

# Groute: An Asynchronous Multi-GPU Programming Model for Irregular Computations

Tal Ben-Nun<sup>†</sup>, **Michael Sutton**<sup>†</sup>, Sreepathi Pai<sup>‡</sup>, and Keshav Pingali<sup>‡</sup>

<sup>†</sup> The Hebrew University of Jerusalem, Israel

<sup>‡</sup> University of Texas at Austin

PPoPP 2017, Austin, TX

# Outline

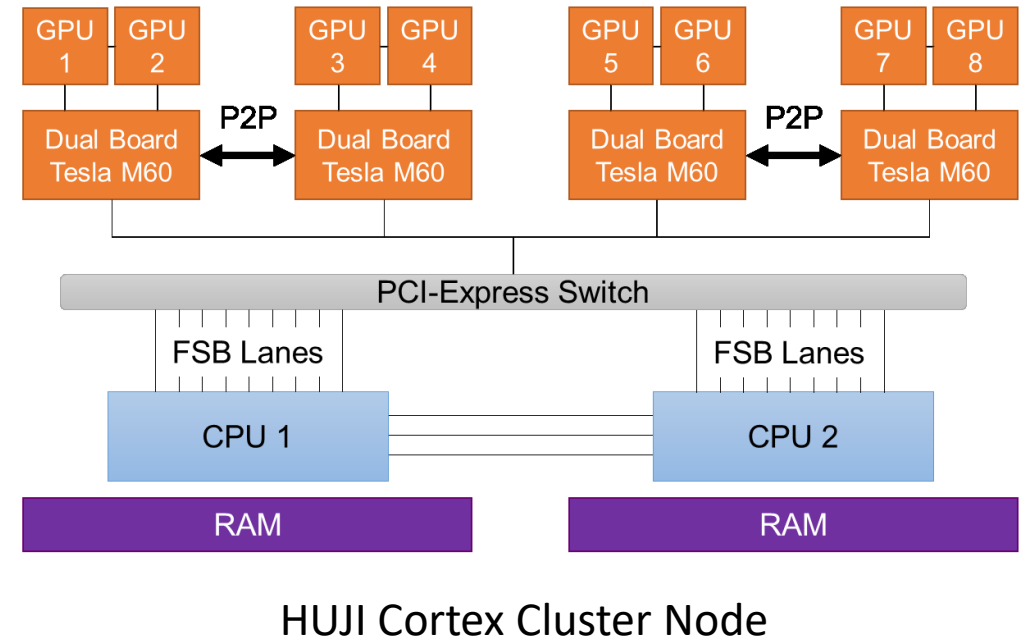
- Motivation
- Groute
  - Programming model
  - Implementation
- Experimental Evaluation
- Conclusions

# Outline

- **Motivation**
- Groute
  - Programming model
  - Implementation
- Experimental Evaluation
- Conclusions

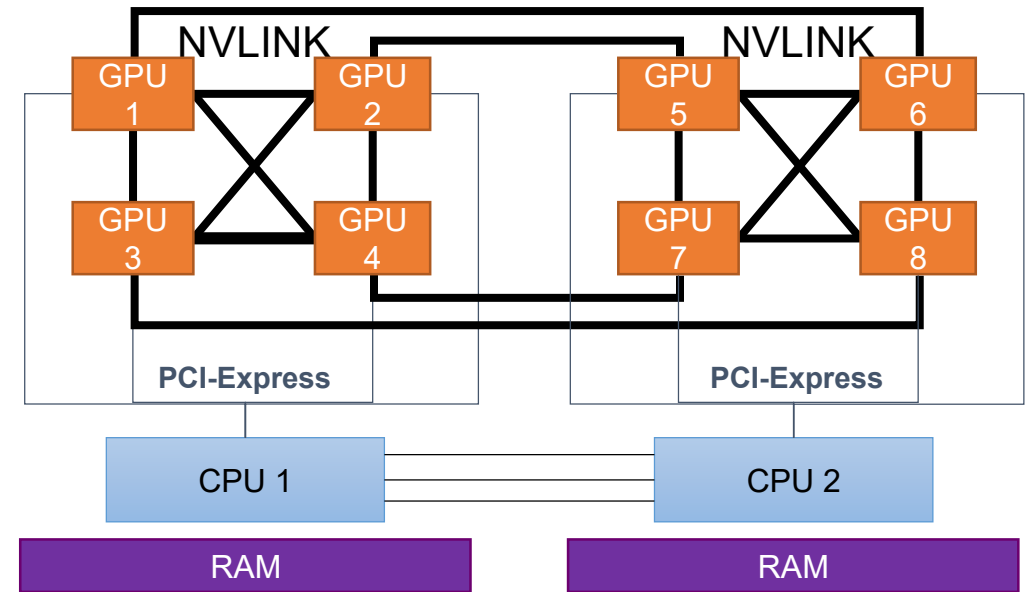
# Multi-GPU Architecture

- Small number of devices (~1-25)
- Low-latency communication
- Large unified memory space
- Wide variation in topology



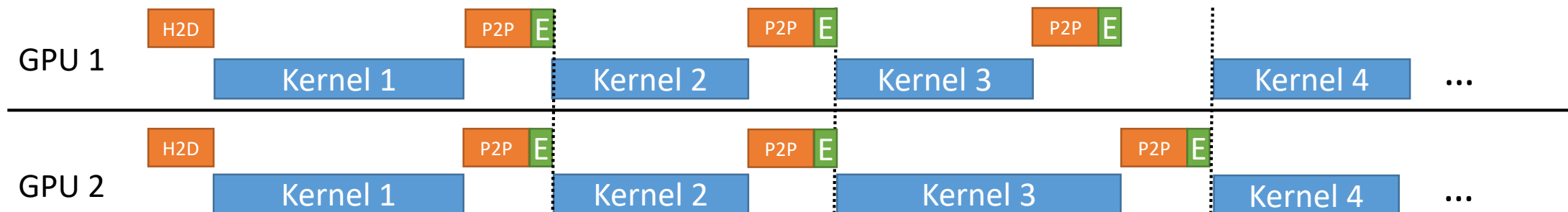
# Multi-GPU Architecture

- Small number of devices (~1-25)
- Low-latency communication
- Large unified memory space
- Wide variation in topology



# Multi-GPU Programming

- Two synchronization constructs: *streams* and *events*
- Current frameworks use BSP-style programming
  - Global synchronization
  - Peer memory transfers: MGPU, MAPS-Multi, Gunrock
  - Inter-GPU direct memory access: Back40Computing (B40C)

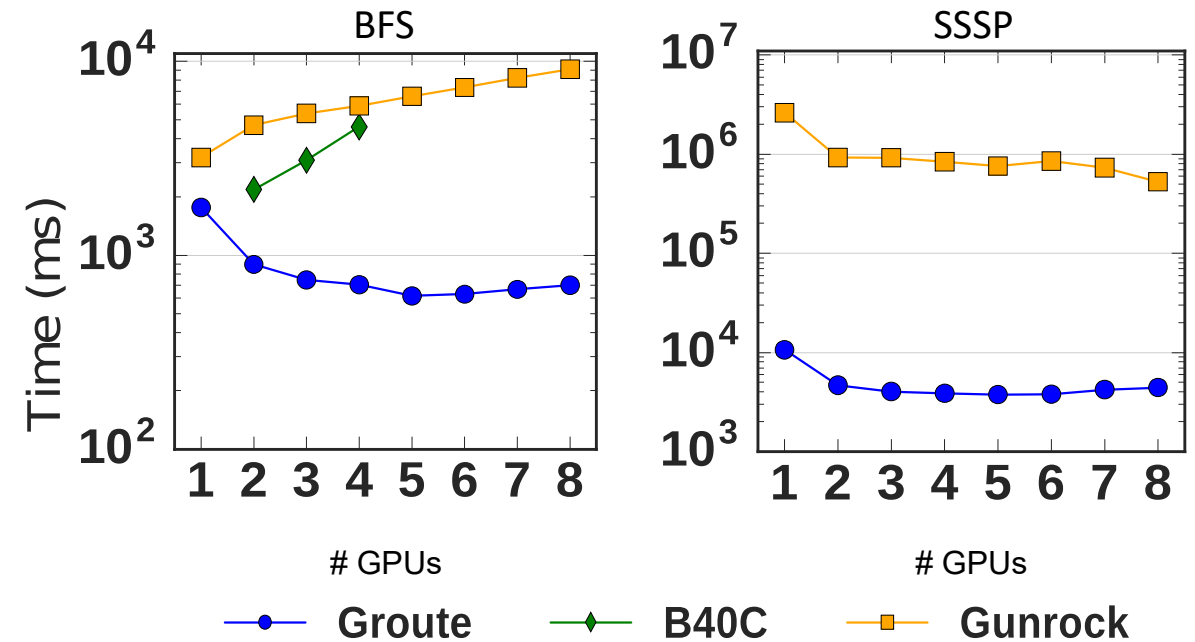


# Irregular Applications

- Characterized by unpredictable memory transfers
- Global synchronization leads to utilization gaps
- Direct memory access does not scale well
  - See paper for micro-benchmark results

# Asynchronous Multi-GPU Programming

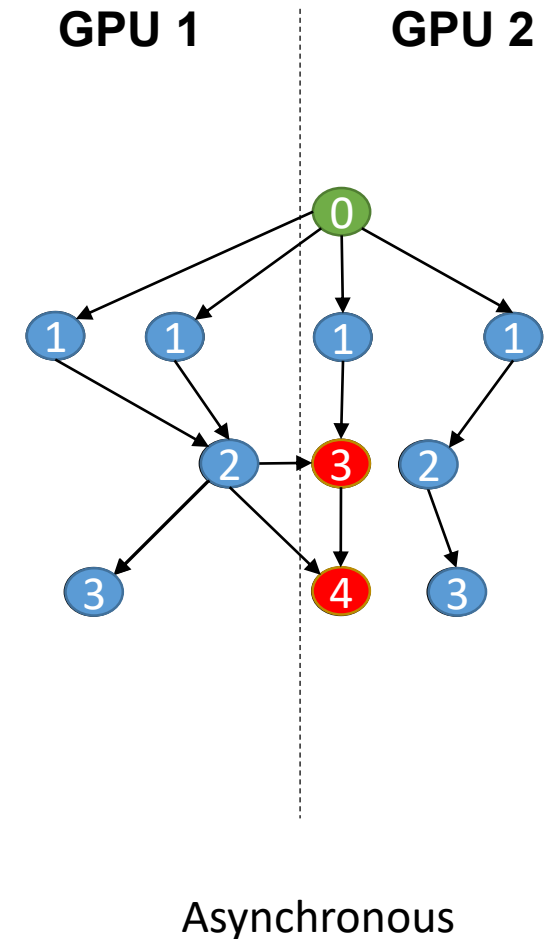
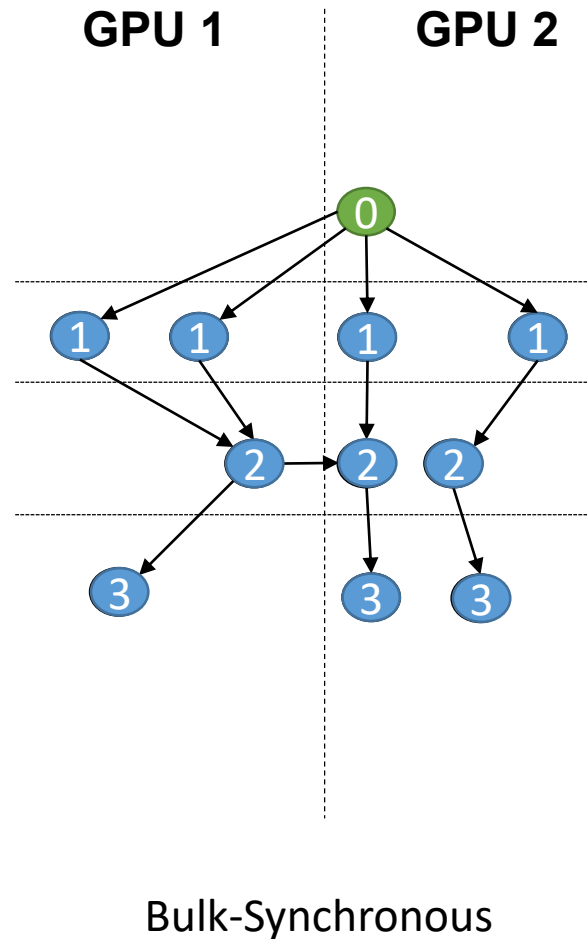
- Asynchronous multi-GPU programming can eliminate utilization gaps
- Requires:
  - Fine-grained synchronization
  - Interleaved communication
  - Device availability indication
- Entry barrier for asynchronous programming is very high, even for skilled developers





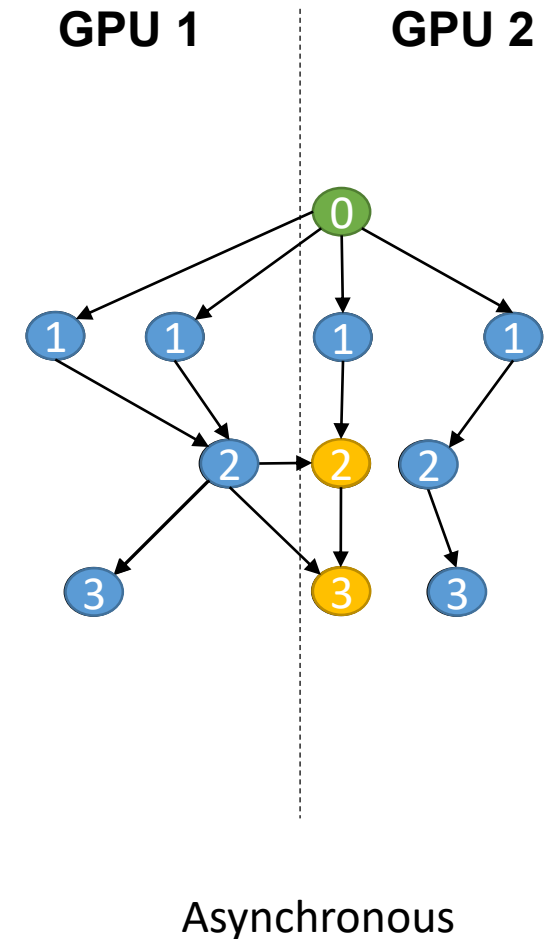
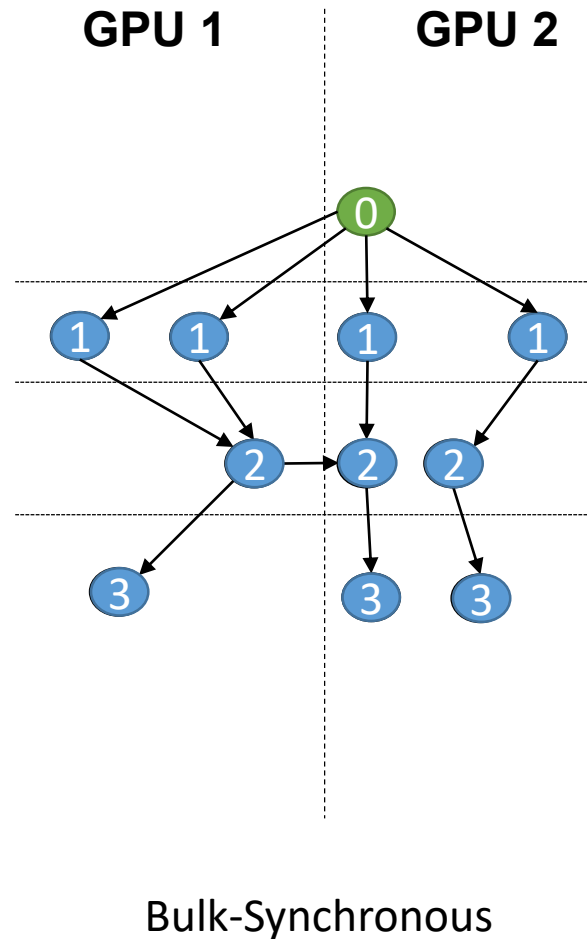
# Asynchronous Scaling Challenges

- Work explosion:  
Asynchronous processing can generate “useless” work
- Causes:
  - Lack of global state management
  - Some devices may be faster
  - Problem structure and partitioning
- May execute billions of additional work items



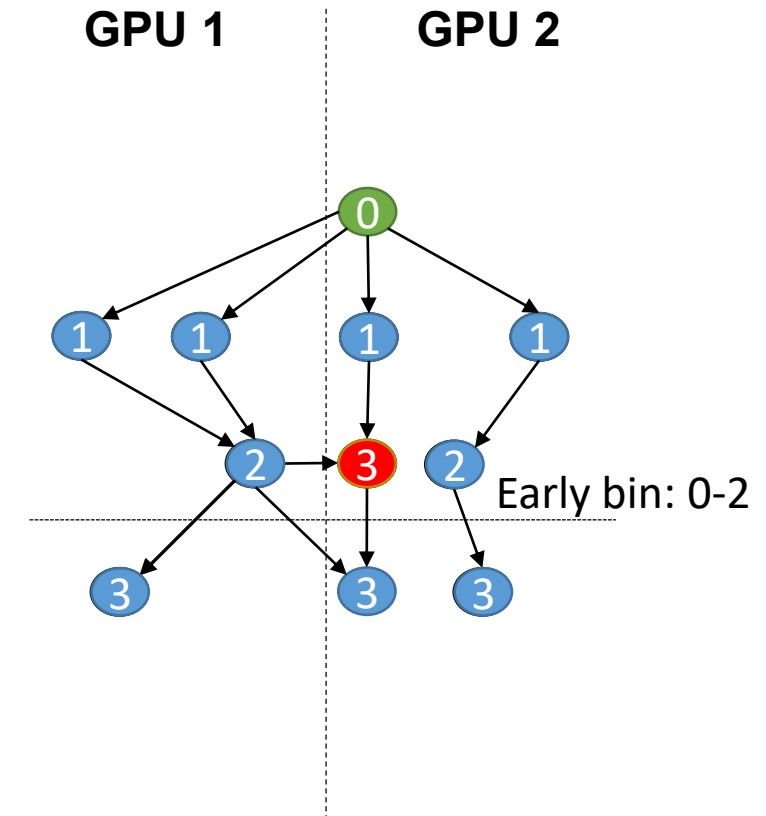
# Asynchronous Scaling Challenges

- Work explosion:  
Asynchronous processing can generate “useless” work
- Causes:
  - Lack of global state management
  - Some devices may be faster
  - Problem structure and partitioning
- May execute billions of additional work items



# Soft Priority Scheduling

- Addresses the work explosion problem
- Assign priorities to work-items
- Bin work-items by priority into early and late bins
- Complete early bin before advancing to late bin
  - Like SSSP delta-stepping
- Requires global consensus on advancing to late bin

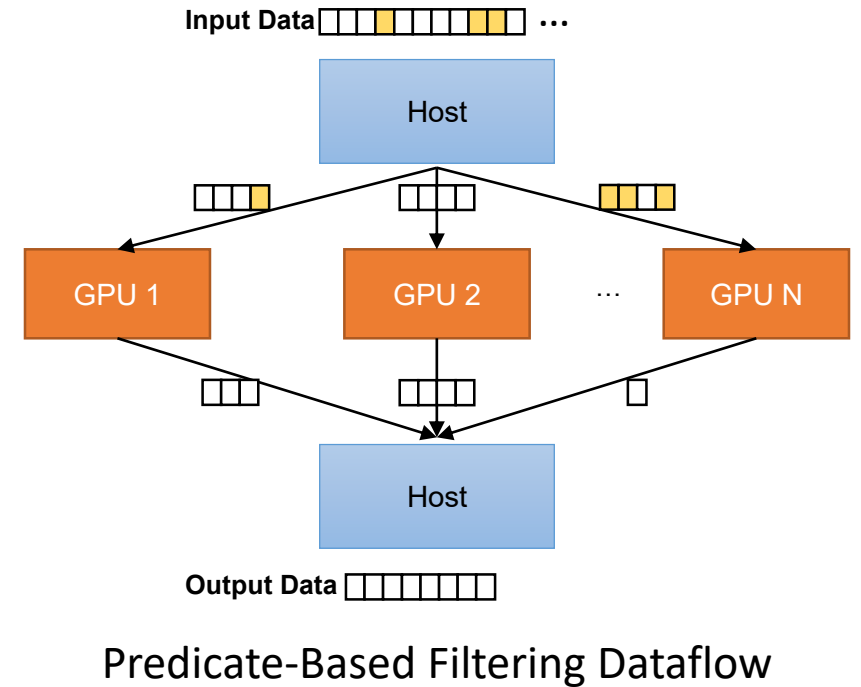


# Outline

- Motivation
- **Groute**
  - Programming model
  - Implementation
- Experimental Evaluation
- Conclusions

# Predicate-Based Filtering (PBF)

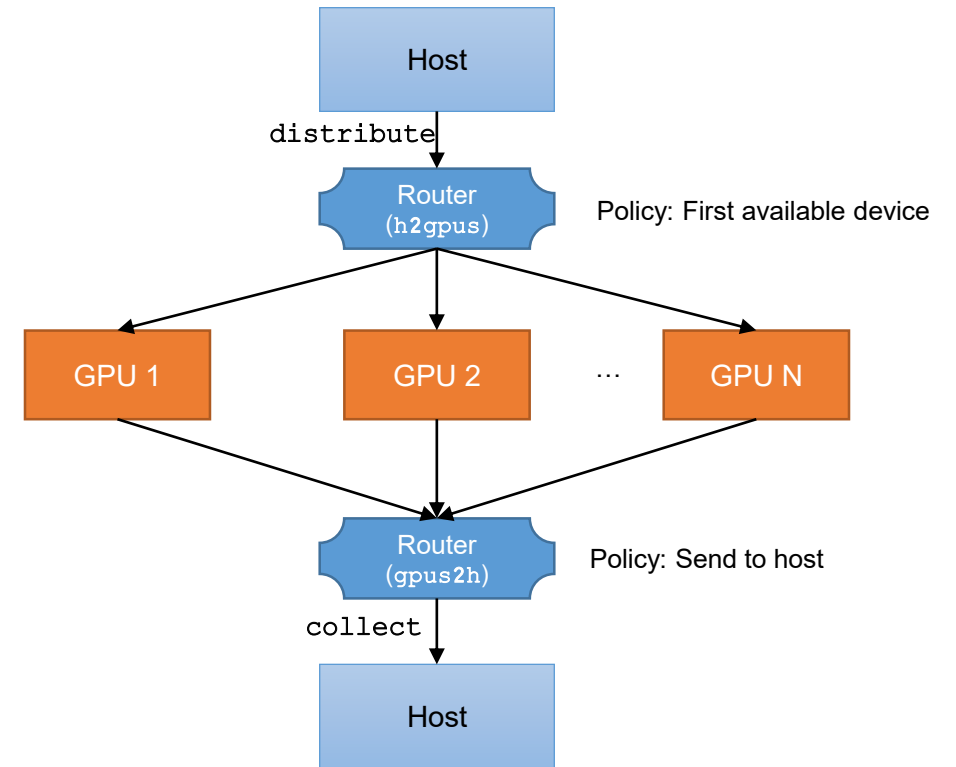
- Host scatters segments of input array to GPUs
- GPUs filter elements that match predicate
- Host gathers filtered data from GPUs
- Challenge: load balancing GPU work



# Groute Programming Model

Groute programs consist of two parts:

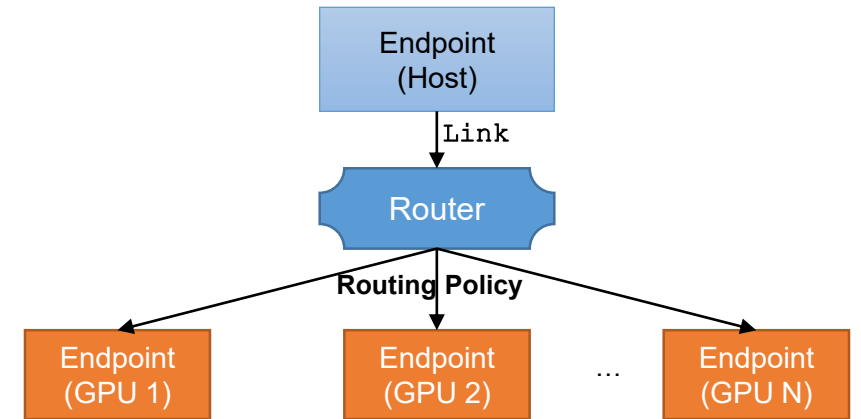
- **Dataflow graph construction**
  - Create links and routers
  - Create management threads
- **Asynchronous execution**
  - Send data through links
  - Process received data asynchronously



Predicate-Based Filtering

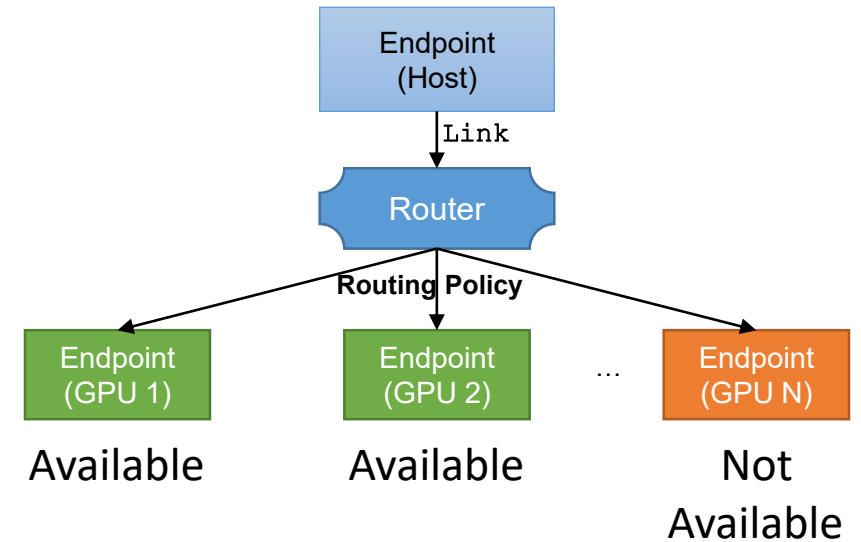
# Groute Dataflow Graph Construction

- **Endpoint:** An entity that can communicate
- **Link:** Connects two Endpoints, data is transferred in pipelined packets
- **Router:** Connects multiple Endpoints for dynamic communication
- **Routing Policy:** Determines destination based on source and contents



# Route Routers

- Routers enable runtime decisions based on policy and availability
- Routing policies express different semantics:
  - Scatter, reduce, all-to-all
- Policy details can be tuned with respect to topology
- Load balancing in PBF implemented by routers

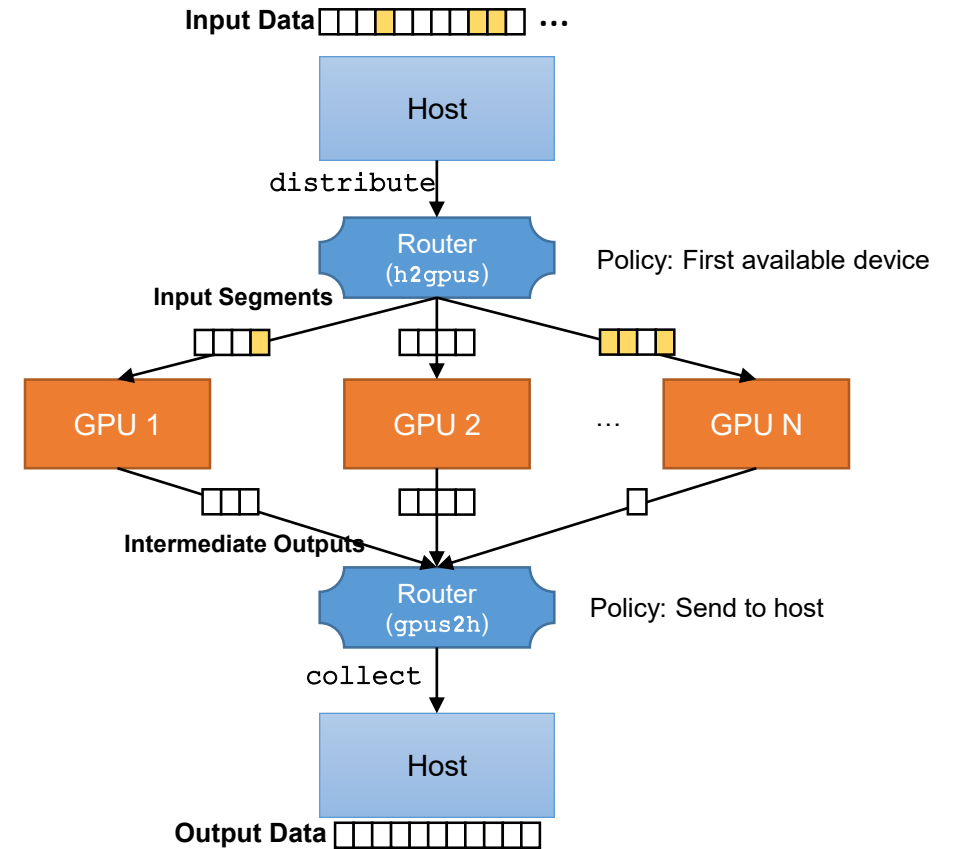




# PBF: Asynchronous Execution

```
void HostThread() {
    std::vector<T> input = ...;
    std::vector<T> output;
    dist.Send(input, input_size);
    dist.Shutdown();
    while(true) {
        PendingSegment output_seg =
            collect.Receive().get();
        if(output_seg.Empty()) break;
        output_seg.Synchronize();
        append(output, output_seg);
        collect.Release(output_seg);
    }
}
```

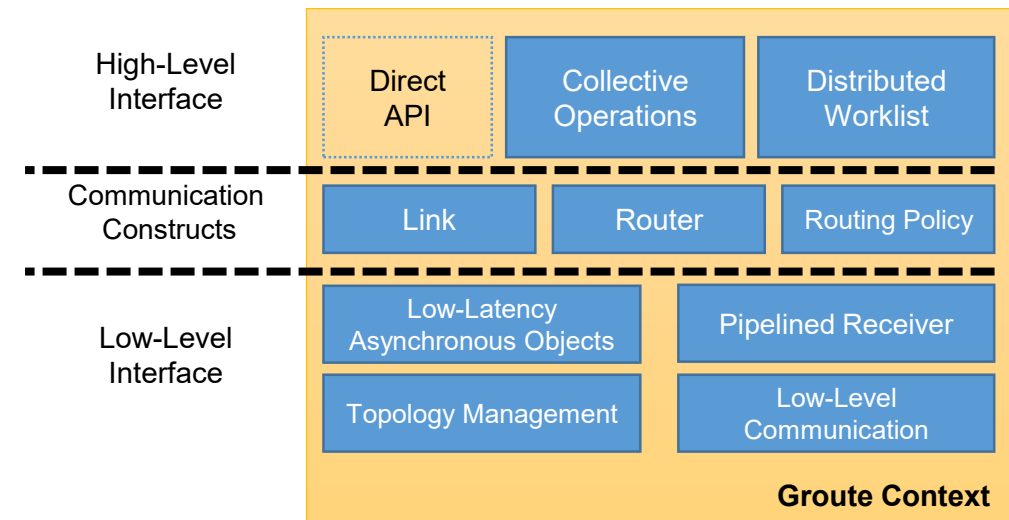
```
void WorkerThread(device_t dev, Link in,
                  Link out) {
    Stream stream(dev);
    T *s_out = ...;
    int *out_size = ...;
    while(true) {
        PendingSegment seg =
            in.Receive().get();
        if(seg.Empty()) break;
        seg.Synchronize(stream);
        Filter<<<..., stream>>>(seg.Ptr(),
                                seg.Size(), s_out, out_size);
        in.Release(seg, stream);
        out.Send(s_out, out_size, stream);
    }
    out.Shutdown();
}
```



Predicate-Based Filtering

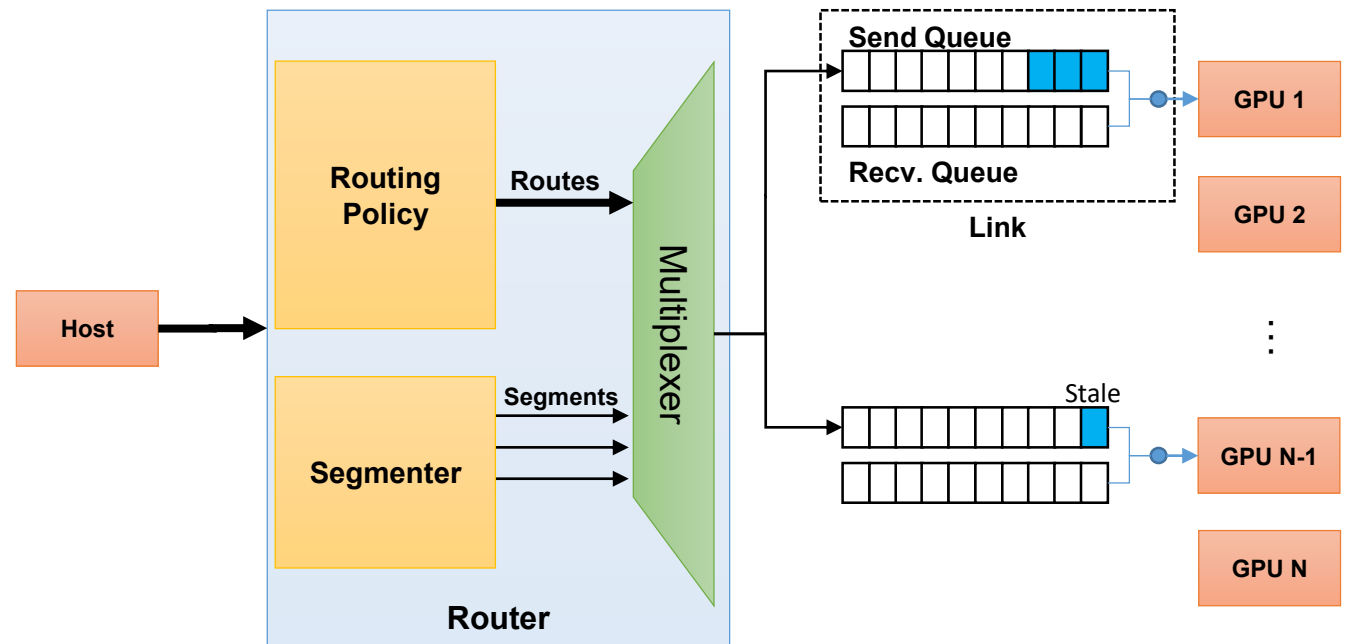
# Groute Implementation

- Three-layered runtime environment
  - Each layer can be skipped for increased programmer control
- Bus topology aware (e.g. PCIe, NVLINK)
- Standard C++11 over CUDA, open source



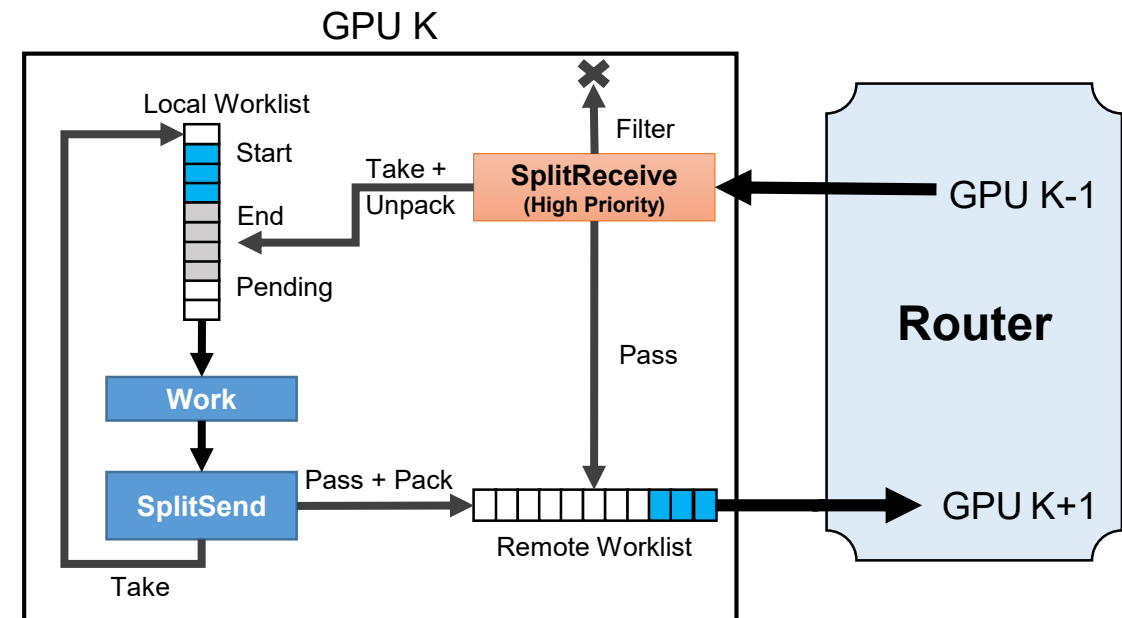
# Router Implementation

- Links enable asynchronous memory transfers
  - Pipelining, workspace allocation, availability indication
- Router packetizes messages and determines destination devices
- Scalable implementation
  - No global locking



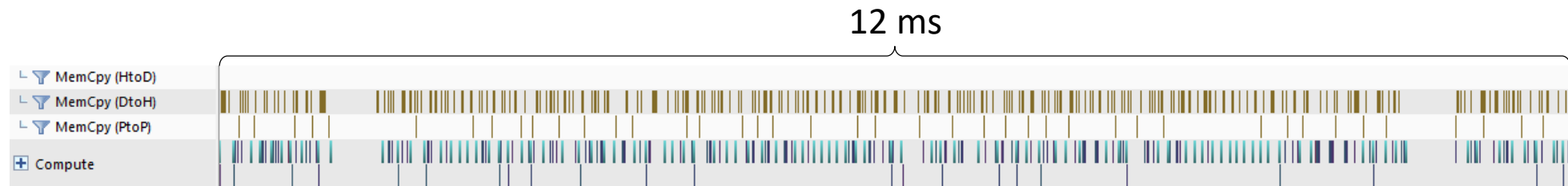
# Worklist-Based Graph Algorithms

- Graph nodes partitioned between GPUs, halo added automatically
- Distributed worklist
  - Implemented over router using ring topology
  - Enables fine-grained all-to-all communication
  - Customized through programmer provided callbacks
  - Supports soft-priority scheduling

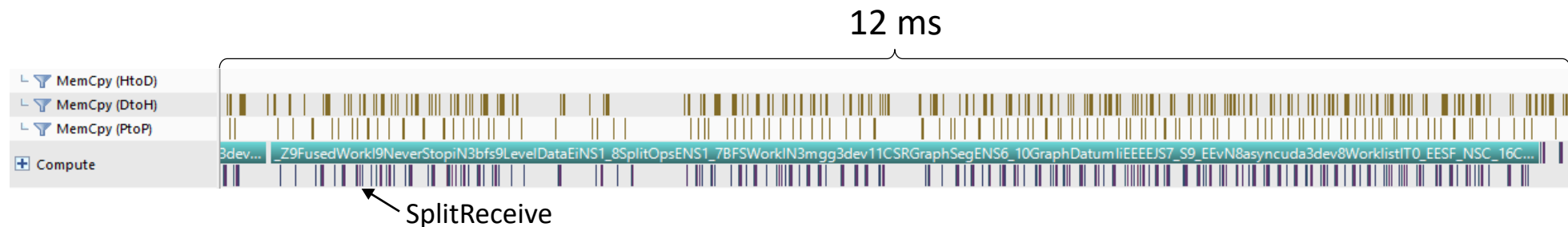


# Distributed Worklist – Kernel Fusion

- Problem: Short kernel lifespan
  - Overhead incurred by kernel launches, CPU-GPU roundtrips



- Solution: Worker Fusion
  - Fuse control-flow, work-item processing and communication
  - One multiprocessor is kept free for receiving data



# Outline

- Motivation
- Groute
  - Programming model
  - Implementation
- **Experimental Evaluation**
- Conclusions

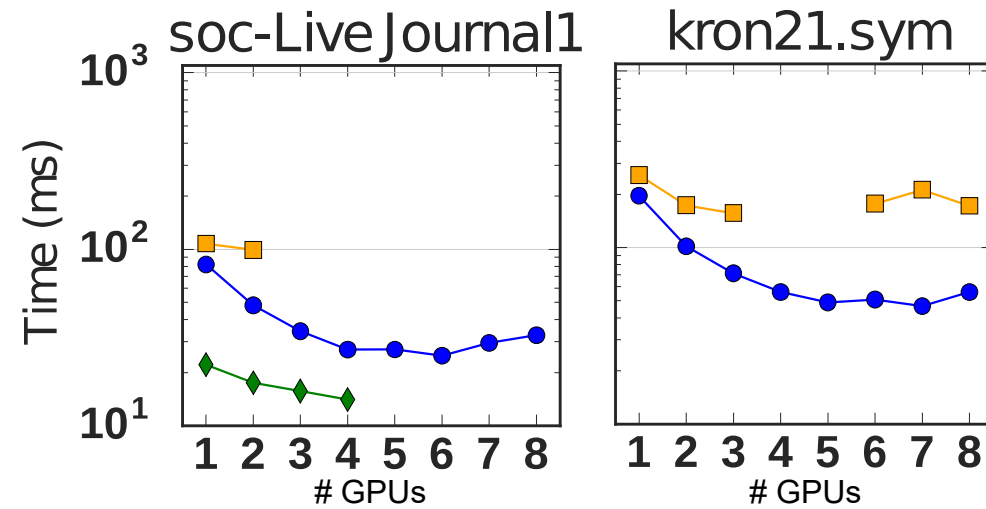
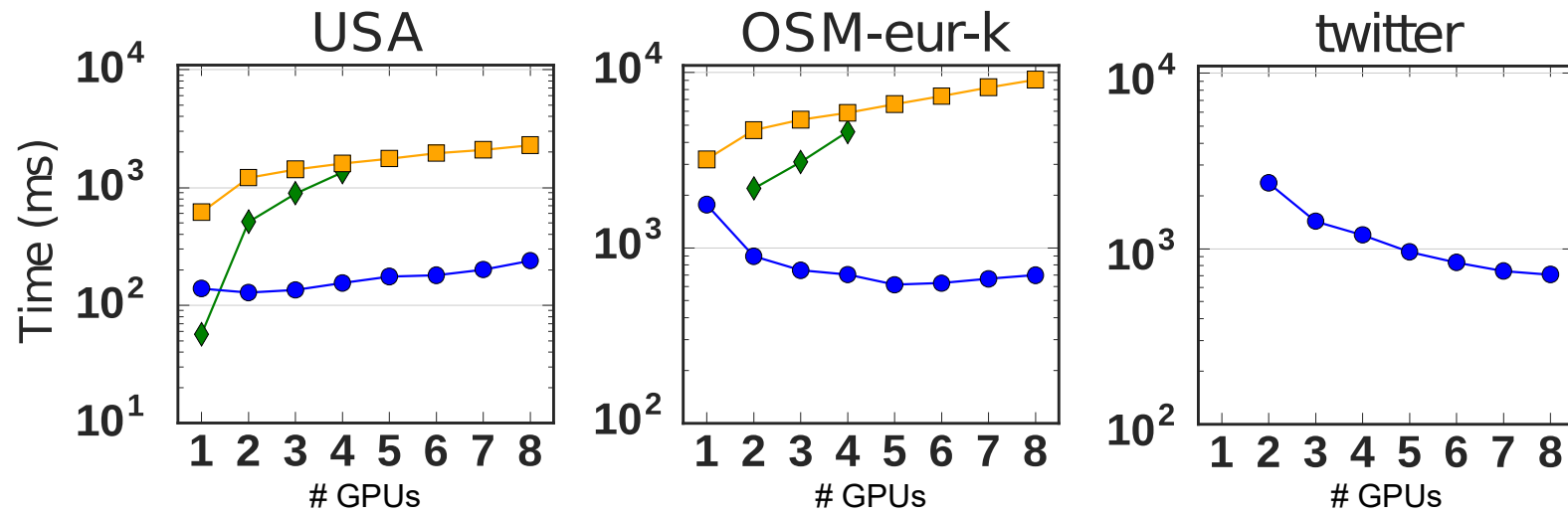
# Experimental Setup

- Systems:
  - HUII Cortex Cluster (8 Tesla M60 GPUs per node)
  - Heterogeneous node with a Quadro M4000 and a Tesla K40c
- Algorithms:
  - Breadth-First Search (data-driven)
  - Single-Source Shortest Path (delta-stepping)
  - PageRank (data-driven, push-based)
  - Connected Components (pointer-jumping)
  - Predicate-Based Filtering
- Datasets:
  - Road Networks: USA, OSM Europe
  - Social Networks: LiveJournal, Twitter
  - Synthetic Graphs: Kronecker Graph (logn21)

# Breadth-First Search Performance

- Implementation is compared with:
  - Gunrock: Frontier-based bulk-synchronous
  - B40C: hardwired BFS implementation, uses direct memory access, limited to 4 GPUs

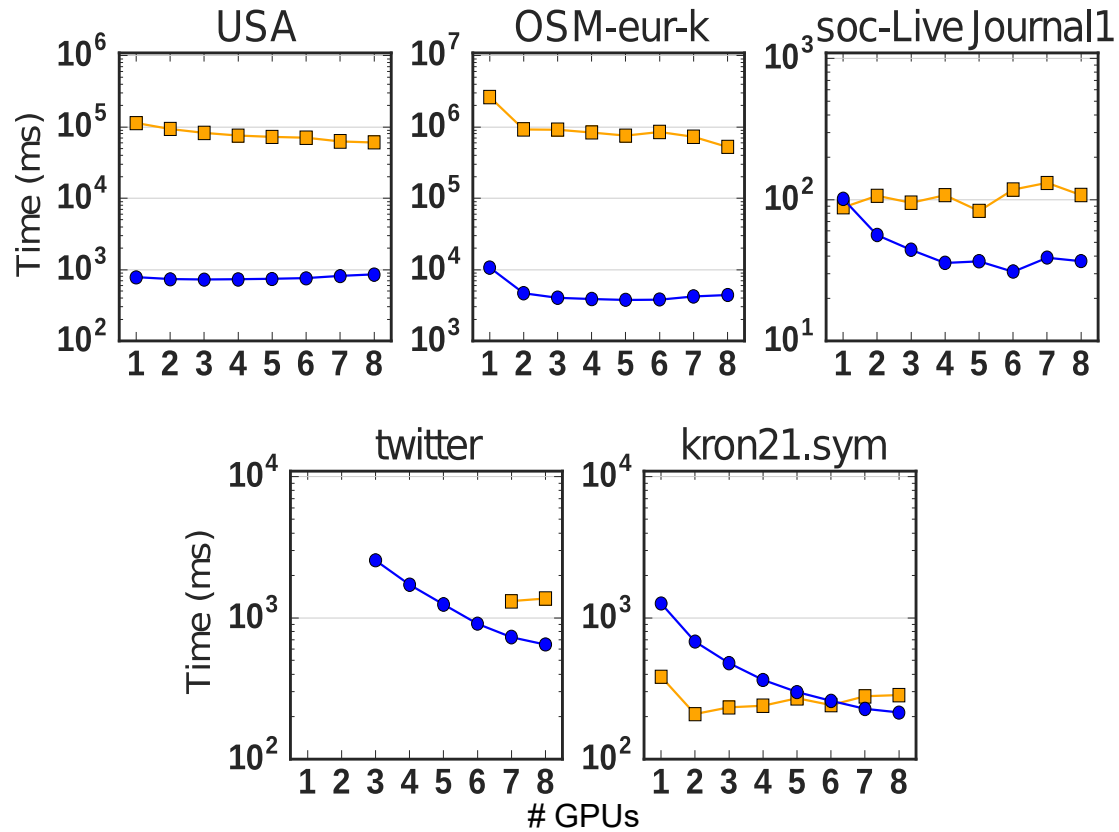




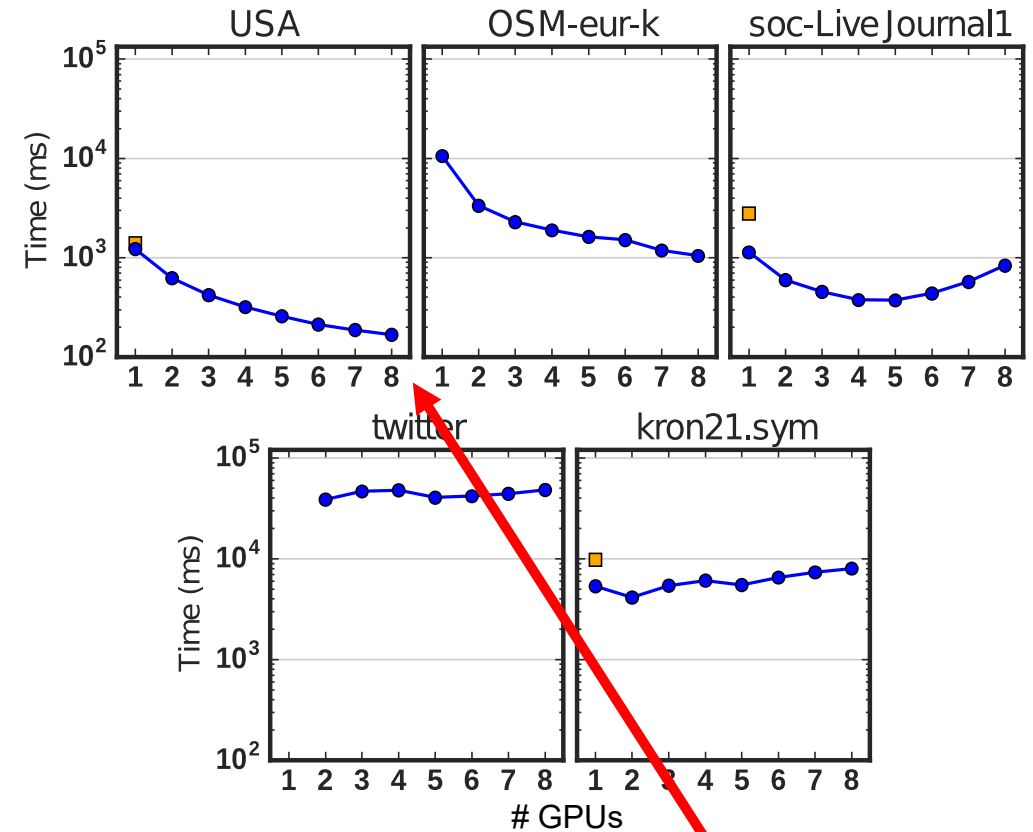
—●— Groute    —◆— B40C    —■— Gunrock

# SSSP and PageRank Performance

## Single-Source Shortest Path



## PageRank

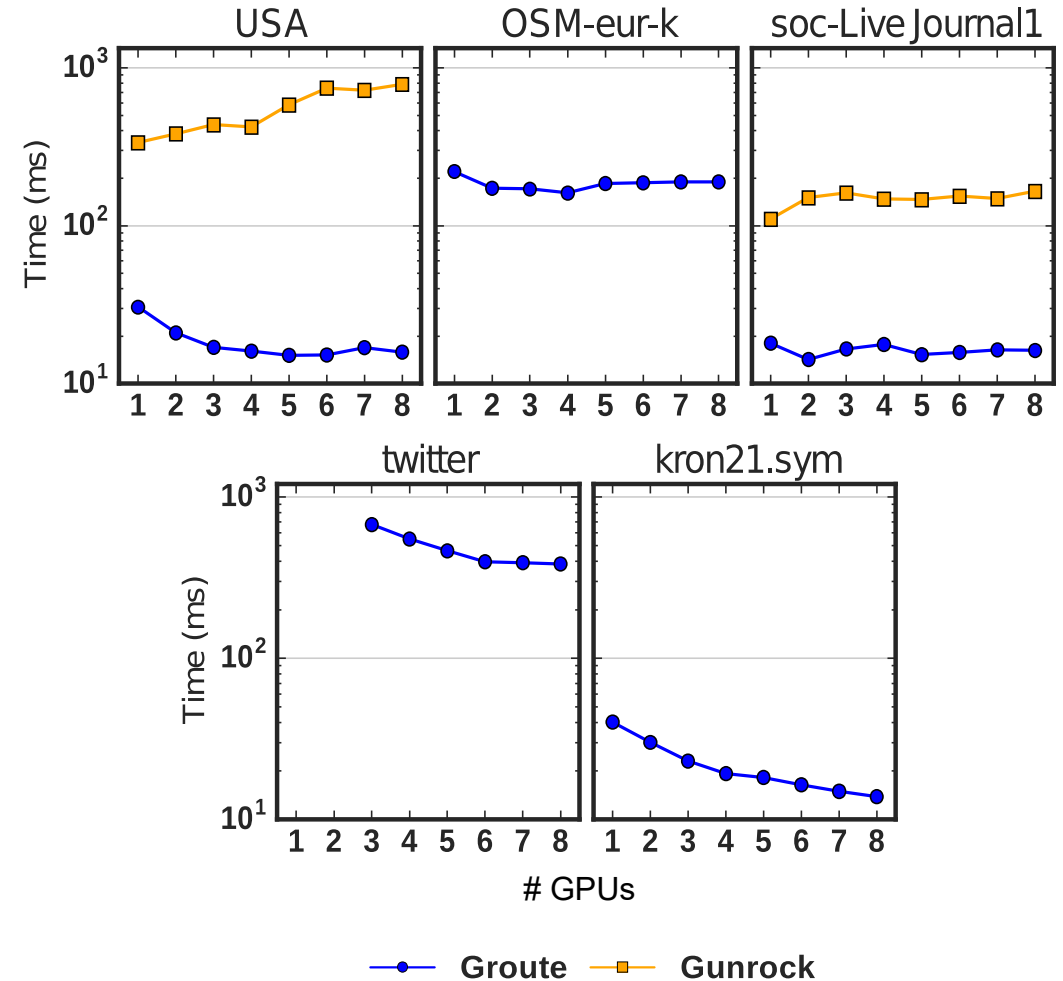
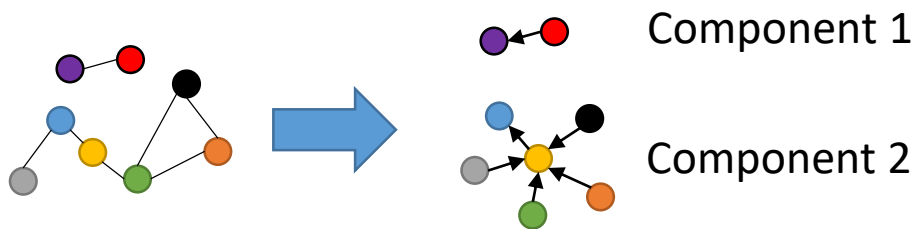


● Groute    ■ Gunrock

7.28x speedup!

# Connected Components Performance

- Pointer-jumping algorithm
  - *Not* a vertex program
  - Topology-driven, highly irregular
- Based on “compression” of graph to depth-1 trees:



# Performance Summary

Best Runtime:

Graph	BFS [ms]			SSSP [ms]		PR [ms]		CC [ms]	
	Gunrock	B40C	Groute	Gunrock	Groute	Gunrock	Groute	Gunrock	Groute
USA	617.85	<b>56.83</b>	128.38	60,656.91	<b>725.93</b>	1,394.25	<b>167.89</b>	335.65	<b>15.11</b>
OSM-eur-k	3,191.78	2,177.1	<b>616.4</b>	874,083.5	<b>3,513.29</b>	—	<b>1,045.33</b>	—	<b>160.96</b>
soc-LiveJournal1	99.11	<b>14.07</b>	24.96	83.36	<b>30.98</b>	2,782.06	<b>371.71</b>	110.05	<b>14.19</b>
Twitter	—	—	<b>713.6</b>	1,310.7	<b>649.2</b>	—	<b>38,549.27</b>	—	<b>384.13</b>
Kron21.sym	156.68	—	<b>46.55</b>	<b>208.43</b>	213.92	9,800.43	<b>5,342.73</b>	—	<b>13.86</b>

Best Groute Scaling:

Graph	BFS	SSSP	PR	CC
USA	1.08x	1.07x	<b>7.28x</b>	1.93x
OSM-eur-k	2.86x	2.83x	3.19x <i>(over 2 GPUs)</i>	1.36x
soc-LiveJournal1	3.28x	3.27x	3.05x	1.27x
Twitter	3.32x <i>(over 2 GPUs)</i>	3.93x <i>(over 3 GPUs)</i>	0.95x <i>(over 2 GPUs)</i>	1.75x <i>(over 3 GPUs)</i>
Kron21.sym	<b>4.21x</b>	<b>5.93x</b>	1.29x	<b>2.90x</b>

# Outline

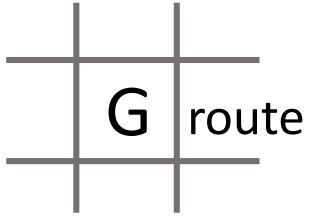
- Motivation
- Groute
  - Programming model
  - Implementation
- Experimental Evaluation
- **Conclusions**

# Conclusions

- Groute programming model is simple and expressive
- Framework hides per-platform tuning
- High-level constructs overcome asynchronous programming pitfalls
- Exceeds state-of-the-art performance

# Future Research

- Programming model application to distributed systems
- Dynamic vertex ownership in traversal algorithms for load balancing
- Complete device-level communication management



Groute is open-source and available at:  
<http://www.github.com/groute/groute>

# Thank you

Questions?