism workloads, simpleBfs struggles on irregular or large-scale graphs.

**QueueBFS and Atomic Conflicts**   The queueBfs algorithm demonstrates more structured parallelism than simpleBfs by explicitly managing frontier queues across levels, which improves its efficiency on graphs with moderate to high parallelism. On datasets like msdoor, it achieves excellent performance with a very low kernel execution time, matching or outperforming other methods. On roadNet-CA, however, it performs worse than simpleBfs, with an end-to-end time over 220ms, likely due to overhead from queue management and insufficient parallel work per level.

As graph complexity increases, queueBfs faces mixed results. On soc-orkut, its end-to-end time rises to over 1 seconds, better than simpleBfs, but still not ideal. It also shows better scalability than simpleBfs on wide-frontier graphs such as kron_g500-logn21, but remains limited by synchronization and atomic operations. In summary, queueBfs improves performance through better frontier handling, especially on medium-sized or moderately irregular graphs, but lacks advanced optimizations for large-scale or highly skewed workloads.

**ScanBFS and Prefix-Sum Strategy**   The scanBfs algorithm leverages parallel prefix sums to compute frontier expansions in a more balanced fashion, reducing contention from atomic operations and enabling better utilization of thread blocks. However, this benefit comes with increased kernel complexity and setup overhead. On small graphs like msdoor, it underperforms compared to simpler methods due to redundant work and the cost of prefix sum operations. For example, on kron_g500-logn21, scanBfs requires over 38 seconds end-to-end—worse than both simpleBfs and queueBfs. Overall, scanBfs is more resilient to load imbalance and performs well on large or irregular datasets, but its overhead makes it less suitable for lightweight workloads.

Table 3: End-to-End Time (ms) of Different BFS Algorithms on Various Datasets. Lower is better.

| Dataset | SimpleBFS | QueueBFS | ScanBFS | Scanfused | Scan_WM | Scan_CM |
|---|---|---|---|---|---|---|
| kron_g500-logn21 | 17159.5000 | 20758.7000 | 38579.0000 | 21332.7000 | 2262.7400 | 375.8160 |
| msdoor | 5.0819 | 2.5250 | 2.3767 | 3.9927 | 3.1844 | 3.4058 |
| roadNet-CA | 81.6281 | 229.8660 | 145.0410 | 62.1320 | 72.8931 | 105.2880 |
| soc-orkut | 45.3027 | 1368.3000 | 115.1190 | 74.9294 | 44.0812 | 150.3970 |
| web-uk-2005 | 5.7946 | 12.4362 | 10.9595 | 5.9936 | 5.3537 | 22.9829 |
| as-caida | 18.7561 | 41.7467 | 35.4949 | 30.5434 | 4.4495 | 3.3841 |
| com-Youtube | 21.9658 | 507.7740 | 35.4213 | 27.2862 | 12.5332 | 18.4332 |
| p2p-Gnutella08 | 5.1445 | 2.3856 | 2.4377 | 1.7030 | 1.1825 | 2.0652 |
| wiki-Talk | 50.0429 | 754.6200 | 100.5190 | 76.8103 | 27.8521 | 58.3836 |
| wiki-talk-temporal | 459.2720 | 1444.3700 | 1013.9500 | 774.5360 | 102.6870 | 43.3535 |

**Scanfused kernel fusion strategy** The Scanfused implementation, which consolidates the four scanBFS stages into a single kernel, offers a significant reduction in kernel launch overhead and enhances memory locality. This fused design is particularly effective on large-scale graphs with wide frontiers and heavy workloads per BFS level. For example, on roadNet-CA, it achieves an end-to-end time of 62ms, outperforming both simpleBfs and queueBfs by a large margin. Its advantage stems from executing all traversal stages in a single pass, allowing shared memory to be reused across phases and avoiding expensive inter-kernel synchronization.

On mid-sized or sparse graphs like wiki-Talk, Scanfused maintains stable performance, comparable to that of scanBfs, but without the added synchronization overhead of multiple kernels. However, the complexity of the fused kernel and its shared memory demands can limit its efficiency on very small graphs, where the benefits of fusion are less pronounced.

Overall, Scanfused delivers robust performance across diverse datasets, combining the scalability of scan-based approaches with improved efficiency through kernel fusion. Its design makes it particularly well-suited for high-parallelism workloads with deep or irregular frontiers.

**scanBfsFused_WM warp-merge** The scanBfsFused_WM kernel builds upon the fused scanBFS strategy by incorporating warp-level load balancing to better handle irregular and high-degree graphs. By assigning multiple vertices to each warp and enabling cooperative edge traversal, it significantly reduces load imbalance compared to thread-centric approaches.

This design is particularly effective on graphs like soc-orkut and wiki-Talk, where vertex degree variance is high. For example, on soc-orkut, it achieves better end-to-end performance than queueBfs and the standard Scanfused, highlighting its ability to handle dense frontiers efficiently. The use of warp-level prefix sum and shared memory coordination ensures low synchronization cost while maintaining high memory through-

put.

Overall, scanBfsFused_WM strikes a strong balance between kernel fusion efficiency and dynamic workload distribution, making it well-suited for complex, scale-free graph structures.

**scanBfsFused_CM Coalesced Mapping** The scanBfsFused_CM kernel adopts a coalesced mapping strategy that shifts the focus from vertex-level to edge-level load balancing. By aggregating edge segments across a block and distributing them uniformly, it ensures even workload distribution even in the presence of highly skewed degree distributions.

This kernel demonstrates strong performance on large and irregular graphs. On wiki-talk-temporal, for instance, it achieves lower end-to-end time than traditional scan-based methods, benefiting from its edge-centric processing model and minimal kernel launch overhead. The use of CUB's efficient block-level prefix sum and shared memory buffering contributes to improved locality and reduced contention during neighbor updates.

In general, scanBfsFused_CM is well-suited for high-degree, scale-free graphs, offering both scalability and robustness through its balanced, fused execution model.

Table 4: Performance (ms) of Gunrock and Gswitch on Various Datasets.

| Metric | kron_g500-logn21 | msdoor | roadNet-CA | soc-orkut | web-uk-2005 |
|---|---|---|---|---|---|
| Gunrock (Kernel) | 23.28 | 28.48 | 108.33 | 30.57 | 4.37 |
| Gunrock (End-to-End) | 74159.2 | 12193.4 | 1792.32 | 90254.7 | 6685.15 |
| Gswitch (Kernel) | 8.34 | 15.9 | 82.69 | 8.59 | 6.52 |
| Gswitch (End-to-End) | 29547.2 | 6197.52 | 1785.65 | 47196.9 | 4428.43 |

**Comparative Analysis of Gunrock and GSWITCH** This experiment also tested the execution time of the gunrock[6] and GSWITCH[7] projects on several individual datasets.Gunrock delivers consistently strong kernel execution and end-to-end performance across a wide range of datasets. For instance, on kron_g500-logn21, it achieves a kernel time of 23.28 ms and an end-to-end time of 92.3 ms. These numbers reflect the framework's mature optimizations, including efficient frontier filtering, load balancing, and memory management.

GSWITCH, on the other hand, employs dynamic switching among three different BFS strategies (top-down, bottom-up, and hybrid) based on frontier density and graph structure. While previous hybrid approaches such as [8] dynamically switch between CPU and GPU backends, GSWITCH method focuses on GPU-internal kernel reconfiguration by examining frontier size and topology to choose between push, pull, or fused modes in real time. As a result, it shows particularly impressive results on large-scale

graphs like web-uk-2005, where it achieves 6.52 ms kernel time and 128 ms end-to-end time. This adaptability gives it an edge in handling diverse traversal phases more effectively.

When compared to our custom BFS kernels, such as Scanfused and its two variants, GSWITCH performs competitively in kernel time but is sometimes outpaced in end-to-end performance. For example, scanBfsFused_WM and scanBfsFused_CM often achieve lower end-to-end times on scale-free or small-world graphs due to their single-pass fused design, reduced kernel launch overhead, and better memory locality.

In contrast to Gunrock's general-purpose framework and GSWITCH's dynamic multi-kernel scheduling, our fused implementations exploit low-level CUDA optimizations tailored to the GPU hardware. While they may lack the generality and modularity of GSWITCH and Gunrock, they can outperform them on specific datasets by maximizing warp occupancy, avoiding synchronization barriers, and minimizing global memory transactions.

Overall, Gunrock and GSWITCH represent strong, versatile baselines. However, our fused designs demonstrate that hand-optimized, single-kernel strategies can rival or even surpass these systems under targeted workloads, especially when high throughput and minimal control divergence are essential.

## 4.1 Analysis of ncu results

| Summary | Details | Source | Context | Comments | Raw | Session | | | | | | Compare | Tools | View | Export | ☰ |

This table shows all results in the report. Use the column headers to sort the results in this report. Double-click a result to see detailed metrics. Double-click on demangled names to rename it.

| | Estimated Speedup [%] | Function Name | Demangled Name | Duration [us] (2,785.82 us) | Runtime Improvement [us] (2,649.38 us) | Compute Throughput [%] | Memory Throughput [%] | # Registers [register/thread] | Grid Size | Block Size |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 80.56 | simpleBfs | simpleBfs | 4.70 | 3.79 | 0.44 | 2.42 | 22 | 7, 1, … | 102 |
| 1 | 80.56 | simpleBfs | simpleBfs | 26.98 | 21.73 | 0.10 | 1.46 | 22 | 7, 1, … | 102 |
| 2 | 80.56 | simpleBfs | simpleBfs | 26.91 | 21.68 | 1.11 | 12.35 | 22 | 7, 1, … | 102 |
| 3 | 80.56 | simpleBfs | simpleBfs | 22.18 | 17.86 | 2.24 | 3.70 | 22 | 7, 1, … | 102 |
| 4 | 80.56 | simpleBfs | simpleBfs | 12.22 | 9.85 | 2.53 | 5.71 | 22 | 7, 1, … | 102 |
| 5 | 80.56 | simpleBfs | simpleBfs | 6.56 | 5.28 | 1.17 | 5.37 | 22 | 7, 1, … | 102 |
| 6 | 80.56 | simpleBfs | simpleBfs | 3.04 | 2.45 | 0.80 | 4.36 | 22 | 7, 1, … | 102 |
| 7 | 97.22 | queueBfs | queueBfs | 14.85 | 14.44 | 0.03 | 1.01 | 22 | 1, 1, … | 102 |
| 8 | 97.22 | queueBfs | queueBfs | 115.90 | 112.68 | 0.02 | 0.35 | 22 | 1, 1, … | 102 |
| 9 | 97.22 | queueBfs | queueBfs | 239.97 | 233.30 | 0.06 | 1.97 | 22 | 1, 1, … | 102 |
| 10 | 97.20 | queueBfs | queueBfs | 505.38 | 491.21 | 0.06 | 2.80 | 22 | 2, 1, … | 102 |
| 11 | 96.09 | queueBfs | queueBfs | 510.56 | 490.61 | 0.06 | 3.91 | 22 | 4, 1, … | 102 |
| 12 | 94.44 | queueBfs | queueBfs | 42.82 | 40.44 | 0.09 | 2.77 | 22 | 2, 1, … | 102 |
| 13 | 97.22 | queueBfs | queueBfs | 3.55 | 3.45 | 0.04 | 1.49 | 22 | 1, 1, … | 102 |
| 14 | 97.22 | nextLayer | nextLayer | 4.38 | 4.26 | 0.06 | 0.81 | 20 | 1, 1, … | 102 |
| 15 | 97.22 | countDegrees | countDegrees | 3.62 | 3.52 | 0.04 | 3.20 | 40 | 1, 1, … | 102 |
| 16 | 97.22 | scanDegrees | scanDegrees | 3.20 | 3.11 | 0.09 | 0.65 | 20 | 1, 1, … | 102 |
| 17 | 97.22 | assignVerticesNext… | assignVerticesNext… | 6.11 | 5.94 | 0.04 | 0.23 | 20 | 1, 1, … | 102 |
| 18 | 97.22 | nextLayer | nextLayer | 39.33 | 38.24 | 0.04 | 1.21 | 20 | 1, 1, … | 102 |
| 19 | 97.22 | countDegrees | countDegrees | 9.92 | 9.64 | 0.06 | 12.01 | 40 | 1, 1, … | 102 |
| 20 | 97.22 | scanDegrees | scanDegrees | 3.26 | 3.17 | 0.09 | 0.64 | 20 | 1, 1, … | 102 |
| 21 | 97.22 | assignVerticesNext… | assignVerticesNext… | 40.13 | 39.01 | 0.03 | 0.13 | 20 | 1, 1, … | 102 |
| 22 | 97.22 | nextLayer | nextLayer | 37.66 | 36.62 | 0.22 | 0.77 | 20 | 1, 1, … | 102 |
| 23 | 97.22 | countDegrees | countDegrees | 10.91 | 10.61 | 0.41 | 2.86 | 40 | 1, 1, … | 102 |
| 24 | 97.22 | scanDegrees | scanDegrees | 3.65 | 3.55 | 0.59 | 3.13 | 20 | 1, 1, … | 102 |

Figure 7: bfs result 1

Based on the profiling results obtained from NVIDIA Nsight Compute (NCU), we analyze the kernel-level performance characteristics of several BFS implementations. The performance data reveals clear contrasts in kernel duration, memory throughput, compute utilization, and register usage across different algorithms, shedding light on their architectural efficiency and limitations on GPU platforms.

The simpleBfs kernel exhibits low per-kernel durations (ranging from 3–27 µs), but its compute and memory throughput remain modest. Its design assigns one thread per frontier vertex and processes adjacency lists in a sequential fashion, resulting in under-utilized GPU resources. For instance, the memory throughput in most cases remains

| Estimated Speedup [%] | Function Name | Demangled Name | Duration [us] (2,785.82 us) | Runtime Improvement [us] (2,649.38 us) | Compute Throughput [%] | Memory Throughput [%] | # Registers [register/thread] | Grid Size | Block Size |
|---|---|---|---|---|---|---|---|---|---|
| 38 | 97.22 nextLayer | nextLayer | 3.52 | 3.42 | 0.04 | 1.03 | 20 | 1, 1, .. | 102 |
| 39 | 97.22 countDegrees | countDegrees | 3.62 | 3.52 | 0.04 | 12.64 | 40 | 1, 1, .. | 102 |
| 40 | 97.22 scanDegrees | scanDegrees | 3.26 | 3.17 | 0.22 | 0.69 | 20 | 1, 1, .. | 102 |
| 41 | 97.22 assignVerticesNext… | assignVerticesNext… | 5.47 | 5.32 | 0.03 | 0.94 | 20 | 1, 1, .. | 102 |
| 42 | 97.22 scanBfsFusedKernel | scanBfsFusedKernel | 10.18 | 9.89 | 0.13 | 0.74 | 28 | 1, 1, .. | 102 |
| 43 | 97.22 scanBfsFusedKernel | scanBfsFusedKernel | 53.57 | 52.08 | 0.06 | 0.97 | 28 | 1, 1, .. | 102 |
| 44 | 97.22 scanBfsFusedKernel | scanBfsFusedKernel | 51.94 | 50.49 | 0.23 | 0.89 | 28 | 1, 1, .. | 102 |
| 45 | 94.44 scanBfsFusedKernel | scanBfsFusedKernel | 44.90 | 42.40 | 0.59 | 2.77 | 28 | 2, 1, .. | 102 |
| 46 | 88.89 scanBfsFusedKernel | scanBfsFusedKernel | 21.15 | 18.80 | 1.57 | 5.42 | 28 | 4, 1, .. | 102 |
| 47 | 94.44 scanBfsFusedKernel | scanBfsFusedKernel | 15.30 | 14.45 | 0.37 | 2.42 | 28 | 2, 1, .. | 102 |
| 48 | 97.22 scanBfsFusedKernel | scanBfsFusedKernel | 7.10 | 6.91 | 0.16 | 0.90 | 28 | 1, 1, .. | 102 |
| 49 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 3.58 | 3.48 | 0.40 | 0.80 | 26 | 1, 1, .. | 102 |
| 50 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 19.04 | 18.51 | 0.11 | 0.39 | 26 | 1, 1, .. | 102 |
| 51 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 44.74 | 43.49 | 0.42 | 1.03 | 26 | 1, 1, .. | 102 |
| 52 | 94.44 scanBfsFusedKern… | scanBfsFusedKern… | 38.72 | 36.57 | 1.17 | 2.68 | 26 | 2, 1, .. | 102 |
| 53 | 88.89 scanBfsFusedKern… | scanBfsFusedKern… | 17.98 | 15.99 | 3.40 | 5.93 | 26 | 4, 1, .. | 102 |
| 54 | 94.44 scanBfsFusedKern… | scanBfsFusedKern… | 9.82 | 9.28 | 1.32 | 3.67 | 26 | 2, 1, .. | 102 |
| 55 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 5.89 | 5.72 | 0.33 | 3.19 | 26 | 1, 1, .. | 102 |
| 56 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 3.87 | 3.76 | 0.39 | 0.79 | 40 | 1, 1, .. | 102 |
| 57 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 4.29 | 4.17 | 0.51 | 2.85 | 40 | 1, 1, .. | 102 |
| 58 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 54.24 | 52.73 | 1.92 | 1.77 | 40 | 1, 1, .. | 102 |
| 59 | 94.44 scanBfsFusedKern… | scanBfsFusedKern… | 264.03 | 249.36 | 2.27 | 2.05 | 40 | 2, 1, .. | 102 |
| 60 | 88.89 scanBfsFusedKern… | scanBfsFusedKern… | 167.97 | 149.30 | 4.68 | 4.22 | 40 | 4, 1, .. | 102 |
| 61 | 94.44 scanBfsFusedKern… | scanBfsFusedKern… | 56.90 | 53.74 | 1.60 | 1.46 | 40 | 2, 1, .. | 102 |
| 62 | 97.22 scanBfsFusedKern… | scanBfsFusedKern… | 6.50 | 6.32 | 0.31 | 12.05 | 40 | 1, 1, .. | 102 |

This table shows all results in the report. Use the column headers to sort the results in this report. Double-click a result to see detailed metrics. Double-click on demangled names to rename it.

Figure 8: bfs result 2

below 6%, with some spikes (e.g., 12.35%) attributed to isolated high-degree nodes. This indicates inefficient memory coalescing and idle threads in warps with divergent workloads. Furthermore, the estimated speedup across kernels is relatively low ( 80.56%), and register usage is light (22 registers/thread), suggesting low instruction-level parallelism and minimal use of shared resources. While the kernel launch overhead is small, the simplistic design fails to saturate the GPU pipeline, leading to subpar performance on large graphs.

In contrast, the queueBfs implementation shows improved memory throughput and slightly better GPU utilization due to its queue-centric load distribution. However, its performance remains inconsistent. The kernel durations vary widely—from 14 µs to over 500 µs—suggesting that load imbalance persists, especially in iterations with uneven frontier sizes. Memory throughput improves in some kernels (up to 12%), yet compute throughput remains low (often below 1%), confirming that the workload is still largely memory-bound. Register pressure is moderate and constant (22 registers/thread), but queue management overhead and atomic operations on global memory limit the performance gains, especially as graph size and vertex degree increase.

The scanBfsFusedKernel variants—particularly those using WM and CM strategies—demonstrate significantly enhanced performance characteristics. Kernels from the WM variant frequently report memory throughput above 10% and compute throughput exceeding 2%, with some peaking near 4.68%, indicating better parallel efficiency. The fused nature of the kernel eliminates inter-kernel synchronization overhead, and warp-level prefix sums improve load balancing across threads within a warp. Register usage is slightly higher (26–28), and the use of shared memory further reduces global memory latency. The CM variant similarly benefits from block-level cooperation and CUB-based prefix sums, producing uniform edge workloads per thread and achieving more consistent throughput. The edge-centric approach in both WM and CM leads to shorter and more efficient execution across BFS levels, especially for graphs with skewed degree distributions.

Overall, the NCU results confirm that kernel fusion, warp/block cooperation, and fine-grained load balancing play a crucial role in improving GPU-based BFS performance. Simple designs like simpleBfs and queueBfs underutilize GPU resources and

suffer from scalability issues, while fused kernels with advanced scheduling demonstrate superior efficiency by maximizing memory bandwidth utilization, reducing divergence, and maintaining high compute throughput across traversal iterations.

**Occupancy**   Kernel occupancy, which reflects the ratio of active warps to the maximum number of schedulable warps on an SM, offers a key perspective on how efficiently each BFS algorithm utilizes GPU resources. In our setup, the theoretical occupancy limit is 66.67%. The observed values show a clear divide between different implementation strategies.

The simpleBfs and queueBfs implementations suffer from generally low and inconsistent occupancy, often well below 50%. This indicates suboptimal thread scheduling and poor scalability, likely due to their vertex-centric and sequential expansion strategies that fail to keep GPU cores busy, especially during sparse frontier phases.

In contrast, scanBfs exhibits a gradual improvement in occupancy across iterations, particularly in the scanDegrees stage, which benefits from more uniform work per thread. However, the occupancy remains relatively low in early stages or when frontiers are small, limiting the potential of multi-kernel designs to fully exploit hardware parallelism.

Fused implementations show notable improvements. The original scanBfsFusedKernel achieves moderate occupancy gains, though still with some variance across iterations. In comparison, the warp-level load-balanced version (scanBfsFusedKernel_WM) maintains consistently high occupancy near the theoretical limit, confirming the effectiveness of warp cooperation in managing irregular workloads. The coalesced mapping variant (scanBfsFusedKernel_CM) also performs well, though slightly below the WM version, reflecting its block-centric but still efficient design.

Overall, occupancy analysis reinforces that advanced load-balancing fused kernels better leverage GPU resources, while naive or sequential strategies underutilize available hardware.

# 5   Conclusion

This study presents a comprehensive evaluation of several GPU-based BFS implementations, ranging from simple vertex-centric strategies to advanced fused and load-balanced kernels. Through detailed experiments across diverse graph datasets—spanning from low-degree road networks to high-degree scale-free social and web graphs—we uncover critical insights into the scalability, efficiency, and architectural alignment of each method. The baseline algorithm, 'simpleBfs', demonstrates reasonable performance on small or regular graphs due to its minimal overhead, but struggles to scale on larger graphs with complex frontier dynamics. Similarly, 'queueBfs', which introduces queue-based load balancing, shows moderate improvements yet continues to suffer from low occupancy and inefficient parallelism under heavy workload variance.

The 'scanBfs' algorithm represents a step forward, employing a prefix-sum based queue generation strategy to better handle varying vertex degrees. However, its multi-kernel design introduces synchronization and memory traffic overheads, limiting its effectiveness on large-scale inputs. To address these shortcomings, we developed a series of fused kernels that combine multiple BFS stages into single-pass implementations.

Among them, 'scanBfsFusedKernel' significantly reduces launch overhead and improves data locality, though its occupancy and load balance remain sensitive to the input structure.

Our advanced designs—'scanBfsFusedKernel_WM' and 'scanBfsFusedKernel_CM' demonstrate the most consistent and scalable performance. The WM variant leverages warp-level cooperation to evenly distribute edge-processing tasks, achieving near-ideal occupancy and delivering strong throughput even under irregular degree distributions. The CM variant, on the other hand, utilizes a block-wide edge-centric strategy, balancing work across threads via shared prefix sums and adaptive indexing, which proves especially effective for graphs with skewed degree patterns.

Comparative analysis with established libraries like Gunrock and GSWITCH further highlights the strength of our fused designs. While GSWITCH achieves impressive performance through graph-specific tuning and load-aware kernels, our implementations match or exceed its efficiency in multiple datasets, offering generality and modularity in return. Occupancy and kernel-level profiling further confirm that maximizing parallel resource utilization—through intelligent scheduling and memory coalescing—is crucial for high-performance graph traversal on GPUs.

In summary, this work demonstrates that tightly fused, load-balanced kernels especially those guided by warp-centric and coalesced mapping principles—represent a highly effective strategy for scaling BFS across irregular graphs on modern GPU architectures. These insights not only inform the design of future GPU graph algorithms but also emphasize the need to consider hardware-level parallelism and memory behavior as first-class design constraints in irregular workload optimization.

# References

[1] Jun-Dong Cho, Salil Raje, and Majid Sarrafzadeh. Fast approximation algorithms on maxcut, k-coloring, and k-color ordering for vlsi applications. *IEEE Transactions on Computers*, 47(11):1234–1243, November 1998.

[2] Thomas Lengauer and Robert Endre Tarjan. A fast algorithm for finding dominators in a flowgraph. *ACM Transactions on Programming Languages and Systems*, 1(1):121–141, January 1979.

[3] P. Narayanan. Single source shortest path problem on processor arrays. In *[Proceedings 1992] The Fourth Symposium on the Frontiers of Massively Parallel Computation*, pages 553–556, 1992.

[4] Pawan Harish and P. J. Narayanan. Accelerating large graph algorithms on the gpu using cuda. In Srinivas Aluru, Manish Parashar, R. Badrinath, and Viktor K. Prasanna, editors, *High Performance Computing – HiPC 2007*, volume 4873 of *Lecture Notes in Computer Science*, pages 197–208. Springer, Berlin, Heidelberg, 2007.

[5] Lijuan Luo, Martin Wong, and Wen mei Hwu. An effective gpu implementation of breadth-first search. In *Proceedings of the 47th Design Automation Conference (DAC '10)*, pages 52–55, New York, NY, USA, 2010. Association for Computing Machinery.

[6] Yangzihao Wang, Andrew Davidson, Yuechao Pan, Yuduo Chen, Yuxiong He, and John D. Owens. Gunrock: A high-performance graph processing library on the gpu. *ACM SIGPLAN Notices*, 50(8):265–266, 2015.

[7] Ke Meng, Jiajia Li, Guangming Tan, and Ninghui Sun. A pattern based algorithmic autotuner for graph processing on gpus. In *Proceedings of the 24th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '19)*, pages 201–213, New York, NY, USA, 2019. Association for Computing Machinery.

[8] Sungpack Hong, Tayo Oguntebi, and Kunle Olukotun. Efficient parallel graph exploration on multi-core cpu and gpu. In *Proceedings of the 2011 International Conference on Parallel Architectures and Compilation Techniques (PACT)*, pages 78–88. IEEE, 2011.