# Compiling Graph Applications for GPUs with GraphIt
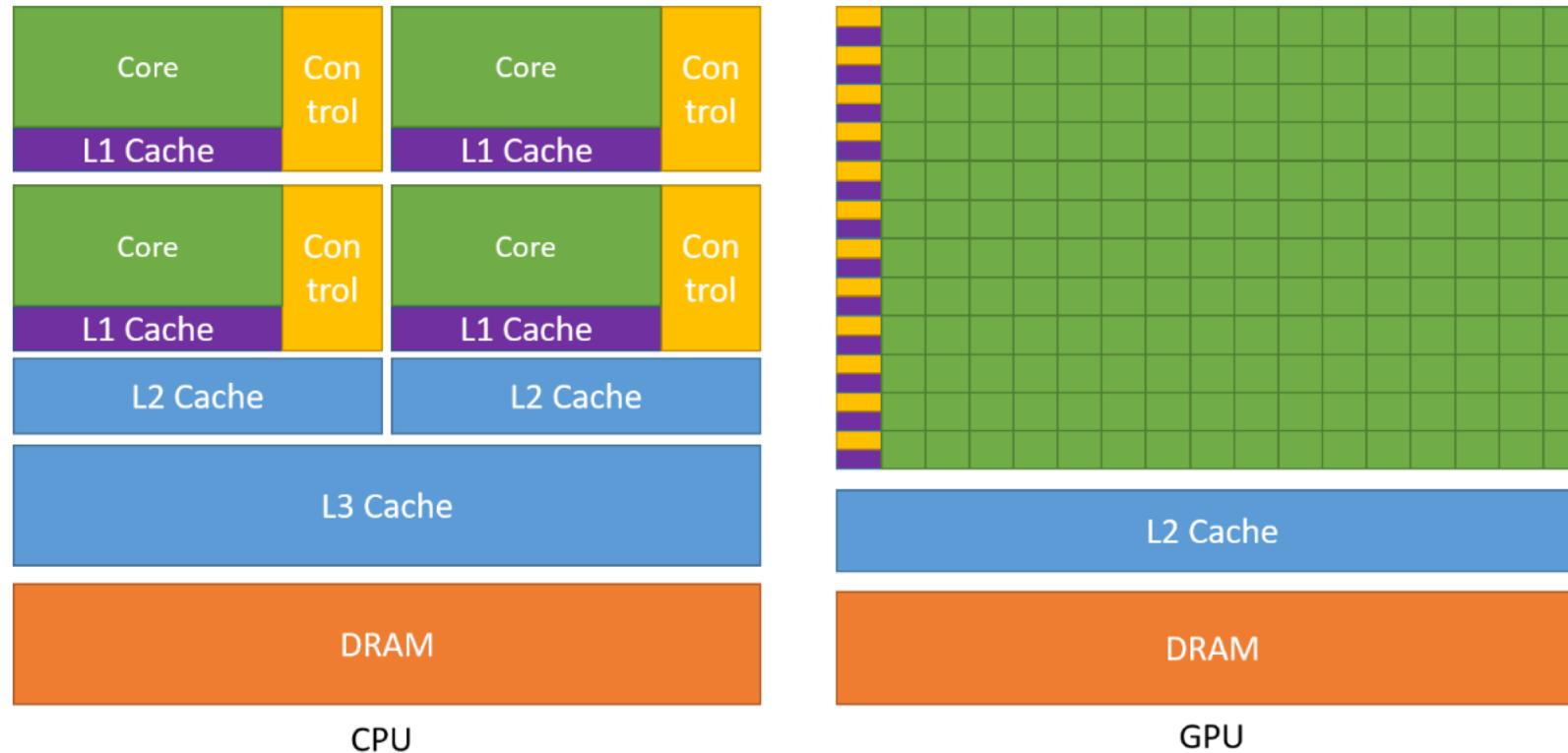
Ajay Brahmakshatriya, Yunming Zhang, Changwan Hong, Shoaib Kamil,
Julian Shun, Saman Amarasinghe
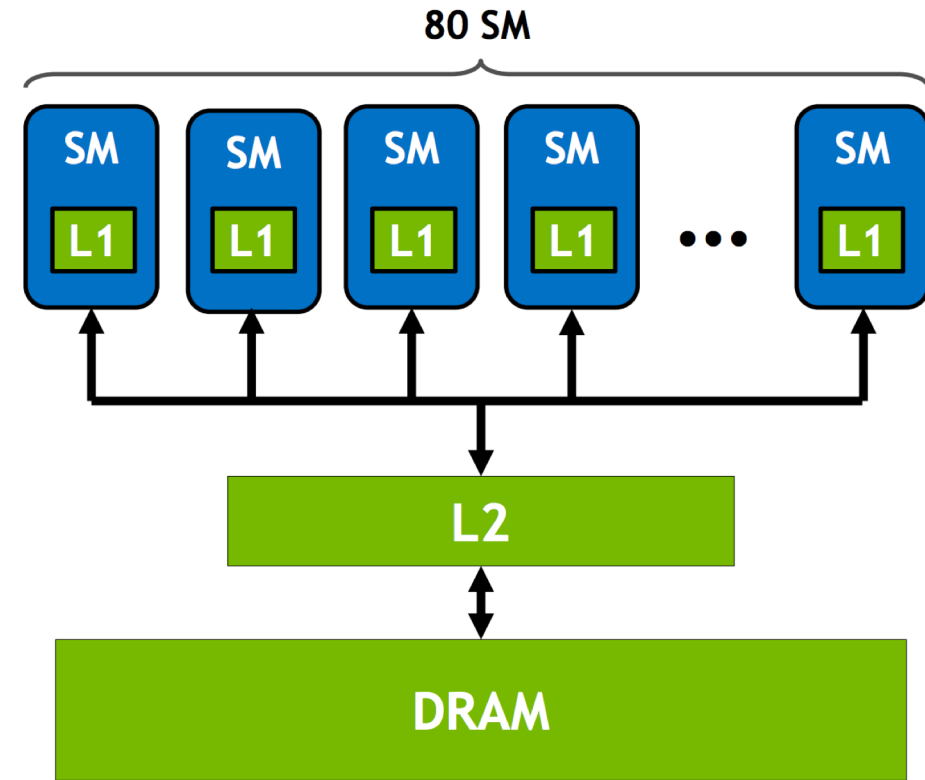
# CPU vs. GPU



Source: https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html

- CPU is designed to execute one thread fast
- GPU is designed to execute many (slower) threads in parallel, achieving higher throughput

# GPU Architecture



- GPU has multiple streaming multiprocessors (SMs)
  - For example, Tesla V100 has 80 SMs

- Work in GPUs is organized into thread blocks (CTAs), and dynamically assigned to SMs

- Each SM has its own data cache, which can be partitioned between L1 cache and shared memory
  - 128KB on Tesla V100

- There is global DRAM (16 or 32 GB on V100) and a global L2 cache (6144 KB on V100)
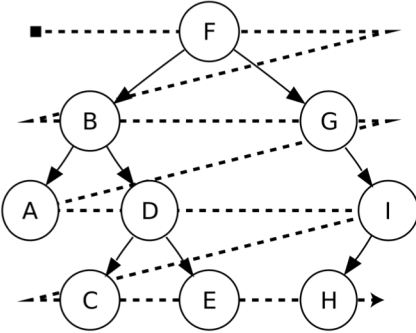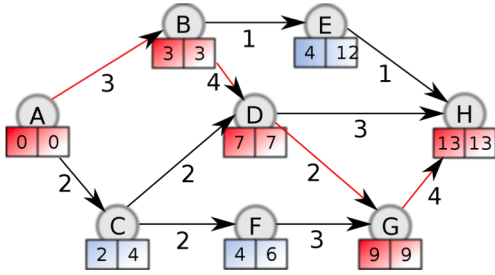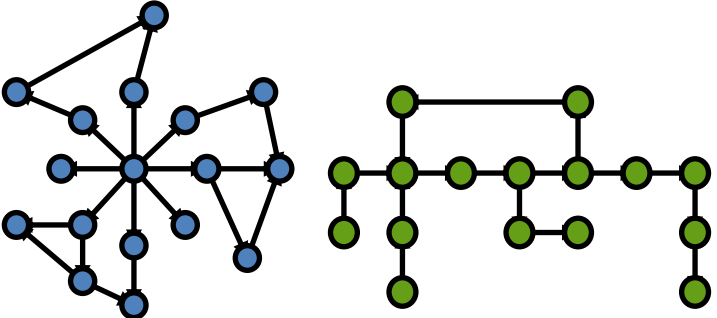
# GPU Architecture

- Each SM schedules warps (groups of 32 threads) to execute the same computation together
  - Tesla V100 has 4 warp schedulers and instruction units, and can execute 4 warps at a time

- Each instruction unit has its own cores for arithmetic, and L0 instruction cache

# Graph Optimization Tradeoff Space



**Graphs**

**Algorithms**

**Optimizations**

**Hardware**

- ## Decouple Algorithm from Optimization
  - Algorithm language: What to compute
  - Scheduling (optimization) language: How to compute

- ## Scheduling representation
  - Easy to use for users to try different combinations of optimizations without changing the algorithm

# GraphIt – A Domain-Specific Language for Graph Analytics



**Algorithm Language**

**Optimization Representation**
- Scheduling Language
- Schedule Representation
  (e.g., Graph Iteration Space)

**Autotuner**

# Hardware Variations?

**PageRank on social networks**



**GPU**



- Massive parallelism
- Coalesced memory accesses

**Vertices Processed**

18   18   18   18



**CPU**



- Limited parallelism
- Large number of iterations

**Vertices Processed**

1   1   1   3

**SSSP on road graphs**

- No existing framework for generating both CPU and GPU code prior to GraphIt

# GraphIt Backend Support



- GraphIt currently has backends for multicore CPUs and GPUs
- First framework to support code generation for both CPUs and GPUs with the same algorithm specification

# Key GPU Optimizations



Direction Optimization???

Load Balancing???

Kernel Fusion???

Vertex Subset Representation??

Edge Blocking

Vertex Subset Deduplication??

# GPU Scheduling Language

- Load Balancing Strategy
- Iteration direction
- Representation of output frontier
- Deduplication of output frontier

- Fusing multiple CUDA kernels
- Graph partitioning for cache utilization
- Runtime combinations of above

```
configLoadBalance(CM|WM|TWC|TWCE|EB|VB|STRICT)
                  configDirection(PUSH|PULL)
configFrontierCreation(FUSED|UNFUSED, BITMAP|
                                       BOOLMAP)
        configDeduplication(ENABLED|DISABLE,
        BITMAP|BOOLMAP|MONOTONIC_COUNTERS)
        configKernelFusion(ENABLED|DISABLED)
       configEdgeBlocking(ENABLED|DISABLED)
                            HybridGPUSchedule
```

# Comparison of Optimizations with Other Frameworks

| Optimization | Gunrock | GSWITCH | SEP-Graph | GraphIt |
|---|---|---|---|---|
| Load Balancing | VERTEX BASED, EDGE BASED, TWC | CM, WM, TWC, STRICT | VERTEX BASED | ETWC, TWC, STRICT, CM, WM, VERTEX BASED, EDGE BASED |
| Edge Blocking | Not supported | Not supported | Not supported | Supported |
| Vertex Set Creation | Fused/Unfused | Fused/Unfused | Fused | Sparse Queue / Bitmap / Boolean Array |
| Direction Optimization | Push/Pull/Hybrid | Push/Pull/Hybrid | Push/Pull/Hybrid | Push/Pull/Hybrid |
| Deduplication | Supported | Not supported | Supported | Supported |
| Vertex Ordering | Supported | Supported | Supported | Supported |
| Kernel Fusion | Supported | Not supported | Supported | Supported |
| Total combinations | 48 | 32 | 16 | **576** |

# New Optimization: ETWC Load Balancing

- Edge-based Thread Warps CTAs load-balancing (ETWC)
  - First, equally partitions vertices across CTAs
  - Then, partitions edges of a vertex into low, medium, high bins to be processed by a thread, warp, and entire CTA, respectively
  - Trades off load balance for reducing load balancing overhead

| Graph | ETWC | TWC | CM |
|-------|--------|--------|--------|
| OK | 43.58 | **40.69** | 42.24 |
| TW | **106.11** | 107.57 | 116.06 |
| LJ | 19.72 | 20.03 | **18.42** |
| SW | **226.35** | 230.00 | 230.03 |
| HW | **4.94** | 5.79 | 8.17 |
| IC | **11.38** | 11.50 | 22.16 |
| RU | **136.64** | 255.89 | 168.90 |
| RC | **91.20** | 162.54 | 109.89 |
| RN | **13.10** | 25.77 | 16.25 |

Times (ms) of ETWC on breadth-first search, compared with existing strategies TWC and CM. Fastest time is **bolded**.

# New Optimization: Edge Blocking

- Edge blocking (EB)
  - Tiles edges into subgraphs such that the random accesses for each subgraph fit in L2 cache
  - Process each subgraph one at a time



| Graph | Without EB | With EB | Speedup |
|---|---|---|---|
| OK | 41.75 | **14.18** | 2.94x |
| TW | 88.25 | **77.86** | 1.13x |
| LJ | 15.67 | **7.68** | 2.04x |
| SW | 144.88 | **102.11** | 1.41x |
| HW | **7.01** | 7.02 | 0.99x |
| IC | **18.24** | 19.55 | 0.93x |
| RU | 8.35 | **6.32** | 1.35x |
| RC | 8.39 | **5.56** | 1.50x |
| RN | 0.44 | **0.43** | 1.02x |

Times (ms) and speedup of
Edge blocking (EB) on PageRank

# GraphIt Compilation



- Whole program analysis / transformations

# Example: Frontier Reuse Analysis

```
...
while (frontier.getVertexSetSize() != 0)
    output = edges.from(frontier).to(toFilter).
        applyModified(updateEdge, parent);
    delete frontier;
    frontier = output;
  end
...
```

```
while (frontier.getVertexSetSize() != 0) {
    cudaMalloc(output, ...);
    ApplyModified<<<, >>>(frontier, output, …);
    ...
    cudaFree(frontier);
    frontier = output;
}
```

**Allocations and freeing on GPUs are costly unlike CPUs**

- Do we really need to allocate and free on every round?
- We should reuse the memory that was allocated for `frontier`
- Is it always safe to do so?
  - What if the old frontier is used again?

**Liveness Analysis!**

# Example: Frontier Reuse Analysis

- Constructs live ranges for each Vertexsubset variable (used to represent frontiers)

```
...
while (frontier.getVertexSetSize() != 0)
    output = edges.from(frontier)…applyModified(…);
    delete frontier;
    frontier = output;
end
...
```



■ Live Range of frontier

■ Live Range of output

- Disjoint live ranges allows memory to be reused

# Experimental Evaluation

- Compare to 3 state-of-the-art frameworks (Gunrock [Wang et al. 2017], GSWITCH [Meng et al. 2019], SEP-graph [Wang et al. 2019])

- 5 algorithms: breadth-first search, single-source shortest paths (Delta-stepping), connected components, betweenness centrality, and PageRank

- 9 datasets: social networks, Web graphs, and road networks

- 2 generations of NVIDIA GPUs: Pascal and Volta

# State of the Art and GraphIt on Titan Xp (Pascal) GPU

Slowdowns relative to the fastest implementation

## GraphIt

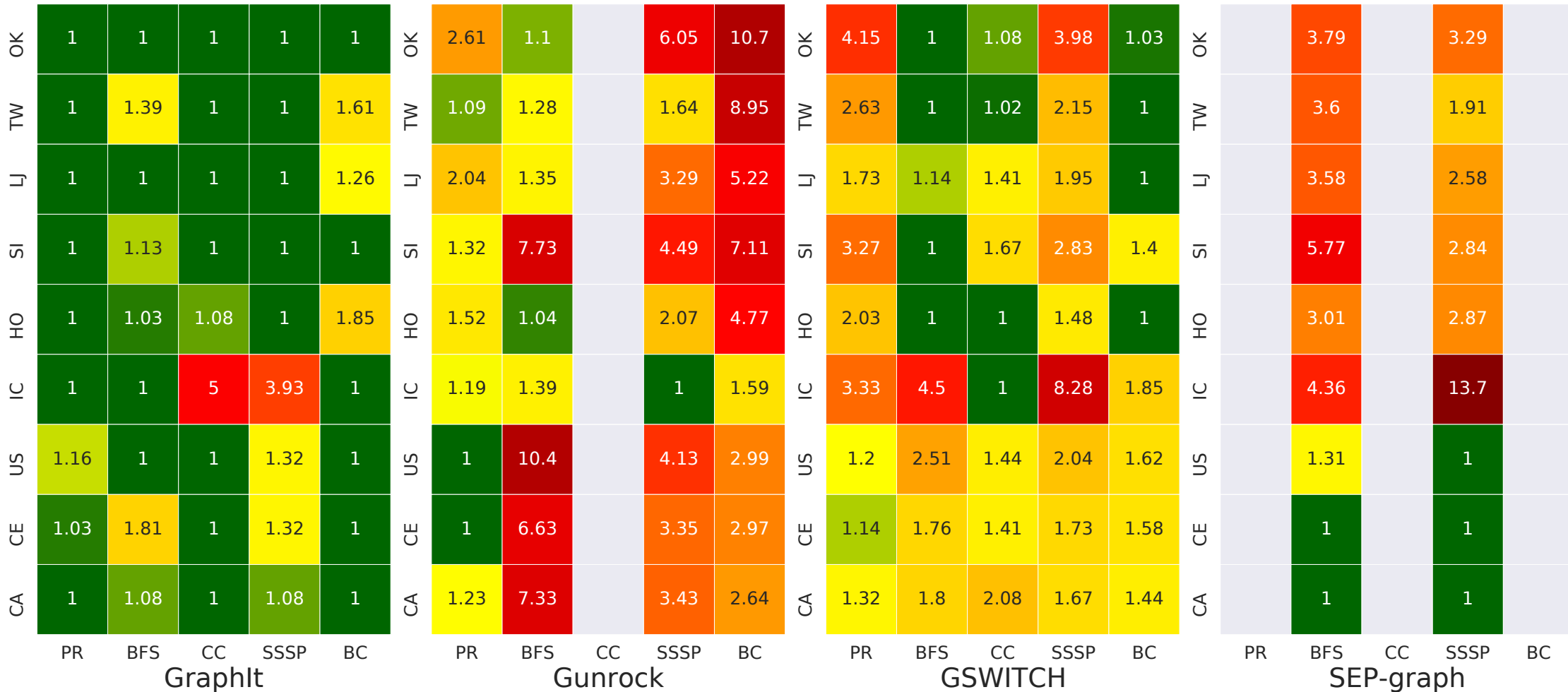| | PR | BFS | CC | SSSP | BC |
|---|---|---|---|---|---|
| OK | 1 | 1 | 1 | 1 | 1 |
| TW | 1 | 1 | 1 | 1 | 1.42 |
| LJ | 1 | 1.27 | 1 | 1 | 1.07 |
| SI | 1 | 1.06 | 1 | 1 | 1 |
| HO | 1 | 1 | 1 | 1 | 1 |
| IC | 1.96 | 1 | 1 | 1 | 3.34 |
| US | 1 | 1 | 1 | 1.37 | 1 |
| CE | 1 | 1 | 1 | 1.18 | 1 |
| CA | 1 | 1 | 1 | 1.03 | 1 |

## Gunrock

| | PR | BFS | CC | SSSP | BC |
|---|---|---|---|---|---|
| OK | 4.28 | 1.1 | 1.14 | 10.4 | 8.16 |
| TW | 1.46 | | 1.9 | 2.31 | 4.11 |
| LJ | 2.3 | 1.3 | 1.46 | 4.78 | 3.26 |
| SI | 1.75 | | 1.59 | 5.06 | 5.36 |
| HO | 3.23 | 1.08 | 3.11 | 4.06 | 2.38 |
| IC | 1.42 | 1.49 | 7.45 | 1.93 | 4.48 |
| US | 1.67 | 6.04 | 3.6 | 122 | 3.27 |
| CE | 1.79 | | 1.78 | 103 | 2.64 |
| CA | 2.19 | | 3.38 | 5.41 | 3.59 |

## GSWITCH

| | PR | BFS | CC | SSSP | BC |
|---|---|---|---|---|---|
| OK | 8.26 | 1.11 | 1.21 | 5.84 | 1.45 |
| TW | 2.71 | 1.07 | | 2.04 | 1 |
| LJ | | 1 | 1.13 | 3.18 | 1 |
| SI | 3.32 | 1 | | 2.82 | 2.03 |
| HO | | 1.22 | 1.51 | 2.58 | 6.51 |
| IC | 1 | 6.62 | 1.36 | 2.23 | 1 |
| US | 1.21 | 1.64 | 1.51 | 2.81 | 1.87 |
| CE | 1.59 | 1.71 | 1 | 2.25 | 1.39 |
| CA | 1.09 | 1.68 | 1.77 | 1.69 | 1.64 |

## SEP-graph

| | PR | BFS | CC | SSSP | BC |
|---|---|---|---|---|---|
| OK | | 3.66 | | 4.61 | |
| TW | | 1.9 | | 2.07 | |
| LJ | | 3.16 | | 3.16 | |
| SI | | 4.94 | | 3.35 | |
| HO | | 2.52 | | 5.04 | |
| IC | | 5.16 | | 4.25 | |
| US | | 1.15 | | 1 | |
| CE | | 1.24 | | 1 | |
| CA | | 1.12 | | 1 | |

# State of the Art and GraphIt on V-100 (Volta) GPU

Slowdowns relative to the fastest implementation



GraphIt achieves a speedup of up to 5.11x due to searching through a much larger space of optimizations

# Programmability

- Lines of code for each algorithm in each framework

| Algorithm | Gunrock | GSWITCH | SEP-Graph | GraphIt (Algorithm+Schedule) |
|---|---|---|---|---|
| Breadth First Search | 2189 | 164 | 481 | **66** |
| PageRank | 2207 | 159 | - | **61** |
| Connected Components | 3014 | 160 | - | **62** |
| Betweenness Centrality | 1792 | 280 | - | **128** |
| SSSP with Delta Stepping | 1438 | 203 | 473 | **50** |

# CPU vs. GPU

- Compared GPU implementations with CPU implementations in GraphIt on a 24-core machine
  - PageRank, BFS, betweenness centrality, and connected components were faster on the GPU
  - Delta-Stepping on road graphs was faster on the CPU
  - CPUs can process much larger graphs

- It is critical to be able to choose between CPU and GPU for each application!

# Summary

- GraphIt DSL and compiler to generate high-performance CPU and GPU code from the same high-level algorithm representation

- New GPU-specific scheduling language options and optimizations

- The GPU algorithms from GraphIt outperform state-of-the-art GPU frameworks while requiring fewer lines of code

- Open source: https://graphit-lang.org