

# McNettle: High Performance Centralized Network Control in a High-Level Language

Yale

Andreas Voellmy Junchang Wang

Yale University University of Science and Technology of China

## MOTIVATION

A centralized SDN controller provides many benefits, but may become a *scaling bottleneck*.

A data center controller may need to provide control for [Tavakoli]: 100s to 1000s of switches, 2 million virtual machines (hosts), and 20 million new TCP sessions per second, but existing sequential controllers cannot meet this performance.

Scaling controller through parallel and concurrent programming may substantially complicate controller programming.

## MCNETTLE: SCALABLE FUNCTIONAL NETWORK CONTROL ON MULTICORES

**Hardware:** Single-node multicore NUMA servers with modern multiqueue NICs.

**Software:** *McNettle API* allows users to write control algorithms in Haskell, a high-level functional, garbage-collected programming language. *McNettle Engine* executes on single-node multicore NUMA servers, automatically scaling performance through 40 cores.

## PROGRAMMING IN MCNETTLE

Each switch is controlled by a user-specified *Switch Controller* accepting messages from the switch and generating messages to switches.

Highly consistent global state supported through shared memory:

- Nonblocking data structures for high performance.
- Software Transactional Memory (STM) for complex shared state manipulation.

Multi-switch state setup supported by allowing one switch controller to setup state in another switch.

## EXAMPLE 1: FINE-GRAINED FLOW MANAGEMENT WITH GLOBAL STATE

```
main = do server ← startServer params
         tbl ← newHostLocTable
         forever (do (features, switch) ← acceptSwitch server
                   startSwitch switch (switchProg tbl))

switchProg tbl msg =
  case msg of
    PacketIn pkt →
      do mport ← lookupAndLearn tbl pkt
         case mport of
           Nothing → forward pkt flood
           Just port → forwardWithExactRule pkt (phyPort port)
    _ → return ()
```

## EXAMPLE 2: BANDWIDTH RESERVATIONS WITH STM

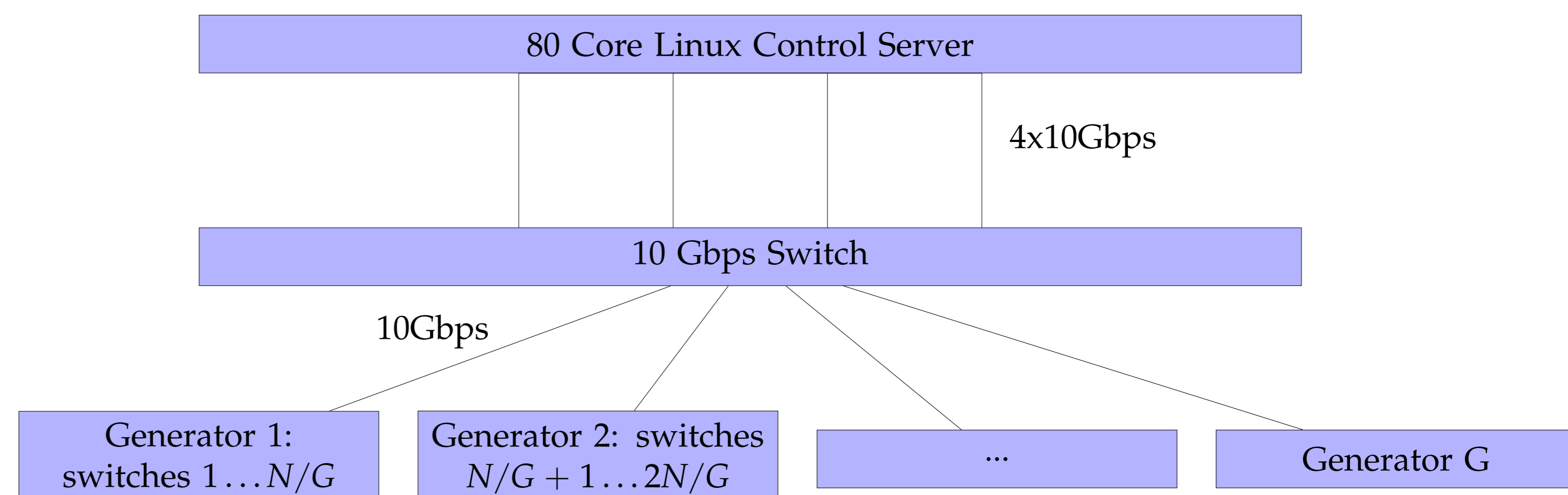
```
atomically (reservePath capTable amt path)
where
  reservePath capTable amt path
    = forM path (\link → reserveLink capTable amt link)
  reserveLink capTable amt link =
    do current ← linkCapacity capTable link
       let remaining = current - amt
         when (remaining < 0) retry
         updateCapacity capTable link remaining
```

## EXAMPLE 3: PARALLEL DIJKSTRA ROUTING

Compute routing for each switch in parallel:

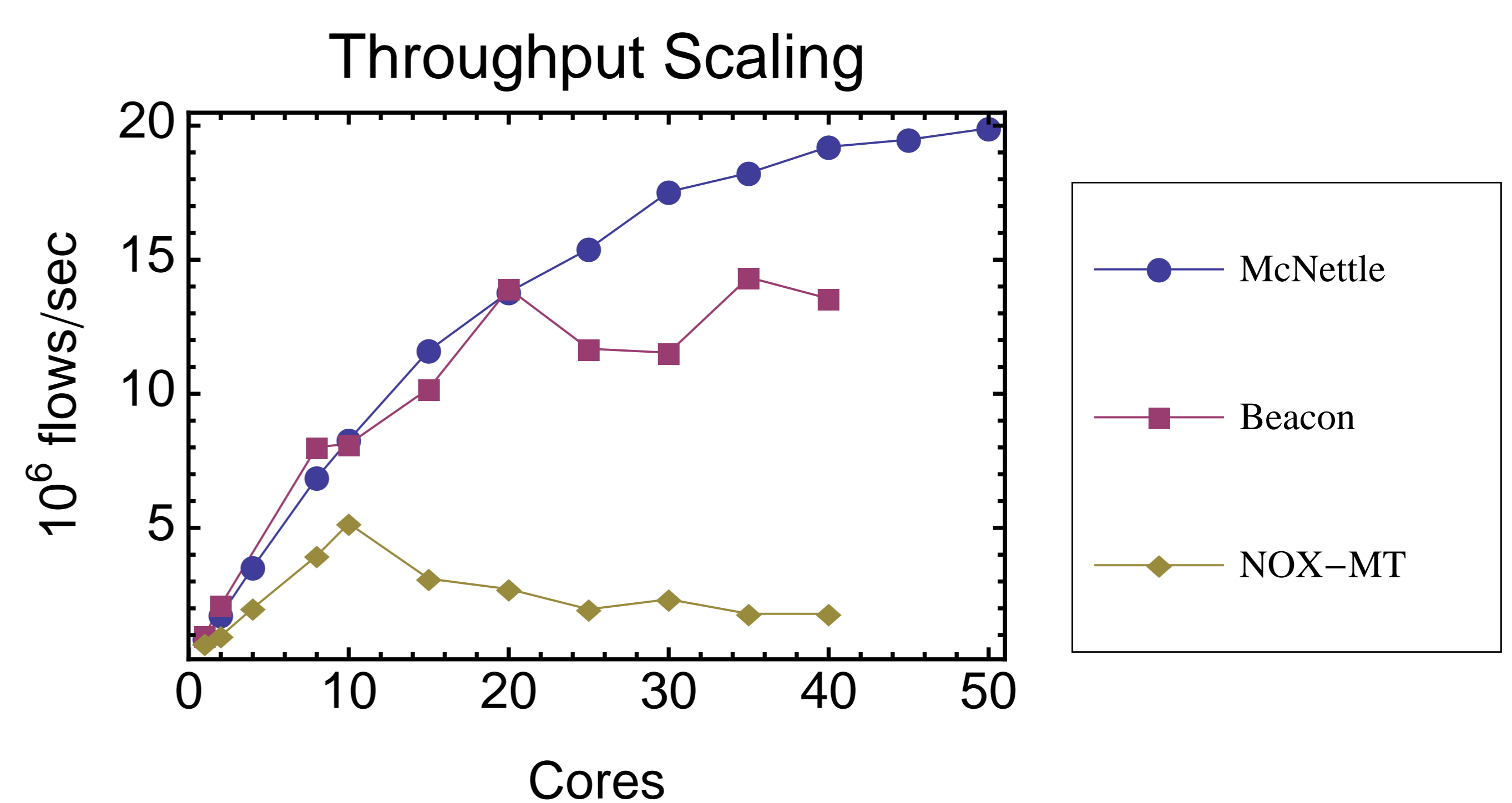
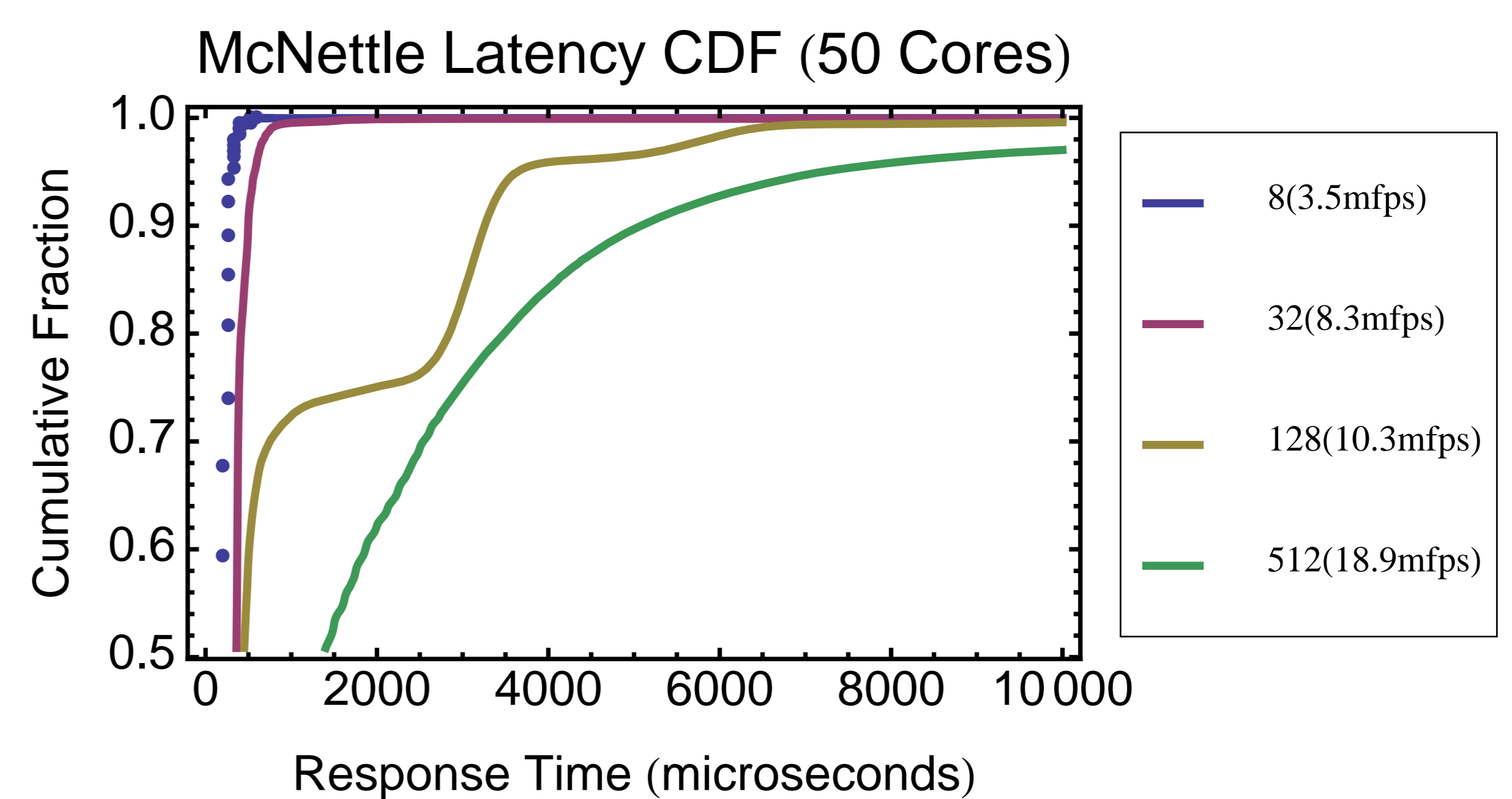
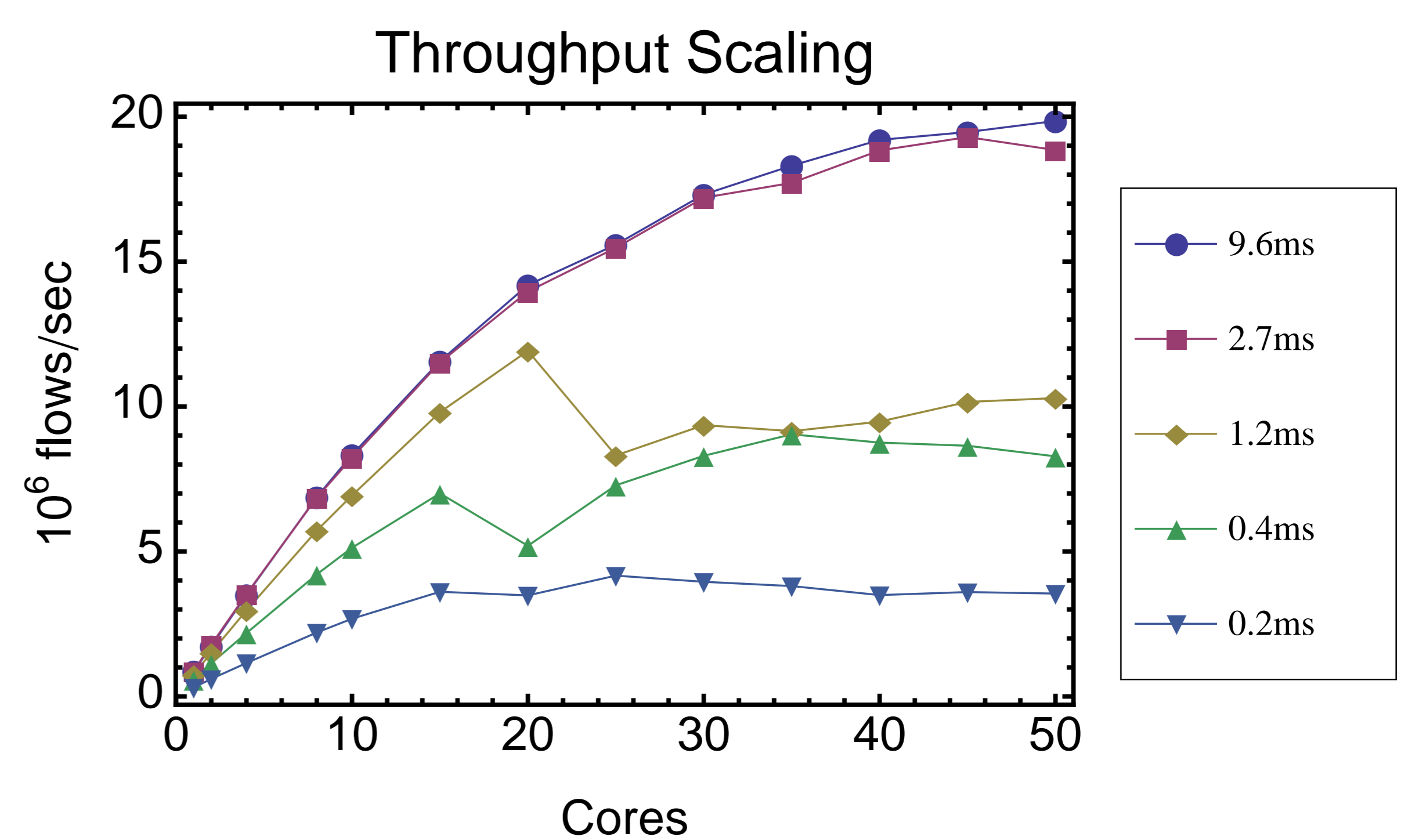
```
parMap (dijkstraRouting topology) switches
```

## EVALUATION SETUP



Generators simulate switches with attached hosts; simulates packet-in messages that would be sent if every host starts a TCP session to every other host.

## THROUGHPUT & LATENCY



## IMPLEMENTATION: KEY POINTS

- Utilize Glasgow Haskell Compiler's (GHC) multicore runtime system, which implements scalable, work-balancing thread scheduler.
- Use batching to minimize system call overhead.
- Custom memory management of critical message buffers to avoid garbage collector & storage allocator overhead.
- Reader-writer queues to efficiently support cross-switch messaging.
- Improve parallelization of GHC garbage collector to reduce pause times.
- Parallelize GHC's IO event loop to reduce response time.
- Separate cores for NIC interrupt handling and McNettle threads.