

Cutting the Cord: A Robust Wireless Facilities Network for Data Centers

Yibo Zhu, Xia Zhou[§], Zengbin Zhang, Lin Zhou,
Amin Vahdat[†], Ben Y. Zhao and Haitao Zheng

U.C. Santa Barbara, [§]Dartmouth College,

[†]U.C. San Diego & Google

yibo@cs.ucsb.edu

Data Center Networks (DCN)

- DCN: key infrastructures for mobile and big data applications

Google™

Walmart 

 JPMorgan

- Large and **dynamic** → management complexity
 - Highly dynamic data traffic
 - Shared by changing customers
 - Frequent failure, maintenance and upgrades

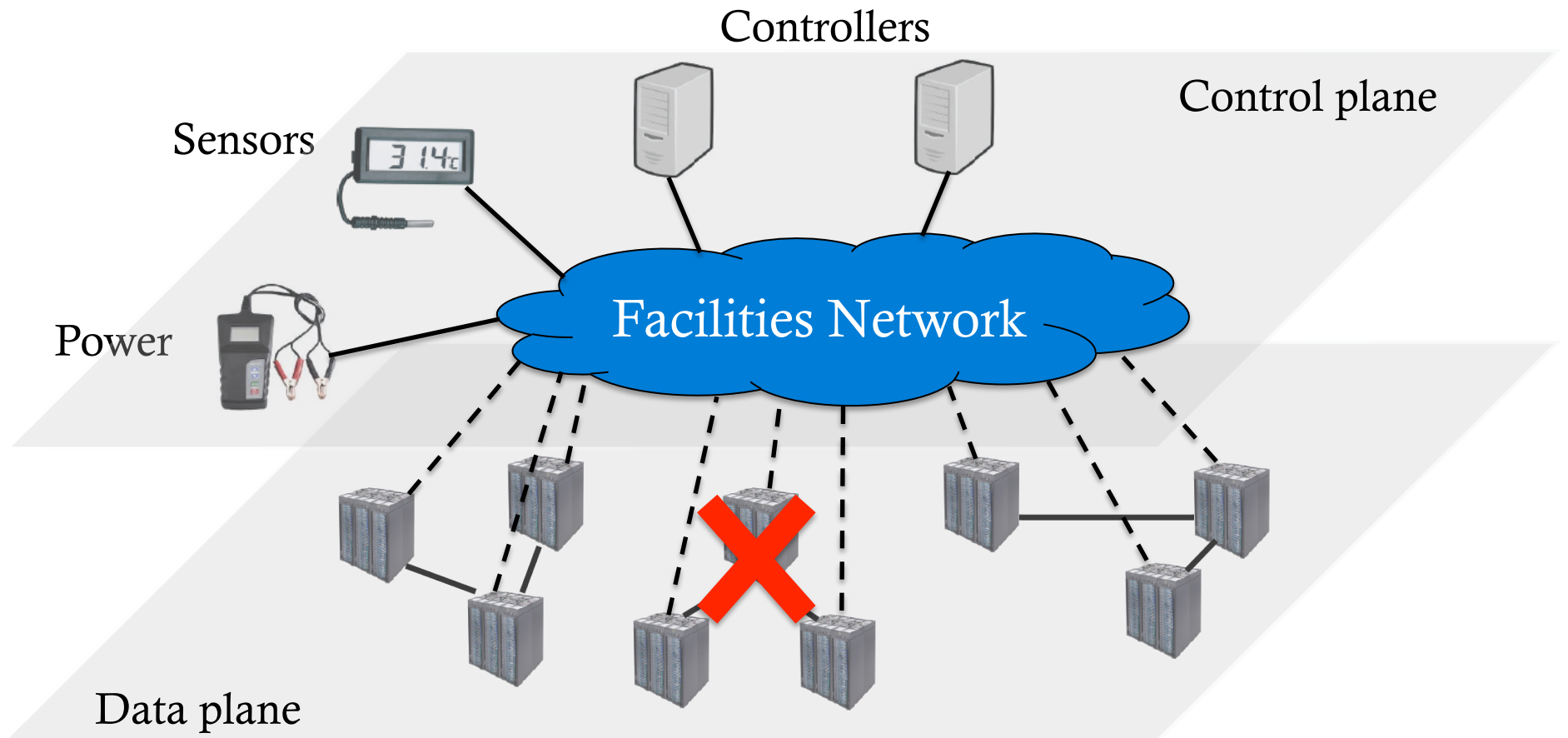
Beyond Data Plane

- Various control messages
 - Flow scheduling
 - Monitoring environment & power
 - Virtual machine imaging and configuration
 - Failure recovery
 - Bootstrap upgraded servers
- Must deliver timely and reliable
 - Not interfered by congested data traffic
 - Even when data plane not working



Upgrade ~100 servers
per day on average

A Facilities Network



Proposed DCN architecture

Requirements of Facilities Network

Performance

Low bandwidth

- 1Gbps enough?

Bounded delay¹

- One packet message <10ms
- 1MB Large message <500ms

¹Devoflow, SIGCOMM'11

Fault isolation

No fate-sharing

- Ideally physically separated

Robustness

Always connected

- Even when large portions down



Must remain working even racks taken off

Option: Wired Facilities Network

- Connect all devices using cables
 - In-band: share w/ data plane
 - Out-of-band



- Advantage: large capacity

- Challenges



- Out-of-band: **high cost, wiring**

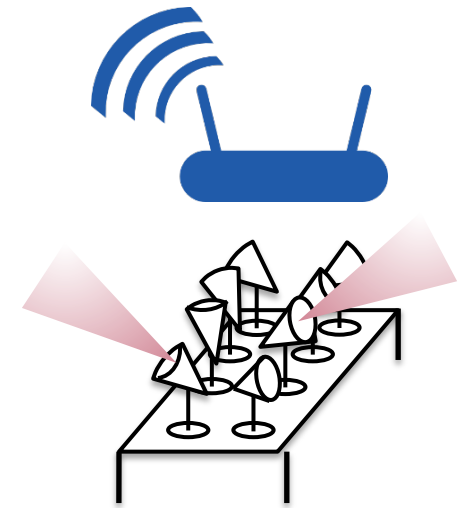


- Poor fault isolation/robustness
 - **Zero fault isolation for in-band**
 - Even out-of-band interrupted by cable tray maintenance



Option: Wireless Facilities Network

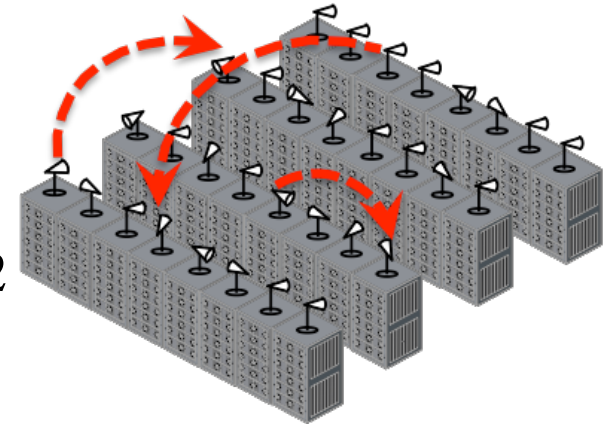
- Place radios on top of racks
 - WiFi (1.3Gbps), 60GHz (6.76Gbps)
 - Enough bandwidth
- Advantages
 - Cost: low (no additional switches/cables)
 - Fault isolation: **physically isolated from data plane**
 - Robustness: automatically reform links
- Challenges
 - Delay from wireless interference
 - Link coordination



Choice of Wireless Technology



60GHz 3D
Beamforming,
SIGCOMM'12



Widely available

Well-understood

- Omni-directional
- Contend for channel

Large interference footprint

- Poor in dense DC
- Unpredictable delay

Recently available

Less-understood

- Highly directional
- Need coordination

Small interference footprint

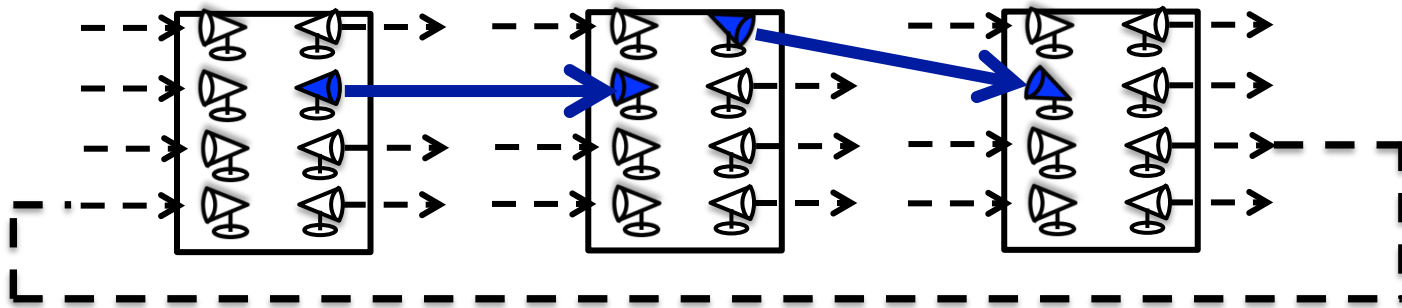
- Good for dense DC

Outline

- Motivation
- System design
 - Angora: a 60GHz facilities network
 - Wireless overlay design
 - Minimizing link interference
 - Fault recovery
- Evaluation
- Conclusion

Angora: a 60GHz Overlay

- Highly directional signal + limited radios per rack → **limited connections per rack**
- Antenna alignment → **extra delay** ☹️



- Angora: fixed topology overlay
 - Multi-hop → any-to-any connectivity
 - Fixed topology → **no antenna coordination** → no extra controllers, minimize delay

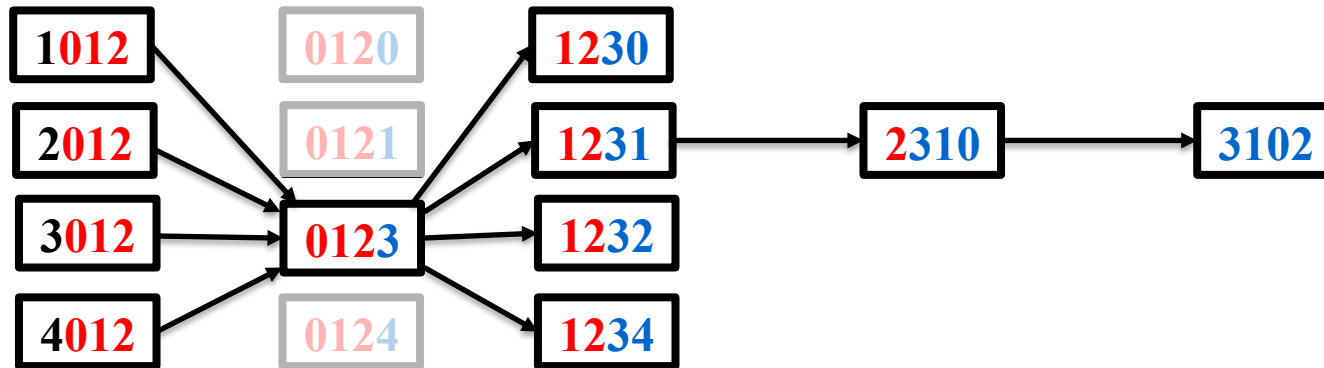
Structured Overlay Graph

- Key goal: minimize delay (hop count)
- The constraint: constant number of radios per rack \rightarrow constant degree graph
- We choose **Kautz** graph
 - Given degree and # of nodes \rightarrow smallest diameter
- Hop count: **Kautz** $<$ **Random**¹ \ll **Fat-tree**
 - Wired networks prefer Fat-tree due to low wiring complexity

¹*Jellyfish*, NSDI'12

Kautz Graph

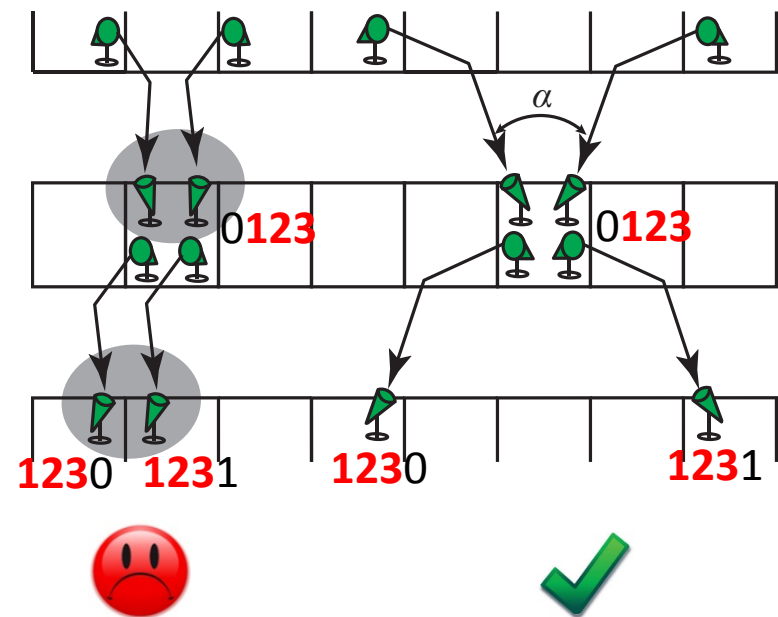
- Simple digit-shift routing



- Graph diameter = length of IDs
- Challenge: Kautz only supports specific node sizes
 - We designed an algorithm to handle arbitrary node size

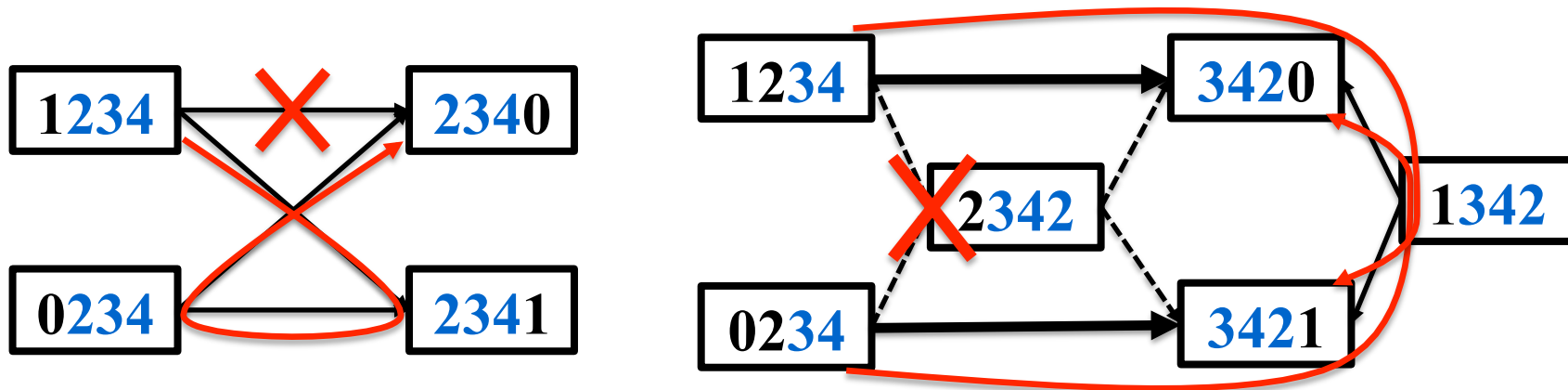
Node Naming and Interference

- Nodes naming affects interference
 - 60GHz interference: function of angular separation
- Goal: maximize angular separation between links
- Designed an optimal naming scheme
 - Achieved 14° angular separation in practice



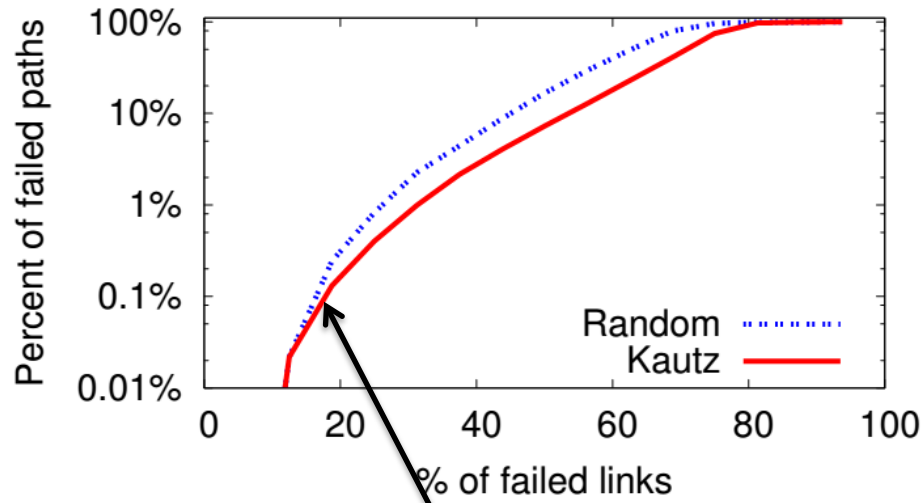
Failure Recovery Algorithms

- Link failure → remove a graph edge
 - May happen when radio fails, or signal blocked
 - Leverage Kautz structure to re-route the traffic
- Rack failure → remove a graph node
 - Similar deterministic algorithm



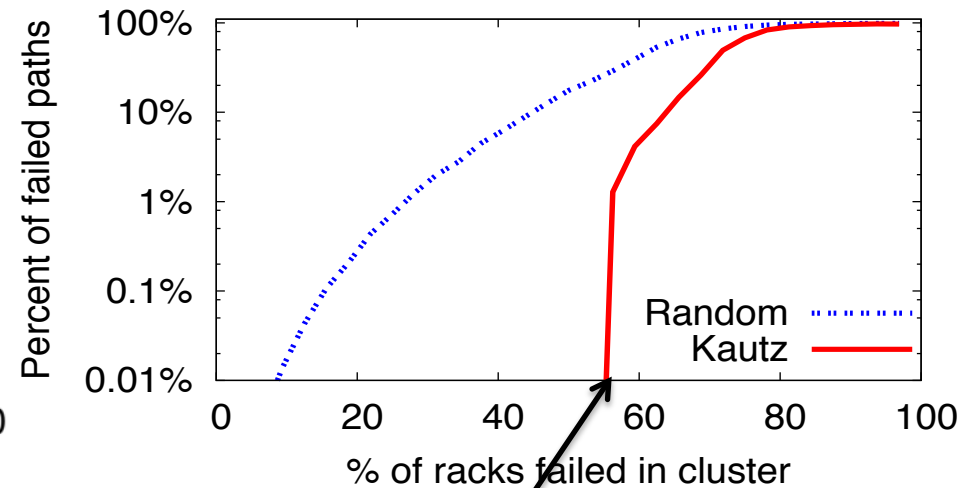
Failure Recovery Results

Random link failures



0.1% paths fail when
20% of links fail

Collocated rack failures



100% connectivity until
>50% collocated racks fail

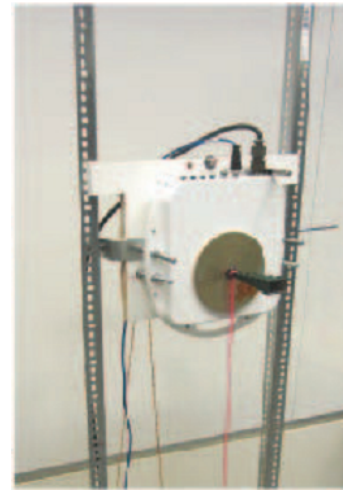
- Structural fault recovery → good robustness
- Deterministic algorithms → no extra coordinator

Outline

- Motivation
- System design
- Evaluation
 - Testbed
 - Simulation
- Conclusion

Testbed Validation

- Two testbeds
 - HXI: horn antennas
 - Wilocity: 2x8 arrays, affordable for multi-hop



HXI testbed
Horn antenna



Wilocity testbed
2x8 array

- Single link performance
 - Measured per-second TCP throughput over **one month**
 - Average **800(HXI)/900Mbps** (capped by 1Gbps NIC)
 - Standard variation **<1%** average throughput → as stable as a wired link

Testbed Validation (Multi-hop)

- Without interference

- Latency is small
- Latency increases with hops

Path Length	10KB Latency
2 hops	2.5ms
3 hops	3.1ms
4 hops	3.5ms

Multi-hop performance

- Path self-interference

- Kautz → at most 4 hops → at most 2 hop-pairs interfere
- Leverage channel allocation (3 channels in 60GHz)
- <1% paths have self-interference

- Cross-path interference mitigated by node naming

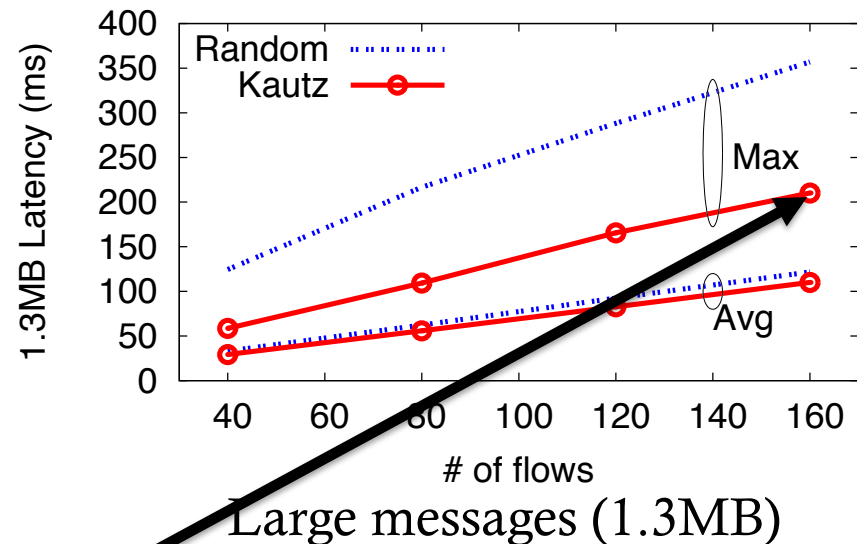
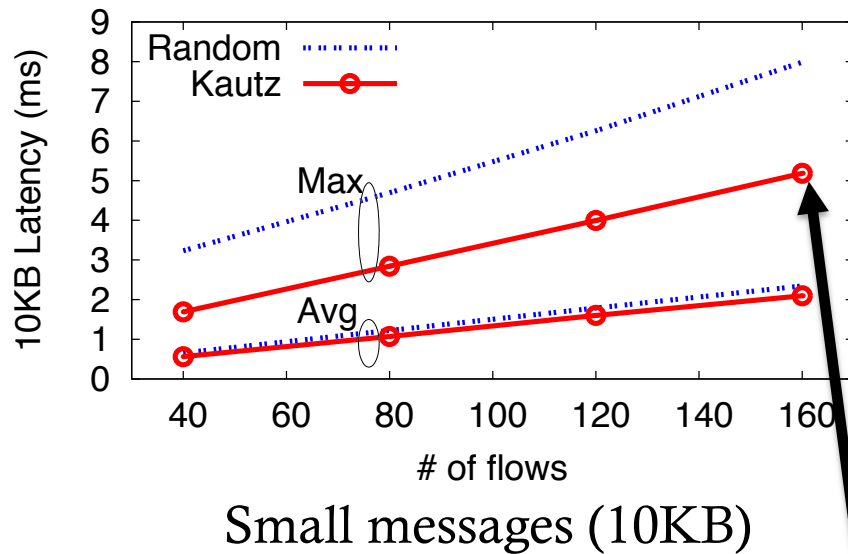
Multi-hop paths have low interference → small and predictable latency.

Large-scale Simulation

- We implement Angora in NS-3
 - Antenna: horns and arrays
 - 3D beamforming signal reflection
 - 802.11ad PHY/MAC
 - Kautz overlay routing
 - Medium size (320~480 racks) DCN layouts
- **Micro-benchmarks**: path hop count, concurrency, fault-tolerance
- **End-to-end performance**: single flow, Poisson flows, synchronized flows

End-to-end Performance

- Worst case: synchronized flows



- Tail delay satisfies facilities network requirements
- Structural (Kautz) \gg random at tails

Conclusion

- Motivation: build an orthogonal facilities network as a core tool for managing DCN.
- We propose Angora, a Kautz overlay built on 60GHz 3D beamforming links.
- Addressed challenges
 - 60GHz link coordination
 - By Kautz graph w/ arbitrary node size
 - Wireless interference
 - Fault tolerance