# Overlay Networks
## CS2510 Guest Lecture

Amy Babay

University of Pittsburgh
School of Computing and Information

# The Internet Revolution
## A Technical Perspective

A single, multi-purpose, IP-based network

- Each additional node increases its reach and usefulness (network effect)
- Each additional application domain increases its economic advantage
- Will therefore absorb/overtake most other networks
  - Already happened: mail to e-mail, fax to PDFs, phone to VoIP
  - Ongoing: TV, various control systems

# The Internet Revolution
# A Technical Perspective

A single, multi-purpose, IP-based network

- The art of design – end-to-end principle
  - Keep it simple in the middle …
    - Best-effort packet switching, routing (intranet, Internet)
  - … and smart at the edge
    - End-to-end reliability, naming

- Enabled dramatic scalability and adaptability
  - Survived for 5 decades and counting
  - Sustained at least 7 orders of magnitude growth

- Standardized and a lot rides on it
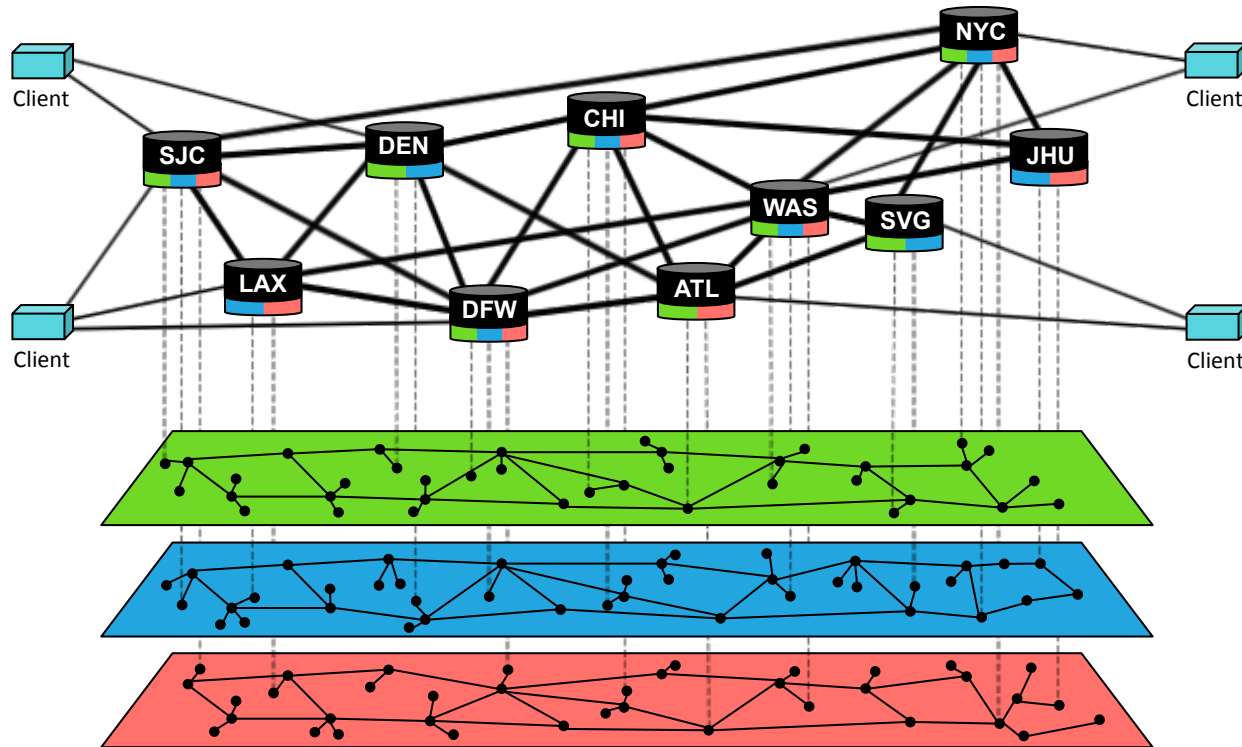  - The basic services are not likely to change

# A New Generation of Internet Applications

- Communication patterns
  - From point-to-point, to point-to-multipoint, to many-to-many
- High performance reliability
  - "Faster than real-time" file transfers
- Low latency interactivity
  - 100ms for VoIP
  - 80-100ms for interactive games
  - 65ms (one way) for remote robotic surgery, remote manipulation
- End-to-end dependability (availability, reliability)
  - From e-mail dependability – to phone service dependability – to remote surgery dependability – to power grid dependability
- System resiliency, security, and access control
  - From e-mail fault tolerance – to financial transaction security – to critical infrastructure (SCADA) intrusion tolerance

Overlay Networks: CS2510

# Addressing New Application Demands: Potential Approaches

- **Build specialized (non-IP) networks**
  - Was done decades before the Internet (e.g. TV Infrastructure)
  - Extremely expensive

- **Build private IP networks**
  - Avoids resource sharing issues, solves some of the scale issues
  - Expensive
  - Still limited by the basic end-to-end principle underlying the IP service

- **Build a better Internet**
  - Improvements and enhancements to IP (or TCP/IP stack)
  - "Clean slate design"
  - Long process of standardization and gradual adoption

- **Build overlay networks**

# Overlay Network Concept



Overlay Concept: use the Internet for underlying transport, but build *overlay networks* with software-based routers that run on top of the Internet to meet the needs of new applications

# The Structured Overlay Network Vision

- Key idea: put processing and context into the middle of the network, providing more flexibility and control
  - At overlay level
  - Underlying network maintains the end-to-end principle
- Three structured overlay network principles:
  1. Resilient network architecture
  2. Overlay node software architecture with global state and unlimited programmability
  3. Flow-based processing

"Structured Overlay Networks for a New Generation of Internet Services",
A. Babay, C. Danilov, J. Lane, M. Miskin-Amir, D. Obenshain, J. Schultz, J. Stanton, T. Tantillo, Y. Amir,
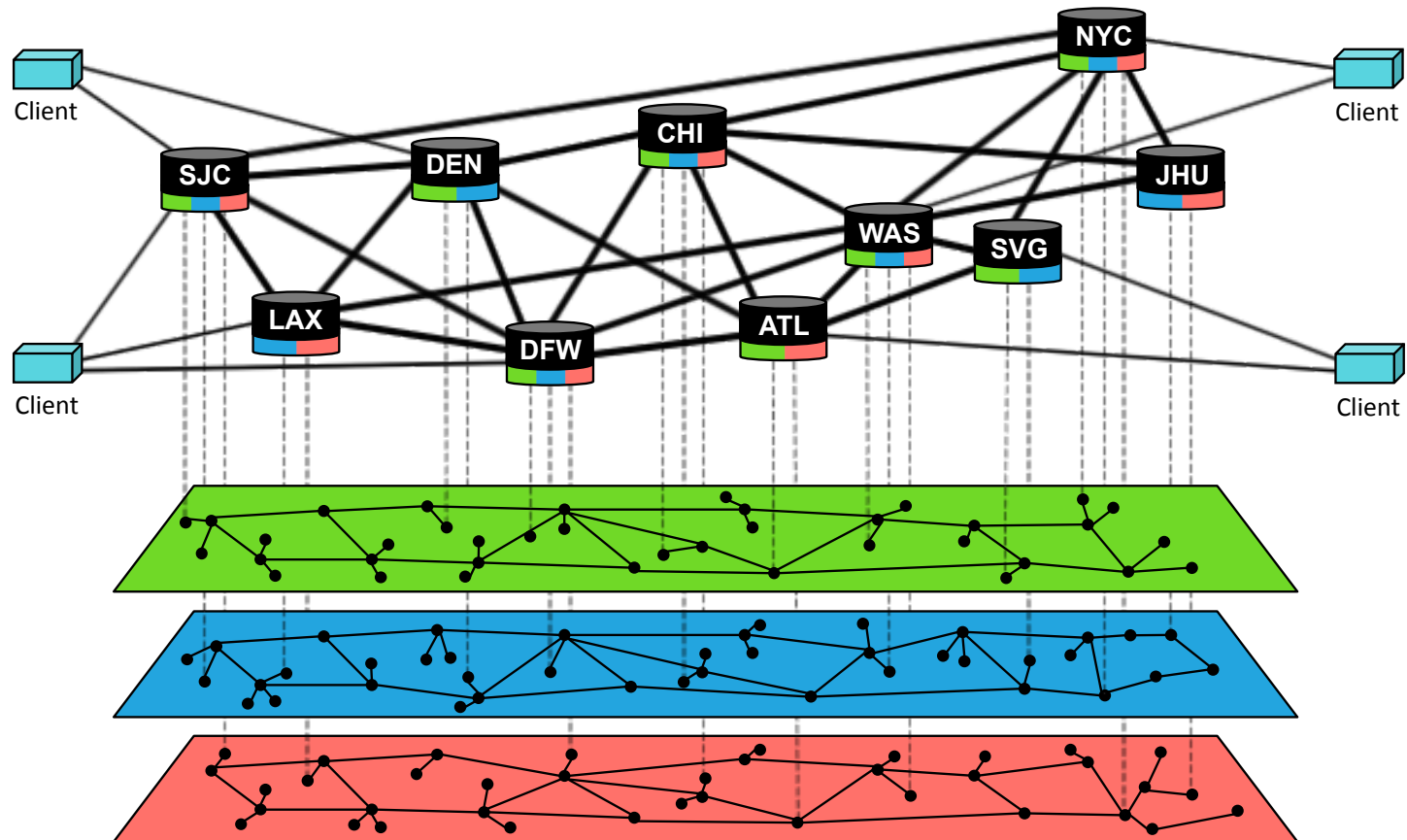*IEEE International Conference on Distributed Computing Systems (ICDCS)*, June 2017.

# Outline

- A New Generation of Internet Services
- The Structured Overlay Network Vision
  - Resilient network architecture
  - Overlay node software architecture with global state and unlimited programmability
  - Flow-based processing
- First Steps and Benefits
  - Responsive overlay routing with a resilient network architecture
  - Hop-by-hop reliability with flow-based processing and unlimited programmability
- The Quest for QoS
  - Almost-reliable real-time protocol for VoIP
  - Almost-reliable real-time protocol for Live TV
- Going even Faster
  - Remote manipulation, remote robotic surgery, collaborative virtual reality
  - Dissemination graphs with targeted redundancy
- Resilient Communication in a Hostile World
  - Intrusion-tolerant networking via structured overlays
  - Critical infrastructure applications
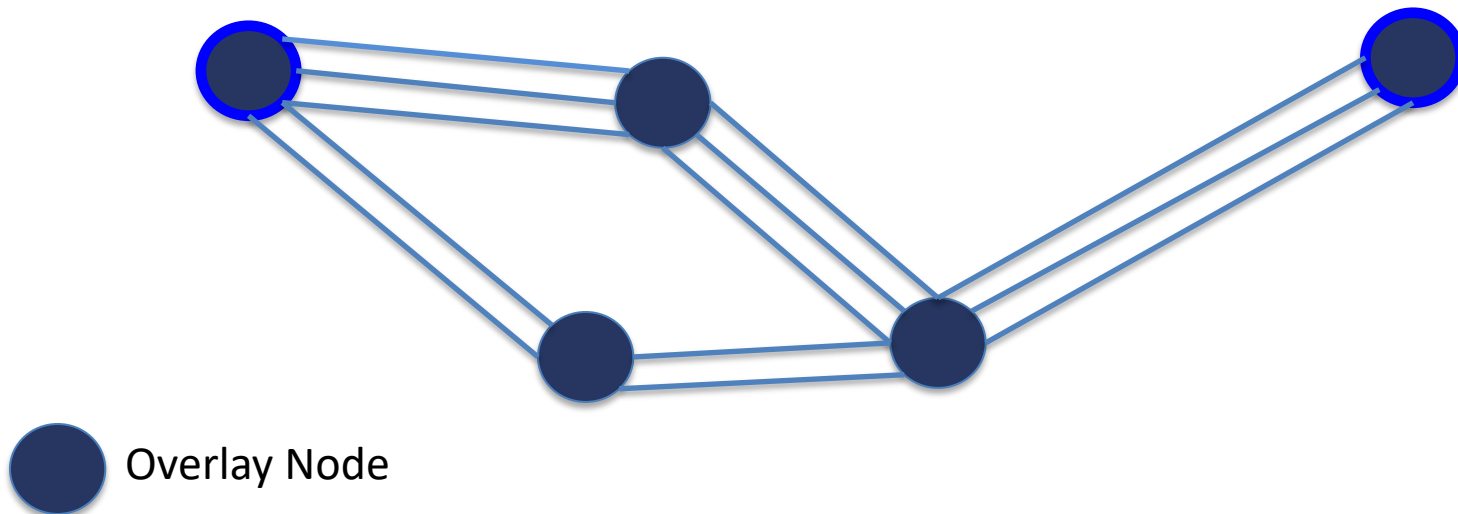- Future Directions

# Outline

- A New Generation of Internet Services
- **The Structured Overlay Network Vision**
  - **Resilient network architecture**
  - **Overlay node software architecture with global state and unlimited programmability**
  - **Flow-based processing**
- **First Steps and Benefits**
  - **Responsive overlay routing with a resilient network architecture**
  - **Hop-by-hop reliability with flow-based processing and unlimited programmability**
- The Quest for QoS
  - Almost-reliable real-time protocol for VoIP
  - Almost-reliable real-time protocol for Live TV
- Going even Faster
  - Remote manipulation, remote robotic surgery, collaborative virtual reality
  - Dissemination graphs with targeted redundancy
- Resilient Communication in a Hostile World
  - Intrusion-tolerant networking via structured overlays
  - Critical infrastructure applications
- Future Directions

# Resilient Network Architecture



U.S. portion of a resilient structured overlay network with overlay nodes located in strategic datacenters
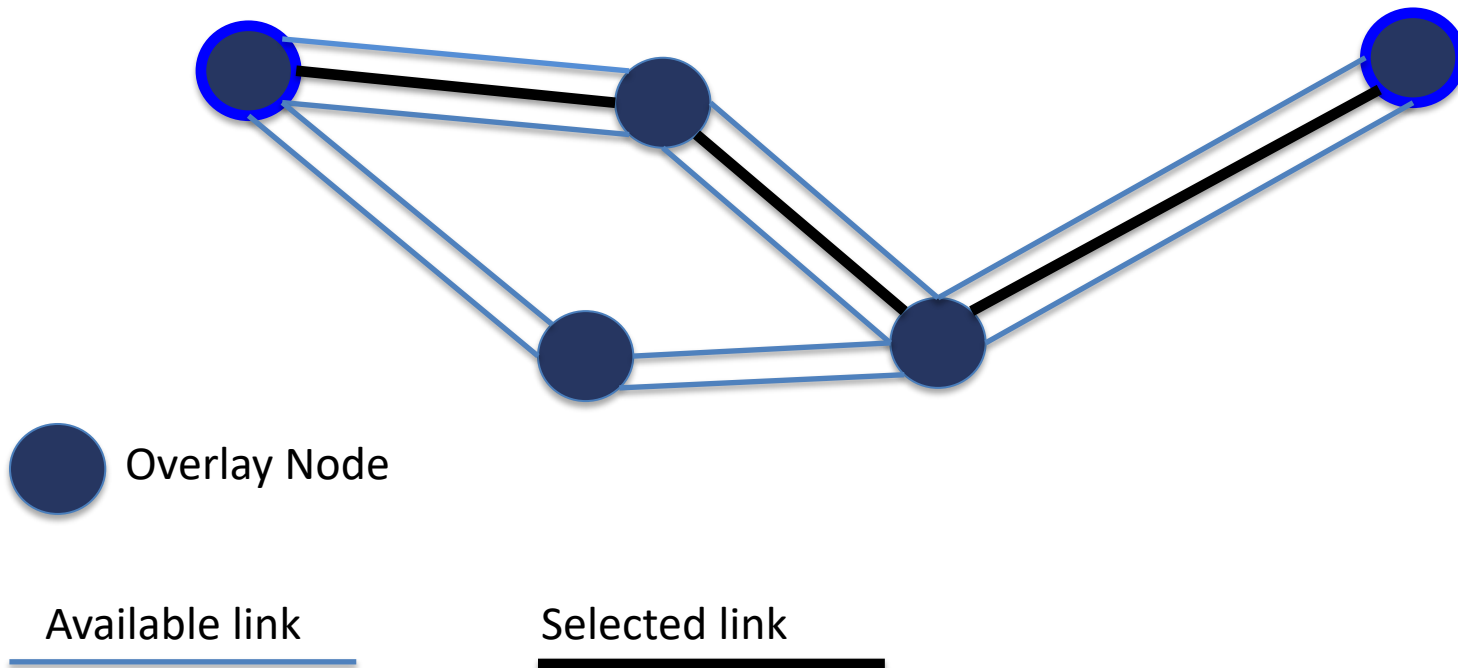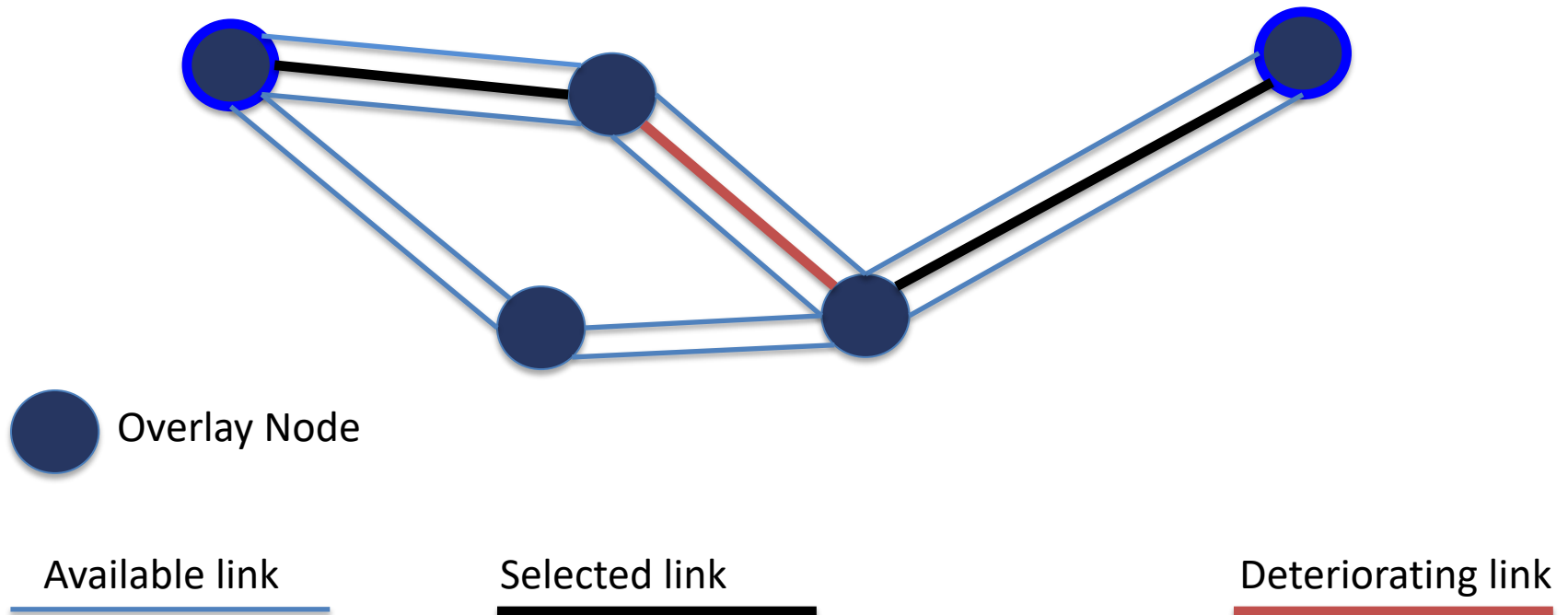
# Responsive Overlay Routing with a Resilient Network Architecture

- Utilizes multiple Tier 1 IP backbones
- Optimized overlay paths determine selected links
- Automatically and instantaneously switch to a better path
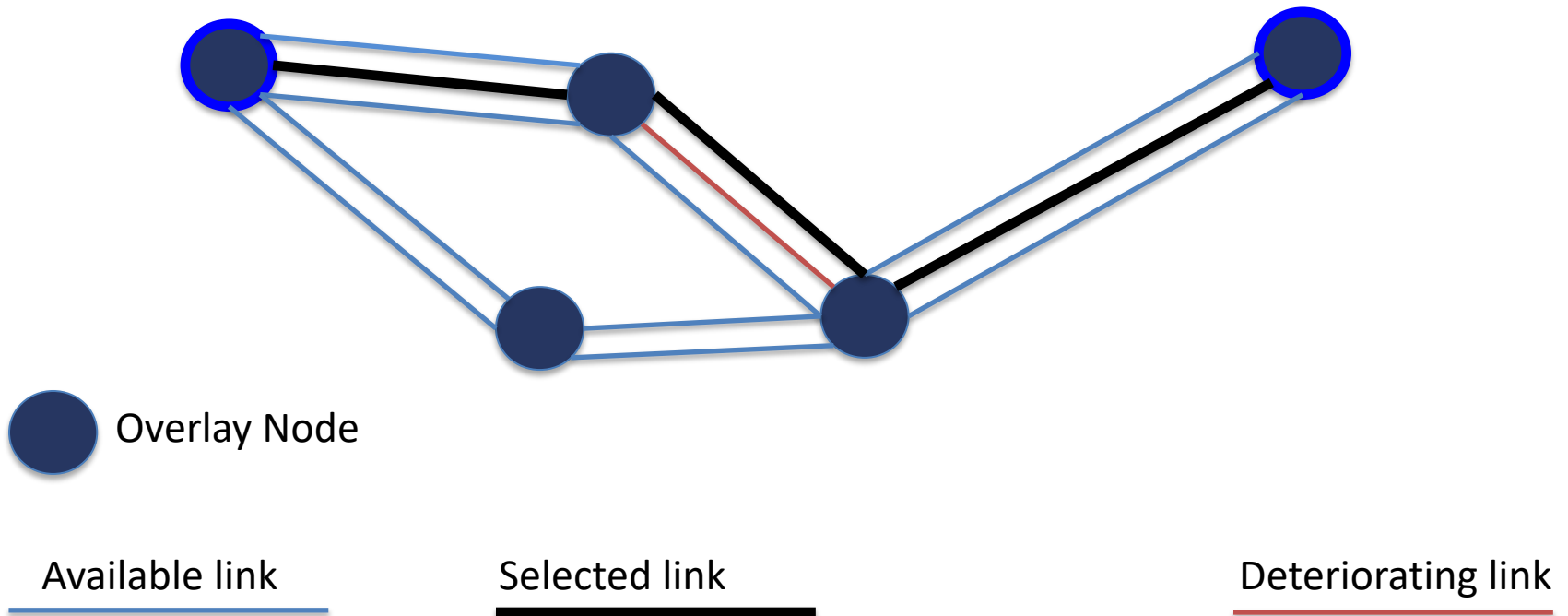
Overlay Node

Available link

# Responsive Overlay Routing with a Resilient Network Architecture

- Utilizes multiple Tier 1 IP backbones
- Optimized overlay paths determine selected links
- Automatically and instantaneously switch to a better path

Overlay Node
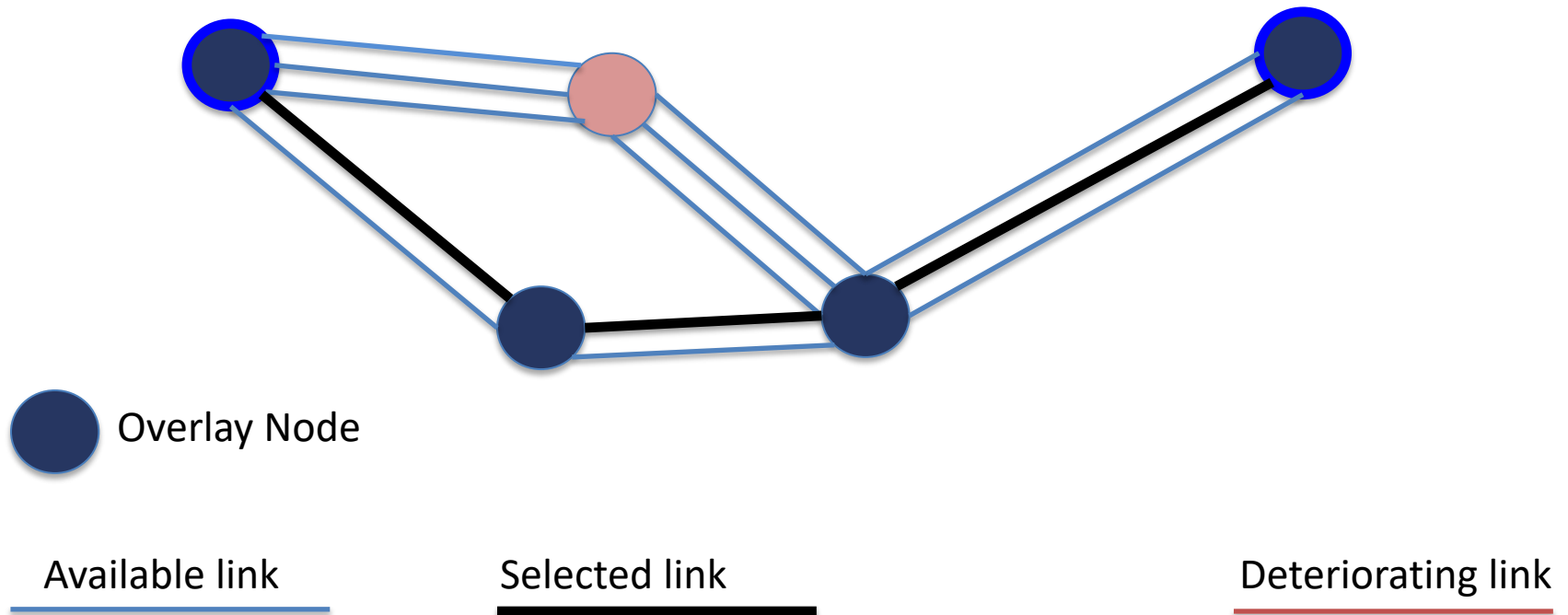
Available link

Selected link

# Responsive Overlay Routing with a Resilient Network Architecture

- Utilizes multiple Tier 1 IP backbones
- Optimized overlay paths determine selected links
- Automatically and instantaneously switch to a better path

Overlay Node

Available link

Selected link

Deteriorating link

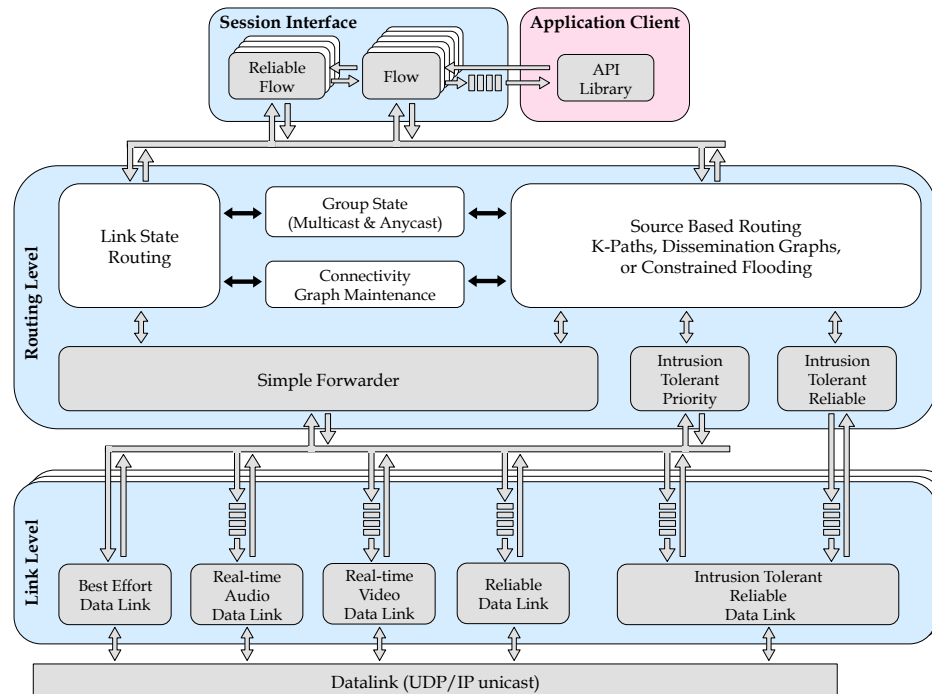# Responsive Overlay Routing with a Resilient Network Architecture

- Utilizes multiple Tier 1 IP backbones
- Optimized overlay paths determine selected links
- Automatically and instantaneously switch to a better path



Overlay Node

Available link          Selected link          Deteriorating link

# Responsive Overlay Routing with a Resilient Network Architecture

- Utilizes multiple Tier 1 IP backbones
- Optimized overlay paths determine selected links
- Automatically and instantaneously switch to a better path



Overlay Node

Available link          Selected link          Deteriorating link

# Overlay Node Software Architecture

- **Structured overlay messaging system**
  - Running overlay software routers on top of UDP as **user-level** internet applications
  - Using **commodity servers** in strategic datacenters
- **Easy-to-use programming platform**
  - API similar to the socket API
  - Additional, seamless API through packet interception
- **Deployable**
  - Vision partially realized by the Spines messaging system (www.spines.org) and its derivatives

# Overlay Node Software Architecture



- **Global State**
  - Possible due to the relatively small number of nodes (e.g. a few tens)
- **Unlimited programmability**
  - General purpose computers (or clusters) in datacenters
  - Flexible and extensible architecture

# Flow-based Processing

- Leverages flow-specific context
  - Flow: source + destination + application
- Enables services like:
  - Hop-by-hop recovery
  - De-duplication of retransmitted or redundantly transmitted packets in the middle of the network
  - Enhanced resiliency through flow-based fairness
- Allows different services to be selected for different application flows
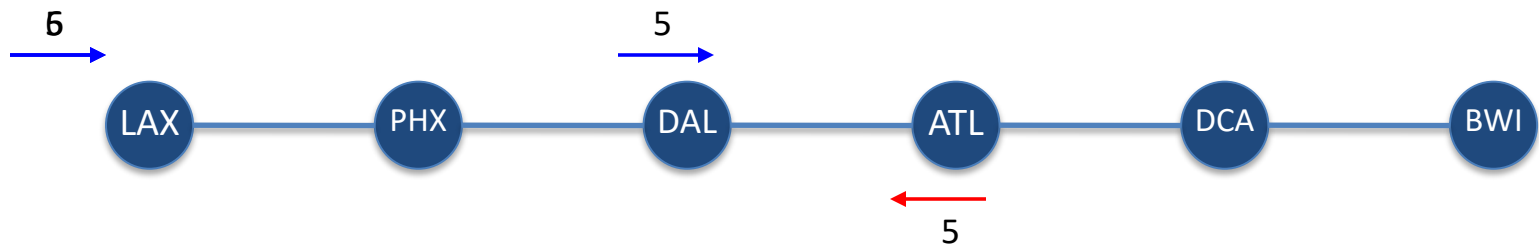
# Example: End-to-End Reliability

- ## 50 millisecond network
  - E.g.  Los Angeles to Baltimore
  - 50 milliseconds to tell the sender about the loss
  - 50 milliseconds to resend the packet
- ## At least 100 milliseconds to recover a lost packet

6

→

LAX ———————————————————————— BWI

←

5

# Example: End-to-End Reliability

- ## 50 millisecond network
  - E.g.  Los Angeles to Baltimore
  - 50 milliseconds to tell the sender about the loss
  - 50 milliseconds to resend the packet
- ## At least 100 milliseconds to recover a lost packet
  - Can we do better ?

LAX ———————————————————— BWI

# Hop-by-Hop Reliability with Flow-based Processing and Unlimited Programmability

- 50 millisecond network, five hops
  - 10 milliseconds to tell node DAL about the loss
  - 10 milliseconds to get the packet back from DAL
- Only 20 milliseconds to recover a lost packet
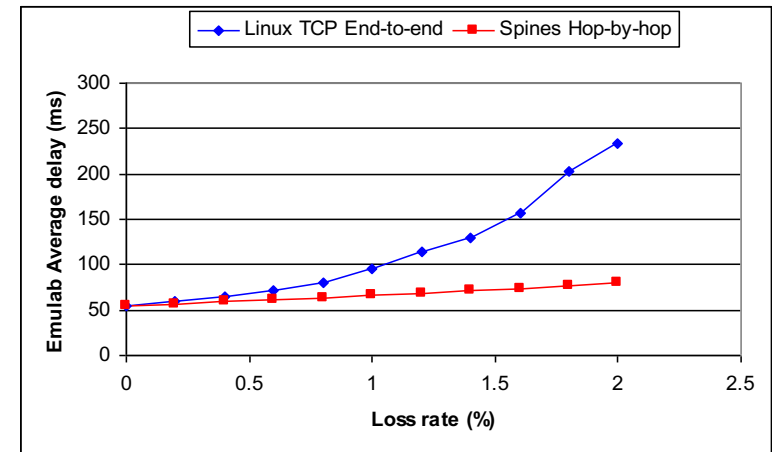  - Lost packet sent twice only on link DAL – ATL

# Average Latency



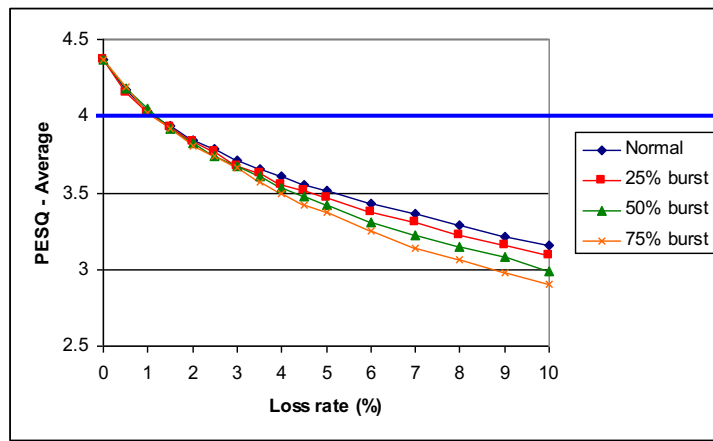Simulation · Spines on Emulab · Latency

"Reliable Communication in Overlay Networks", Y. Amir, C. Danilov,
*IEEE International Conference on Dependable Systems and Networks*, 2003.
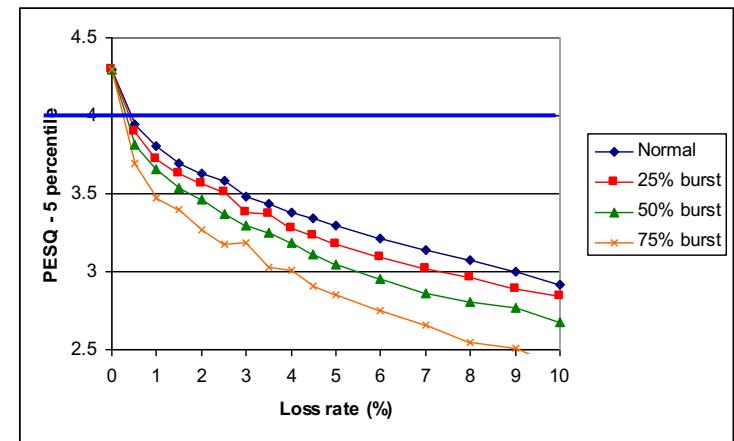
# Outline

- A New Generation of Internet Services
- The Structured Overlay Network Vision
  - Resilient network architecture
  - Overlay node software architecture with global state and unlimited programmability
  - Flow-based processing
- First Steps and Benefits
  - Responsive overlay routing with a resilient network architecture
  - Hop-by-hop reliability with flow-based processing and unlimited programmability
- **The Quest for QoS**
  - **Almost-reliable real-time protocol for VoIP**
  - **Almost-reliable real-time protocol for Live TV**
- Going even Faster
  - Remote manipulation, remote robotic surgery, collaborative virtual reality
  - Dissemination graphs with targeted redundancy
- Resilient Communication in a Hostile World
  - Intrusion-tolerant networking via structured overlays
  - Critical infrastructure applications
- Future Directions

# Siemens VoIP Challenge

- Can we maintain a "good enough" phone call quality over the Internet?
- High quality calls demand predictable performance
  - VoIP is interactive. Humans perceive delays at 100ms
  - The best-effort service offered by the Internet was not designed to offer any quality guarantees
  - Communication subject to dynamic loss, delay, jitter, path failures



50ms network delay

# A Structured Overlay Approach to VoIP

- Real-time "almost-reliable" hop-by-hop recovery protocol
  - Retransmission is attempted only once
  - Packets are only stored until delivery deadline (100ms) expires
- Responsive overlay routing with tailored routing metric
  - Cost metric based on measured latency and loss rate of the links
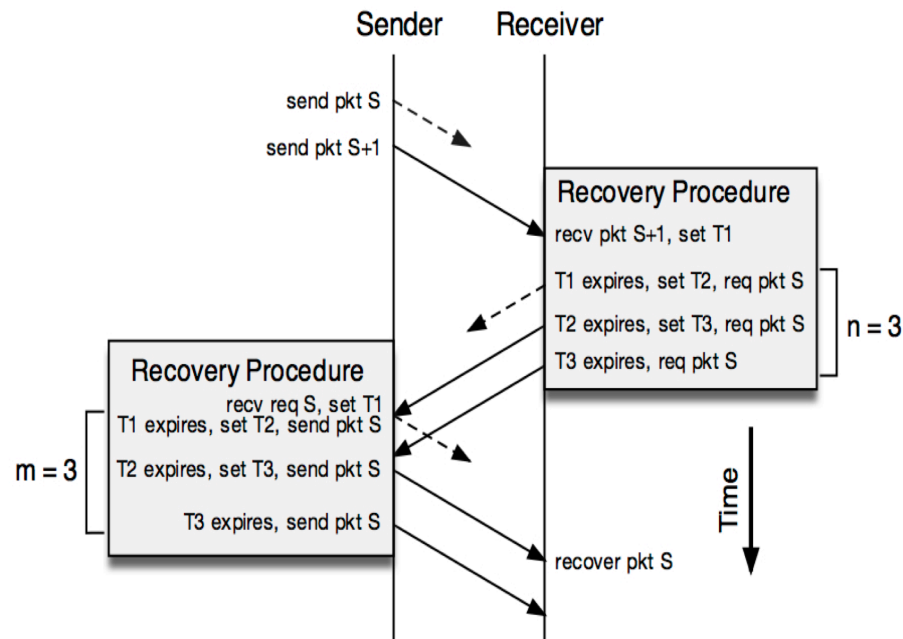  - Link cost equivalent to the expected packet latency when retransmissions are considered

"An Overlay Architecture for High Quality VoIP Streams", Y. Amir, C. Danilov, S. Goose, D. Hedqvist, A. Terzis, *IEEE Transactions on Multimedia*, 2006.
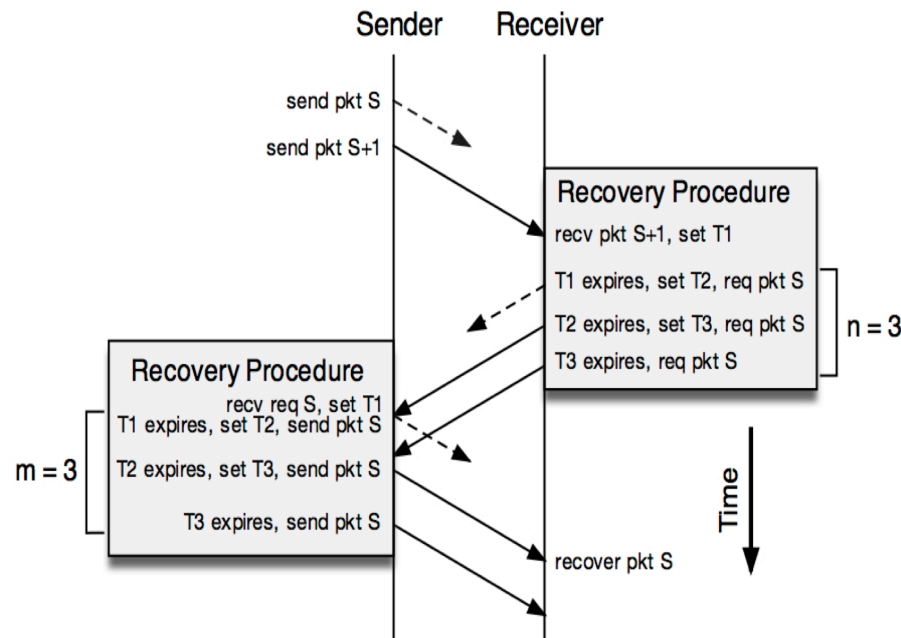
# The LTN TV Challenge



200ms one-way latency requirement, 99.999% reliability guarantee
40ms one-way propagation delay across North America

# Almost-Reliable Real-Time Protocol for Live TV



NM-strikes overlay link protocol: guaranteed timeliness, "almost reliable" delivery

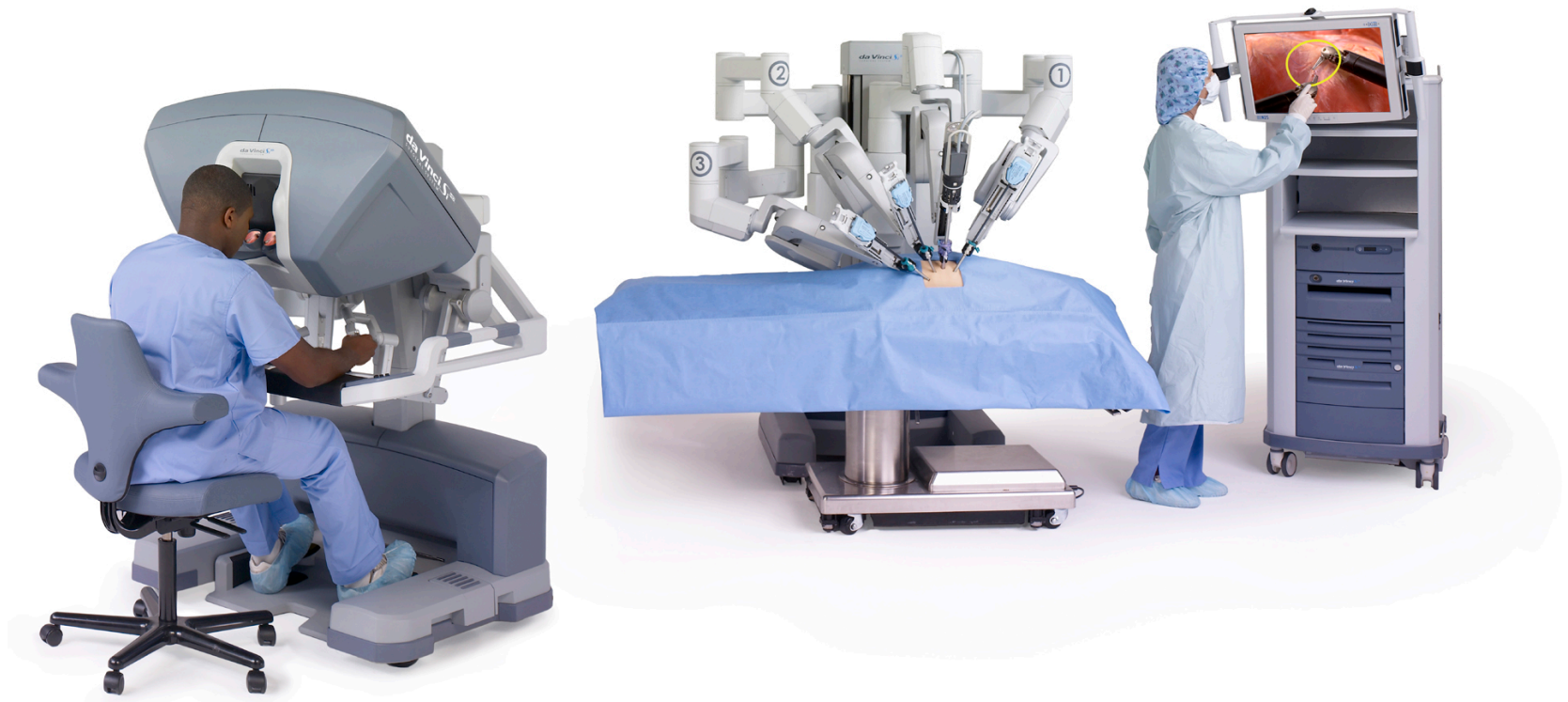# Almost-Reliable Real-Time Protocol for Live TV



| Network packet loss on one link (assuming 66% burstiness) | Loss experienced by flows on the LTN Network |
|---|---|
| 2% | < 0.0003% |
| 5% | < 0.003% |
| 10% | < 0.03% |

# Outline

- A New Generation of Internet Services
- The Structured Overlay Network Vision
  - Resilient network architecture
  - Overlay node software architecture with global state and unlimited programmability
  - Flow-based processing
- First Steps and Benefits
  - Responsive overlay routing with a resilient network architecture
  - Hop-by-hop reliability with flow-based processing and unlimited programmability
- The Quest for QoS
  - Almost-reliable real-time protocol for VoIP
  - Almost-reliable real-time protocol for Live TV
- **Going even Faster**
  - **Remote manipulation, remote robotic surgery, collaborative virtual reality**
  - **Dissemination graphs with targeted redundancy**
- Resilient Communication in a Hostile World
  - Intrusion-tolerant networking via structured overlays
  - Critical infrastructure applications
- Future Directions

# The Remote Surgery Challenge



130ms **round-trip** latency requirement

# The Remote Surgery Challenge



65ms **one-way** latency requirement

40ms one-way propagation delay across North America

# The Remote Surgery Challenge



65ms latency constraint – 40ms propagation delay
only 25ms available for recovery of lost packets

# Addressing the Challenge:
## Dissemination Graphs with Targeted Redundancy

- Stringent latency requirements give less flexibility for buffering and recovery

- Core idea: Send packets redundantly over a subgraph of the network (a dissemination graph) to maximize the probability that at least one copy arrives on time
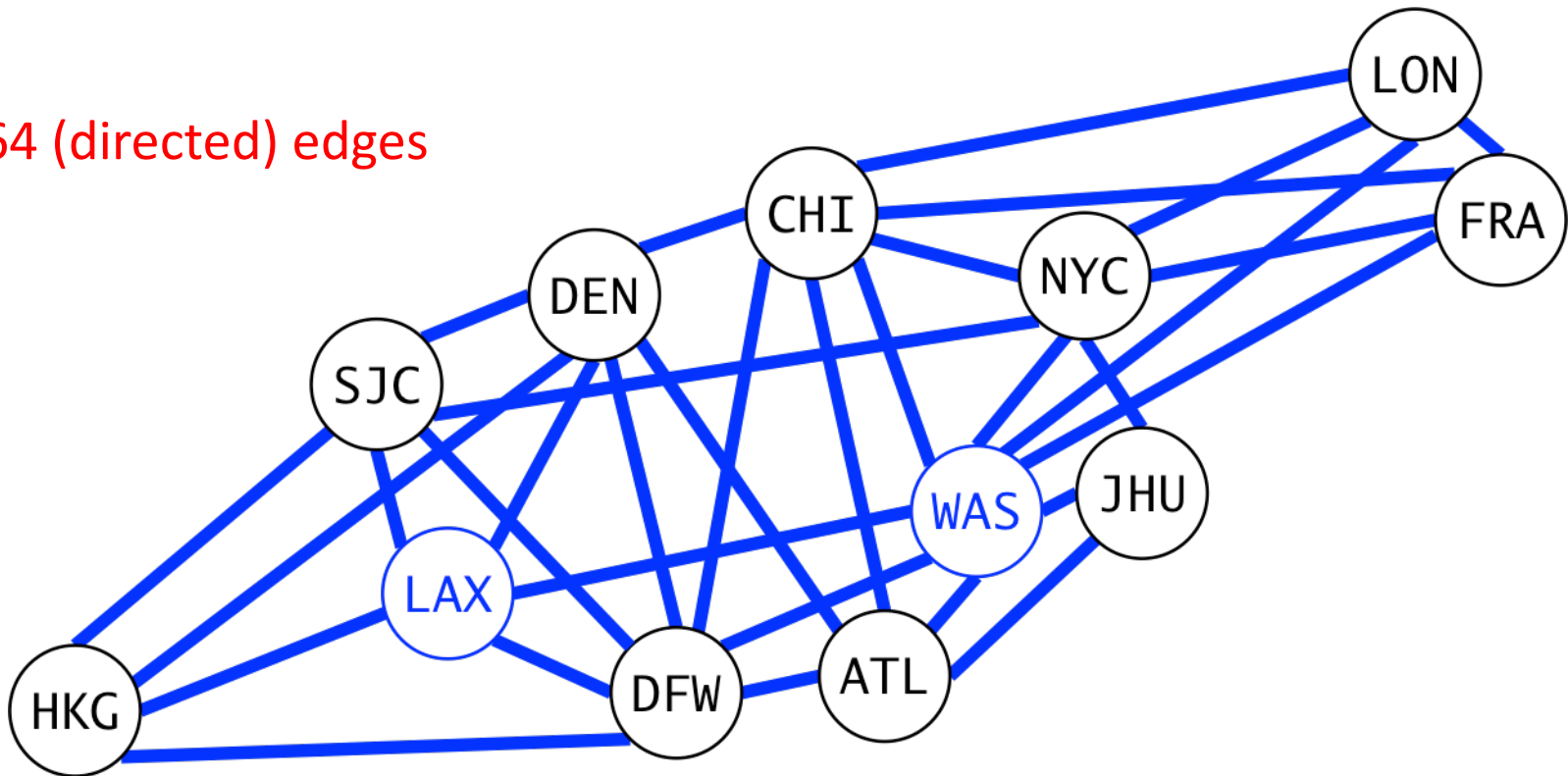
How do we select the subgraph (subset of overlay links) on which to send each packet?

"Timely, Reliable, and Cost-effective Internet Transport Service using Dissemination Graphs", Amy Babay, Emily Wagner, Michael Dinitz, and Yair Amir, *IEEE International Conference on Distributed Computing Systems (ICDCS),* 2017

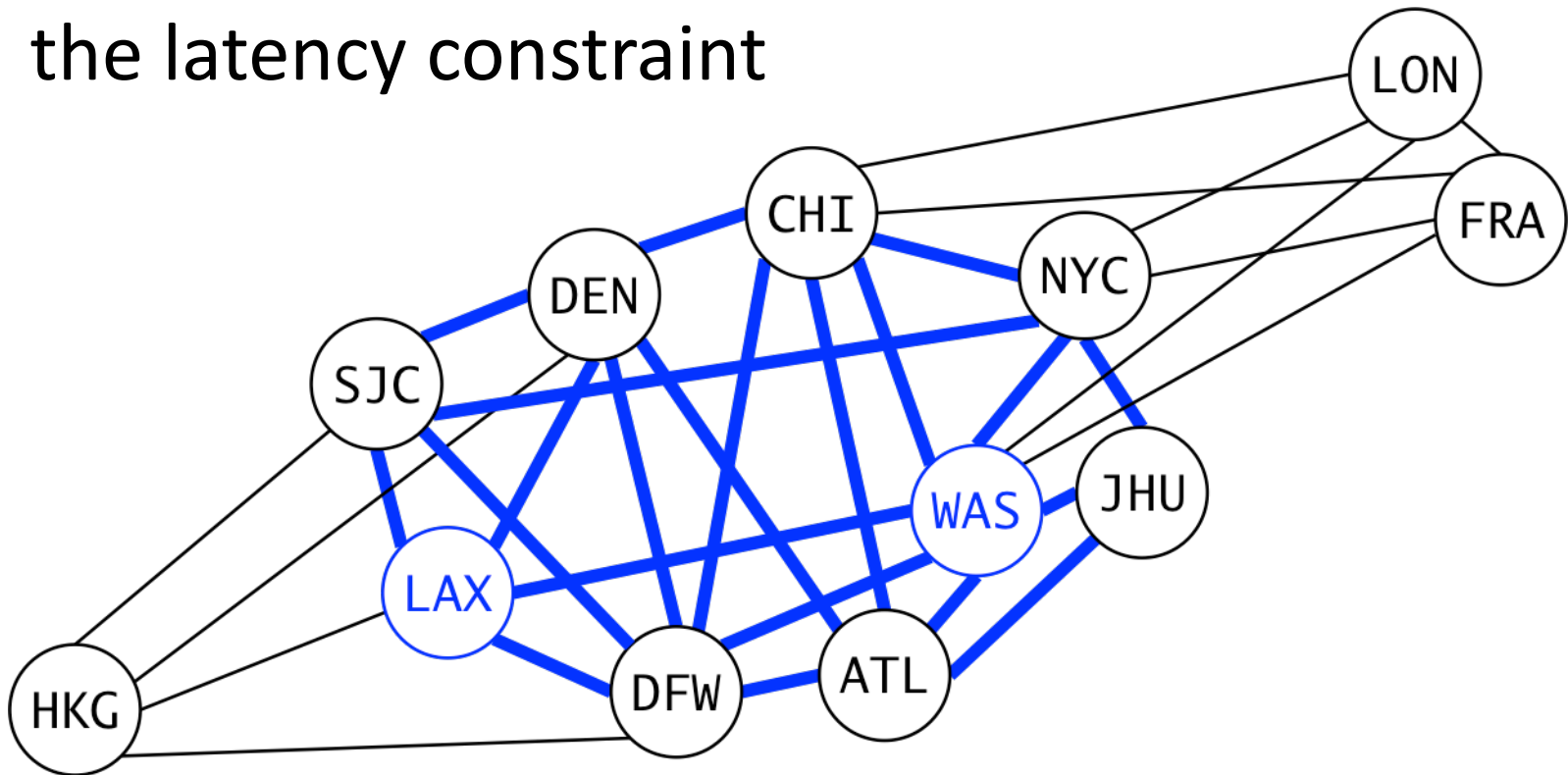# Initial Approaches to Selecting a Dissemination Graph

- Overlay Flooding: send on all overlay links
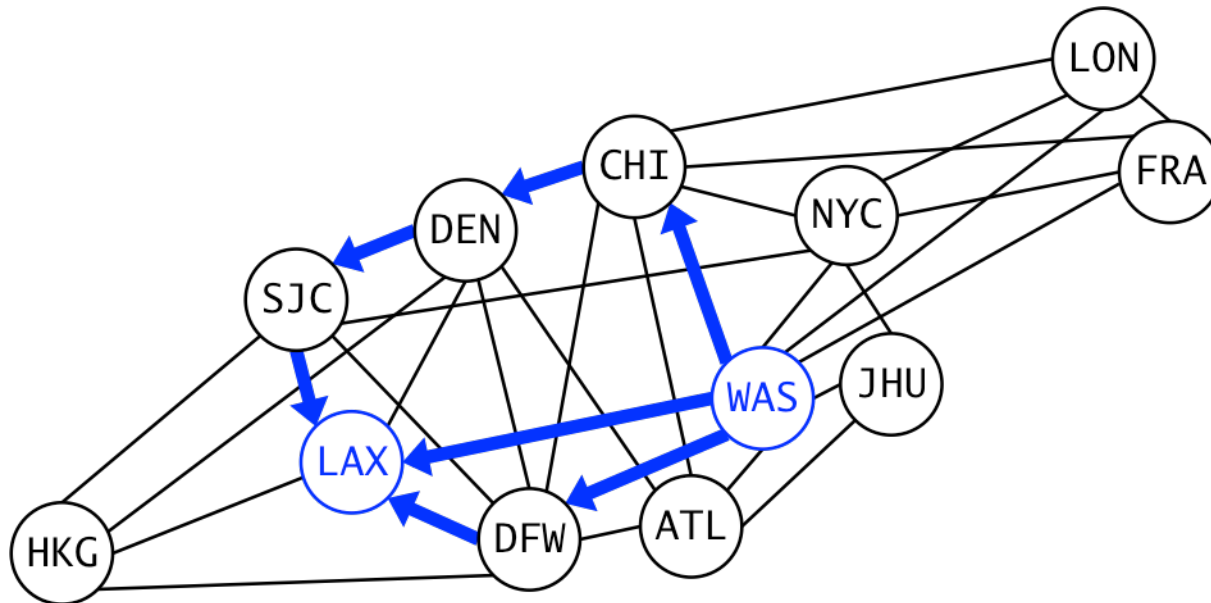  - Optimal in timeliness and reliability but expensive

64 (directed) edges

# Initial Approaches to Selecting a Dissemination Graph

- Time-Constrained Flooding: flood only on edges that can reach the destination within the latency constraint

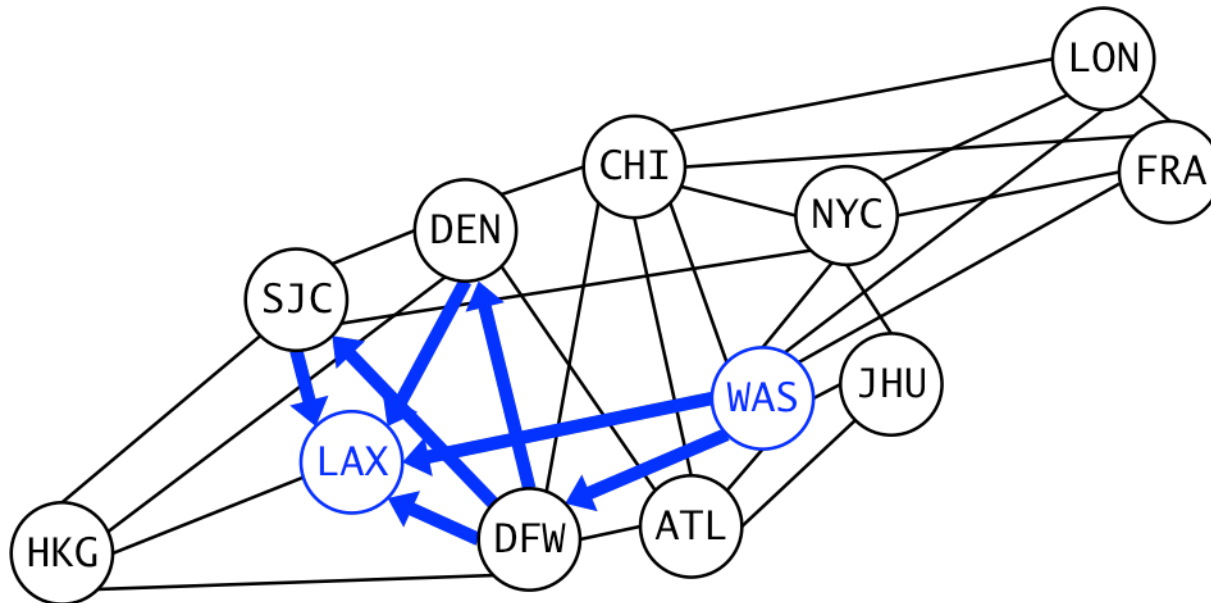# Initial Approaches to Selecting a Dissemination Graph

- Disjoint Paths: send on several paths that do not share any nodes (or edges)
  - Good trade-off between cost and timeliness/reliability
  - Uniformly invests resources across the network

# Selecting an Optimal Dissemination Graph

Can we use knowledge of the network characteristics to do better?

Invest more resources in more problematic regions:

# Problem Definition: Selecting an Optimal Dissemination Graph

- We want to find the best trade-off between cost and reliability (subject to timeliness)
  - Cost: # of times a packet is sent (= # of edges used)
  - Reliability: probability that a packet reaches its destination within its application-specific latency constraint (e.g. 65ms)
- **Service provider perspective**: minimize cost of providing an agreed upon level of reliability (SLA)

# Selecting an Optimal Dissemination Graph

- Solving the proposed problem is NP-hard
  - Without the latency constraint, computing reliability is the two-terminal reliability problem (which is #P-complete) [Val79]
  - Computing optimal dissemination graphs in terms of cost and reliability is also NP-hard
  - Exact calculations (via exhaustive search) can take on the order of tens of seconds for practical topologies – cannot support fast rerouting

# Data-Informed Dissemination Graphs

- Goal: Learn about the types of problems that occur in the field and tailor dissemination graphs to address common problem types

- Collected data on a commercial overlay topology (www.ltnglobal.com) over 4 months

- Analyzed how different dissemination-graph-based routing approaches (time-constrained flooding, single path, two disjoint paths) would perform (Playback Overlay Network Simulator)
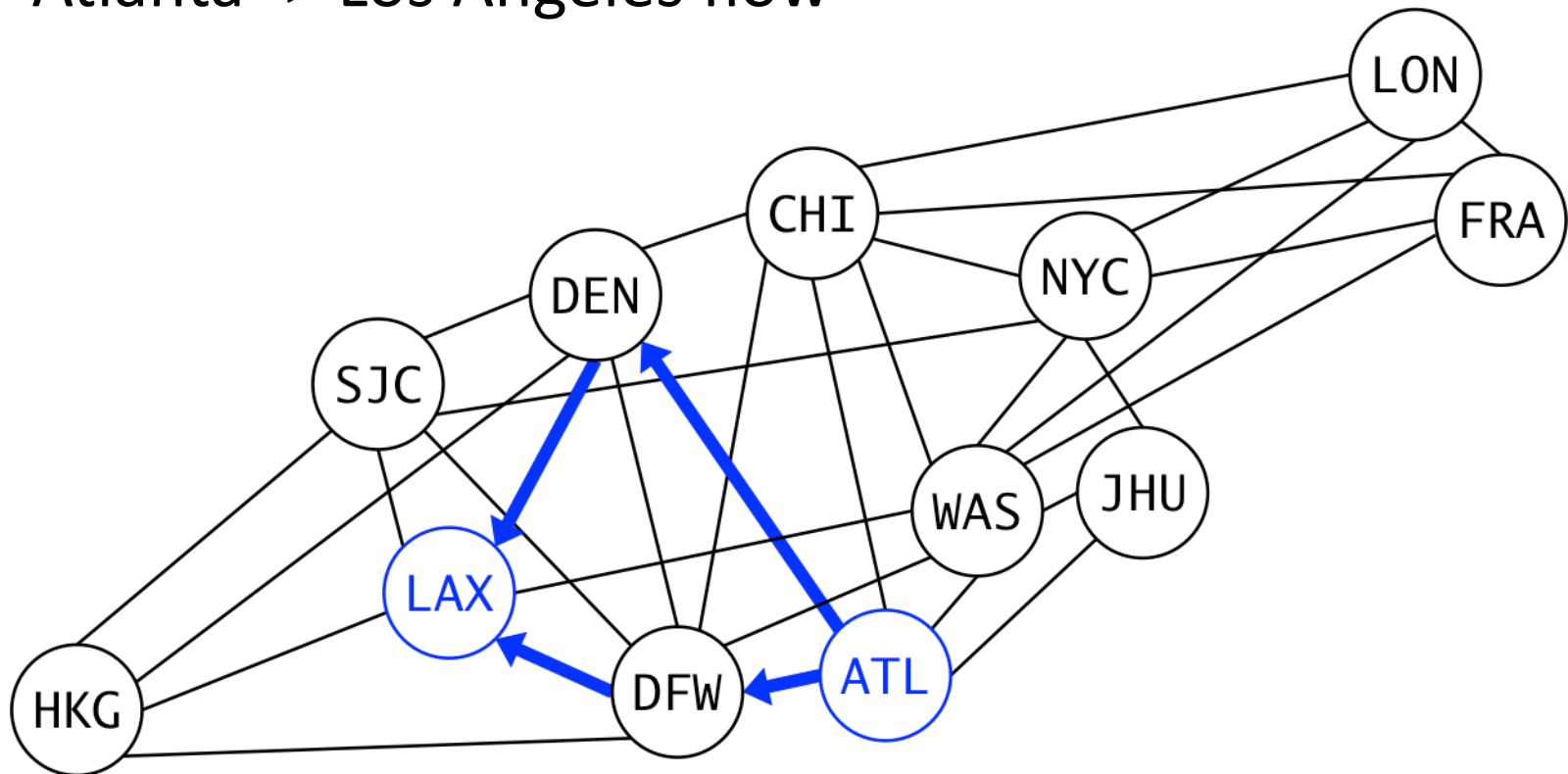
# Data-Informed Dissemination Graphs

- Key findings:

  - Two disjoint paths provide relatively high reliability overall
    - Good building block for most cases
  - Almost all problems not addressed by two disjoint paths involve either:
    - A problem at the source
    - A problem at the destination
    - Problems at both the source and the destination

# Dissemination Graphs with Targeted Redundancy

- Our approach:
  - Use two (dynamic) disjoint paths graph in the normal case
  - Pre-compute three additional graphs per flow:
    - Source-problem graph
    - Destination-problem graph
    - Robust source-destination problem graph (dynamically combined with two disjoint paths)
  - If a problem is detected at the source and/or destination of a flow, switch to the appropriate dissemination graph
  - Converts hard optimization problem into easy classification problem

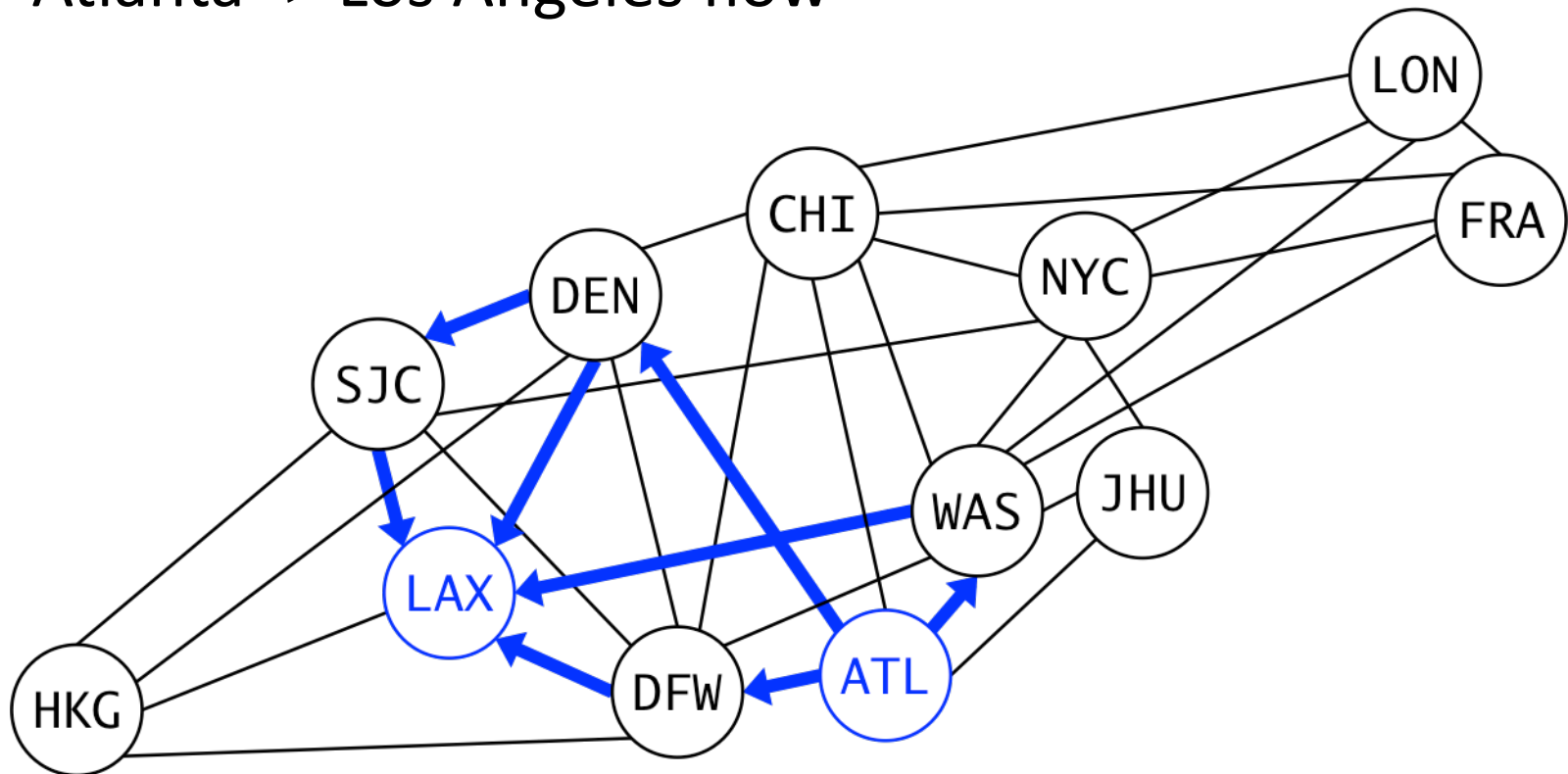# Dissemination Graphs with Targeted Redundancy: Example

- Atlanta -> Los Angeles flow



Two node-disjoint paths dissemination graph (4 edges)

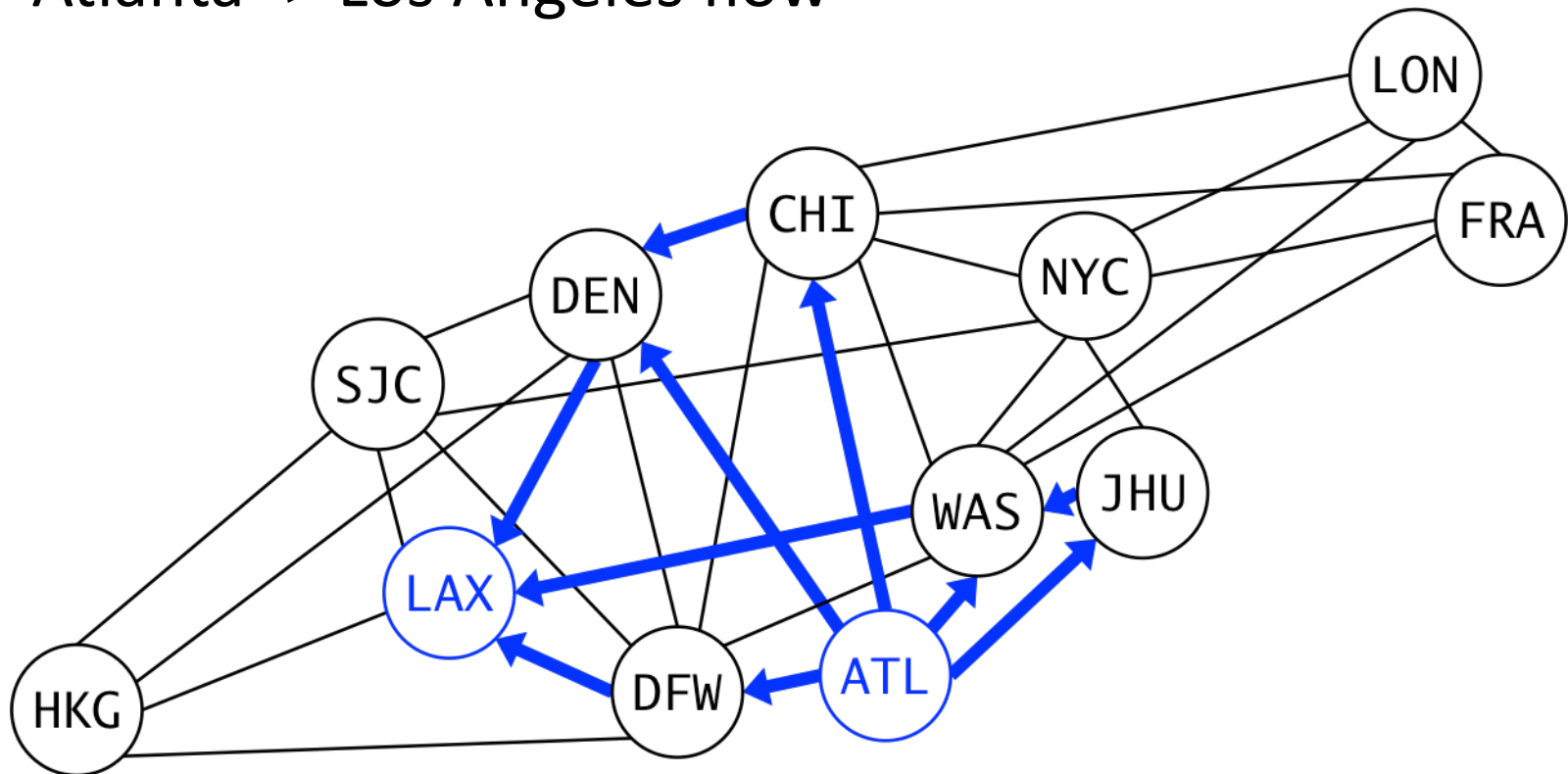# Dissemination Graphs with Targeted Redundancy: Example

- Atlanta -> Los Angeles flow



Destination-problem dissemination graph (8 edges)

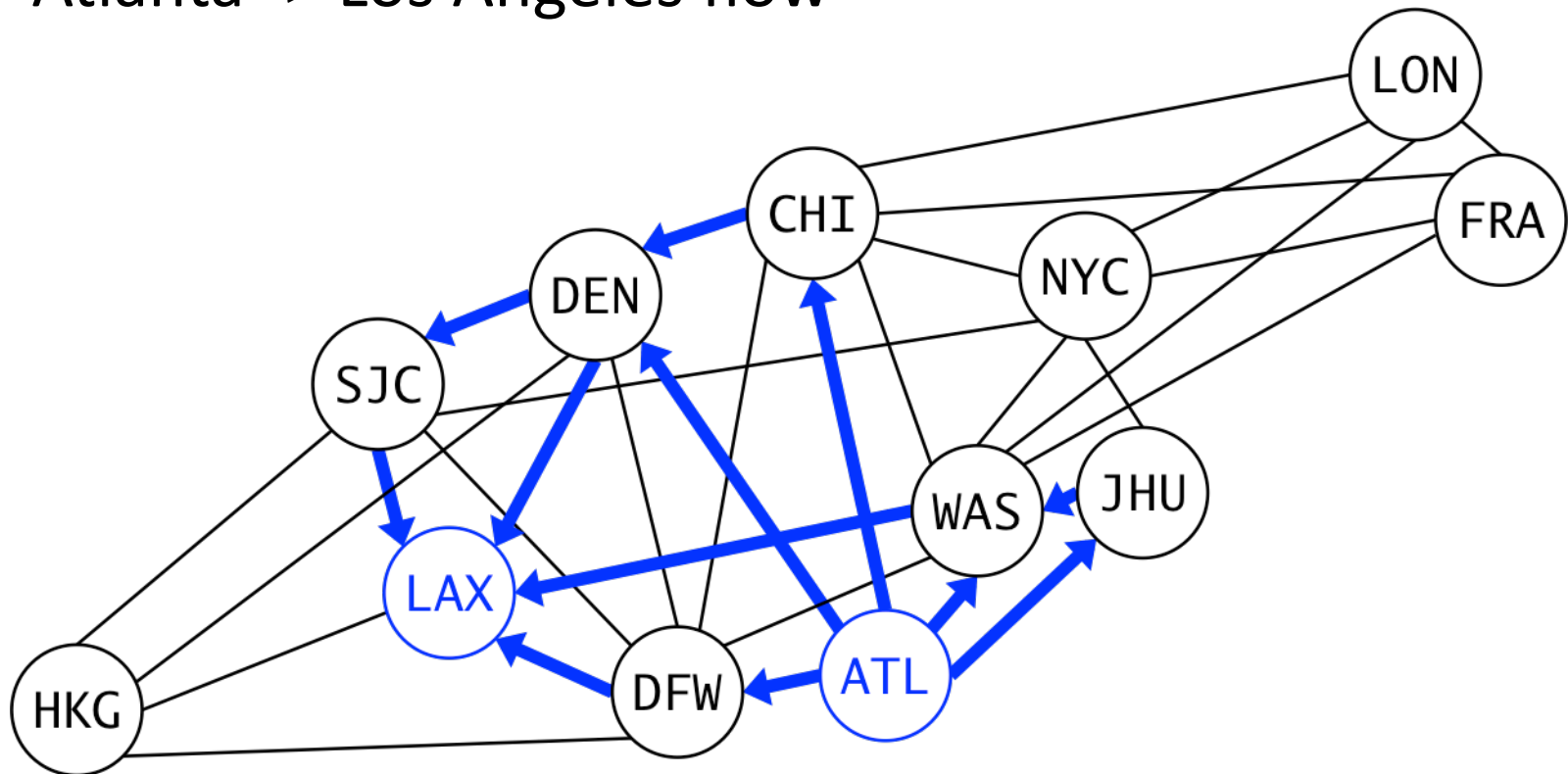# Dissemination Graphs with Targeted Redundancy: Example

- Atlanta -> Los Angeles flow



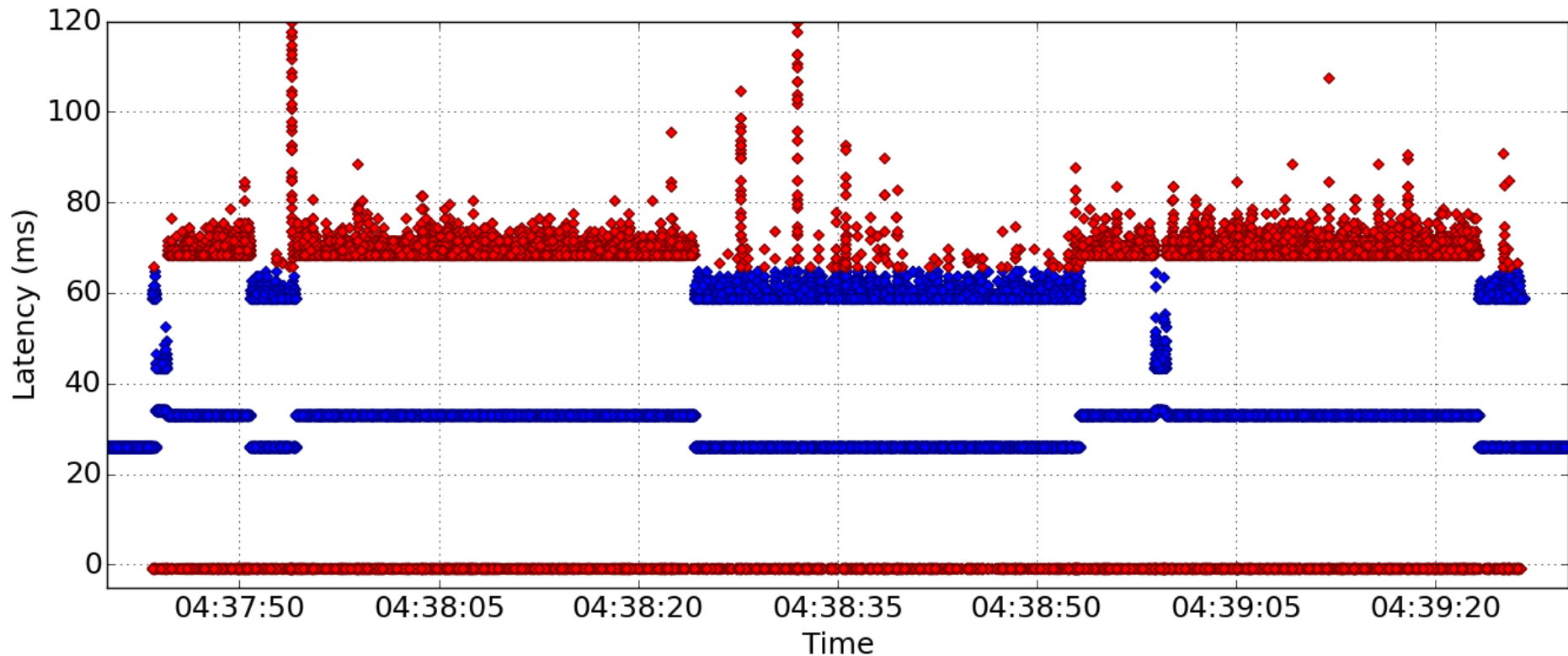Source-problem dissemination graph (10 edges)

- Atlanta -> Los Angeles flow



Robust source-destination-problem dissemination graph (12 edges)

# Dissemination Graphs Case Study: Single Path

- Case study: Atlanta -> Los Angeles; August 15, 2016



Packets received and dropped over a 110-second interval using dynamic single path (27,353 lost/late packets, 5 packets with latency over 120ms not shown)

# Dissemination Graphs Case Study: Single Path
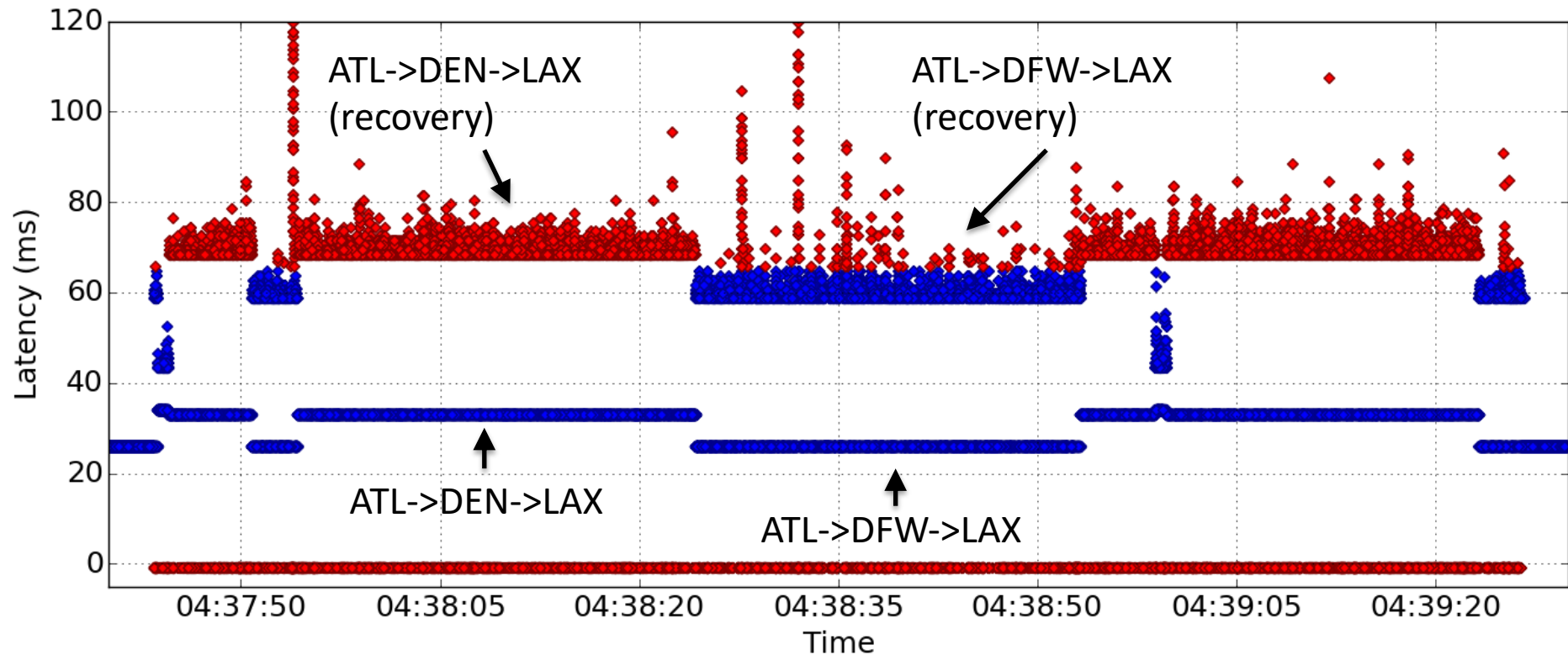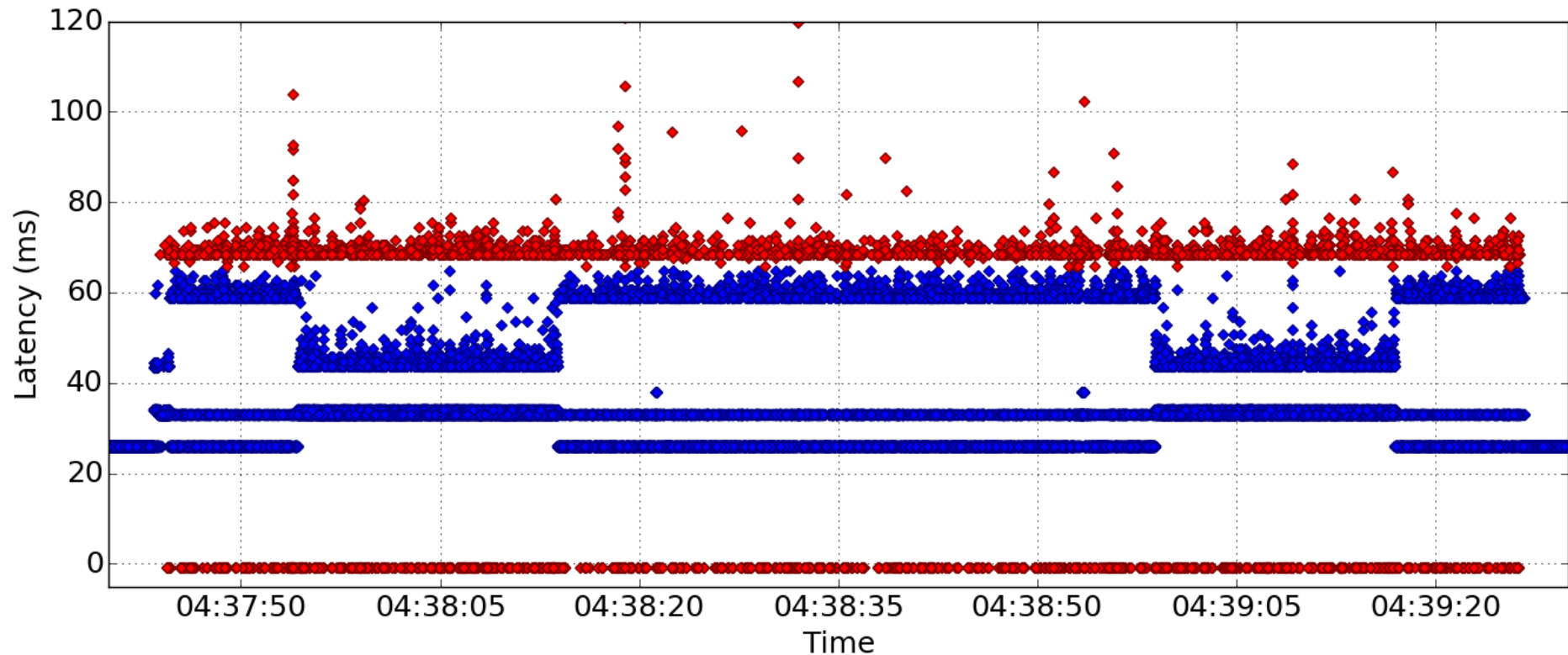
- Case study: Atlanta -> Los Angeles; August 15, 2016



Packets received and dropped over a 110-second interval using dynamic single path (27,353 lost/late packets, 5 packets with latency over 120ms not shown)

# Dissemination Graphs Case Study: Two Node-Disjoint Paths
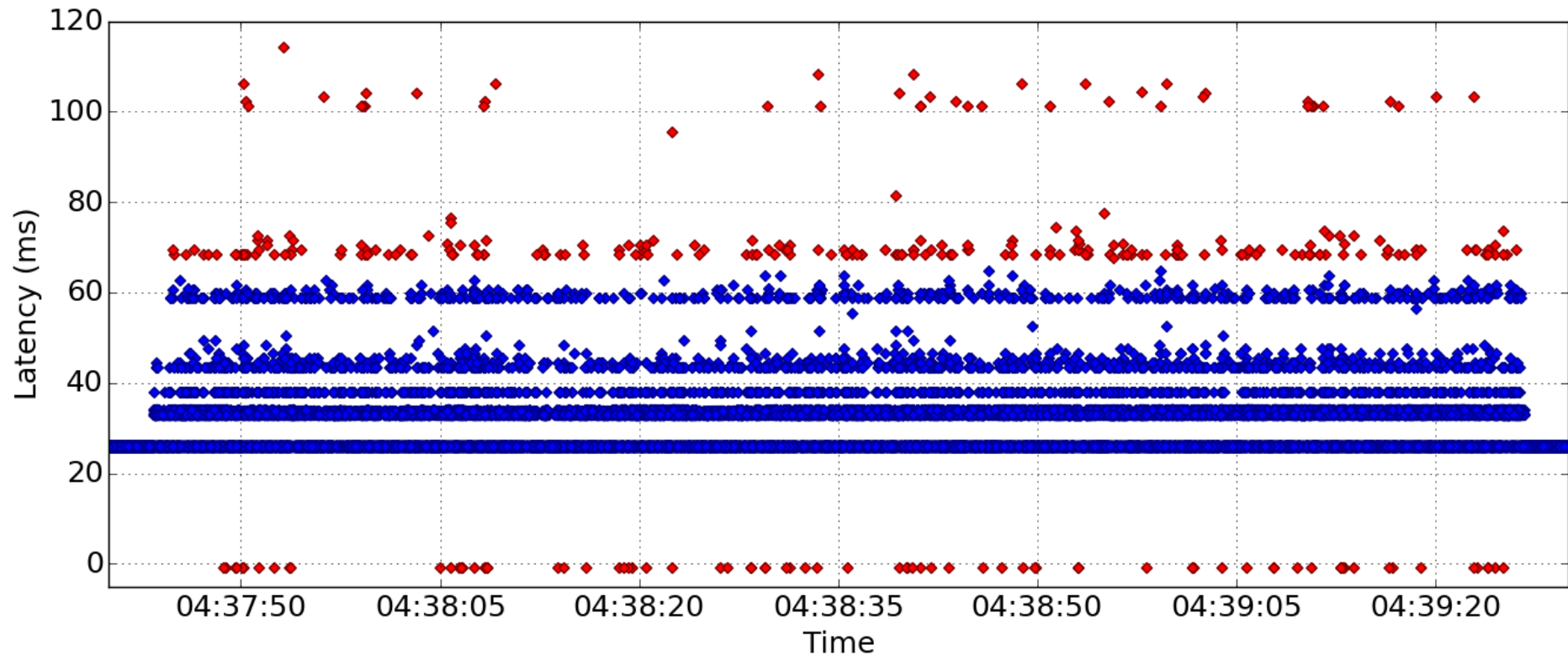
- Case study: Atlanta -> Los Angeles; August 15, 2016



Packets received and dropped over a 110-second interval using dynamic two disjoint paths (5,100 lost/late packets, 15 packets with latency over 120ms not shown)

# Dissemination Graphs Case Study: Targeted Redundancy

- Case study: Atlanta -> Los Angeles; August 15, 2016



Packets received and dropped over a 110-second interval using our dissemination-graph-based approach to add targeted redundancy at the destination (338 lost/late packets)

# Dissemination Graphs with Targeted Redundancy: Results

- 4 weeks of data collected over 4 months

  - Packets sent on each link in the overlay topology every 10ms

- Analyzed 16 transcontinental flows

  - All combinations of 4 cities on the East Coast of the US (NYC, JHU, WAS, ATL) and 2 cities on the West Coast of the US (SJC, LAX)

  - 1 packet/ms simulated sending rate

# Dissemination Graphs with Targeted Redundancy: Results

| Routing Approach | Availability (%) | Unavailability (seconds per flow per week) | Reliability (%) | Reliability (packets lost/ late per million) |
|---|---|---|---|---|
| Time-Constrained Flooding | 99.995883% | 24.90 | 99.999863% | 1.37 |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| Single Path | 99.994286% | 34.56 | 99.997710% | 22.90 |

# Dissemination Graphs with Targeted Redundancy: Results

| Routing Approach | Availability (%) | Unavailability (seconds per flow per week) | Reliability (%) | Reliability (packets lost/ late per million) |
|---|---|---|---|---|
| Time-Constrained Flooding | 99.995883% | 24.90 | 99.999863% | 1.37 |
| | | | | |
| | | | | |
| Static Two Disjoint Paths | 99.995266% | 28.63 | 99.998438% | 15.62 |
| Redundant Single Path | 99.995223% | 28.89 | 99.998715% | 12.85 |
| Single Path | 99.994286% | 34.56 | 99.997710% | 22.90 |

# Dissemination Graphs with Targeted Redundancy: Results

| Routing Approach | Availability (%) | Unavailability (seconds per flow per week) | Reliability (%) | Reliability (packets lost/ late per million) |
|---|---|---|---|---|
| Time-Constrained Flooding | 99.995883% | 24.90 | 99.999863% | 1.37 |
| | | | | |
| Dynamic Two Disjoint Paths | 99.995676% | 26.15 | 99.999103% | 8.97 |
| Static Two Disjoint Paths | 99.995266% | 28.63 | 99.998438% | 15.62 |
| Redundant Single Path | 99.995223% | 28.89 | 99.998715% | 12.85 |
| Single Path | 99.994286% | 34.56 | 99.997710% | 22.90 |

# Dissemination Graphs with Targeted Redundancy: Results

| Routing Approach | Availability (%) | Unavailability (seconds per flow per week) | Reliability (%) | Reliability (packets lost/ late per million) |
|---|---|---|---|---|
| Time-Constrained Flooding | 99.995883% | 24.90 | 99.999863% | 1.37 |
| Dissemination Graphs with Targeted Redundancy | 99.995864% | 25.02 | 99.999849% | 1.51 |
| Dynamic Two Disjoint Paths | 99.995676% | 26.15 | 99.999103% | 8.97 |
| Static Two Disjoint Paths | 99.995266% | 28.63 | 99.998438% | 15.62 |
| Redundant Single Path | 99.995223% | 28.89 | 99.998715% | 12.85 |
| Single Path | 99.994286% | 34.56 | 99.997710% | 22.90 |

# Results: % of Performance Gap Covered (between TCF and Single Path)

| Routing Approach | Week 1 2016-07-19 | Week 2 2016-08-08 | Week 3 2016-09-01 | Week 4 2016-10-13 | Overall | Scaled Cost |
|---|---|---|---|---|---|---|
| Time-Constrained Flooding | 100.00% | 100.00% | 100.00% | 100.00% | **100.00%** | **14.350** |
| Dissem. Graphs with Targeted Redundancy | 94.19% | 99.19% | 98.00% | 99.50% | **98.97%** | **2.203** |
| Dynamic Two Disjoint Paths | 80.91% | 71.34% | 47.73% | 73.46% | **70.74%** | **2.197** |
| Static Two Disjoint Paths | -76.72% | 50.89% | 53.58% | 40.79% | **39.50%** | **2.194** |
| Redundant Single Path | 54.12% | 37.25% | 4.89% | 59.10% | **45.75%** | **2.000** |
| Single Path | 0.00% | 0.00% | 0.00% | 0.00% | **0.00%** | **1.000** |

# Applications: Remote Manipulation



Video demonstration: [www.dsn.jhu.edu/~babay/Robot_video.mp4](www.dsn.jhu.edu/~babay/Robot_video.mp4)

# Applications: Remote Robotic Ultrasound

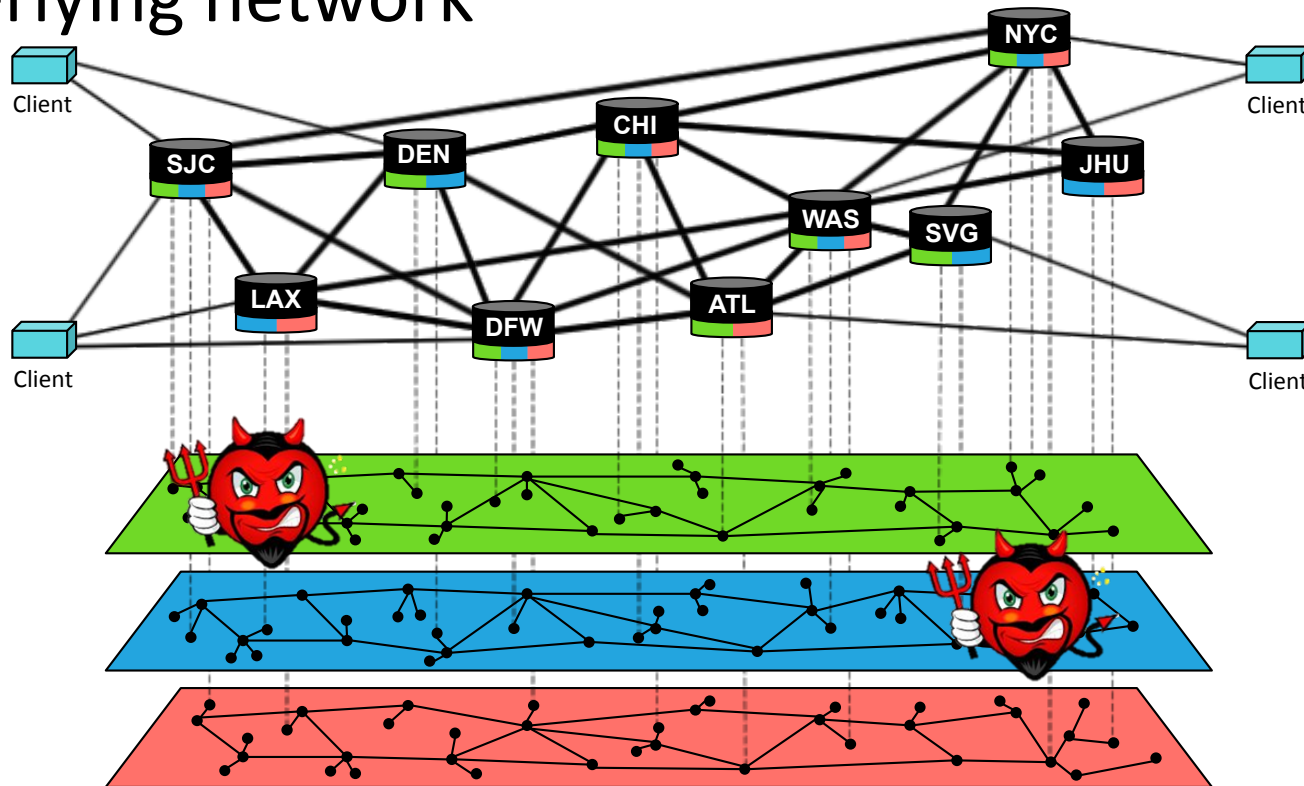- Collaboration with JHU/TUM CAMP lab (https://camp.lcsr.jhu.edu/)

# Outline

- A New Generation of Internet Services
- The Structured Overlay Network Vision
  – Resilient network architecture
  – Overlay node software architecture with global state and unlimited programmability
  – Flow-based processing
- First Steps and Benefits
  – Responsive overlay routing with a resilient network architecture
  – Hop-by-hop reliability with flow-based processing and unlimited programmability
- The Quest for QoS
  – Almost-reliable real-time protocol for VoIP
  – Almost-reliable real-time protocol for Live TV
- Going even Faster
  – Remote manipulation, remote robotic surgery, collaborative virtual reality
  – Dissemination graphs with targeted redundancy
- **Resilient Communication in a Hostile World**
  – **Intrusion-tolerant networking via structured overlays**
  – **Critical infrastructure applications**
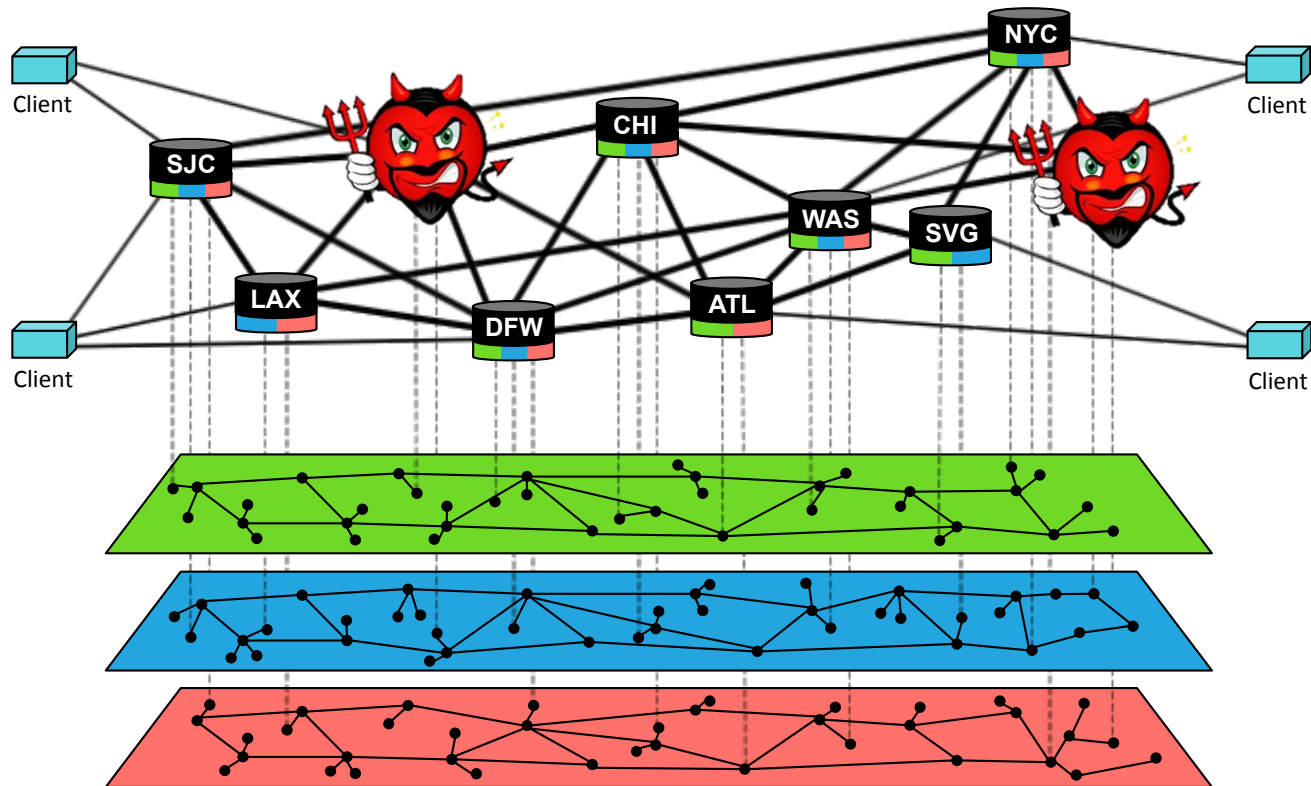- Future Directions

# Intrusion-Tolerant Networks via Structured Overlays

- Resilient network architecture + responsive overlay routing protects against compromises in the underlying network

# Intrusion-Tolerant Networks via Structured Overlays

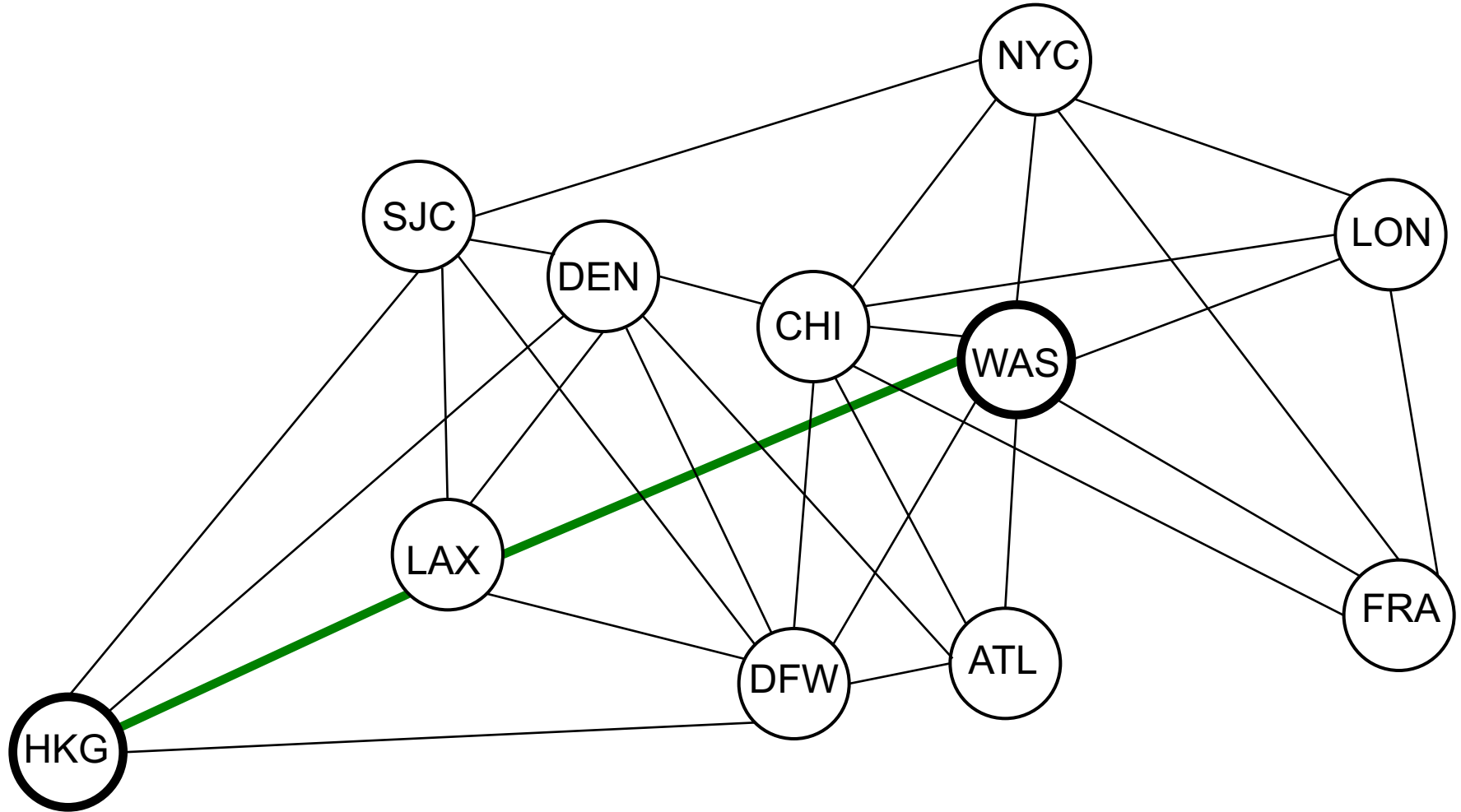- Intrusion-tolerant overlay protocols protect against overlay node compromises

# Intrusion-Tolerant Networks
# via Structured Overlays

- Intrusion-tolerant overlay protocols protect against overlay node compromises
  - Authorized nodes are known in advance and authenticated (maximal topology with minimal weights)
  - Redundant dissemination (k node-disjoint paths or constrained flooding)
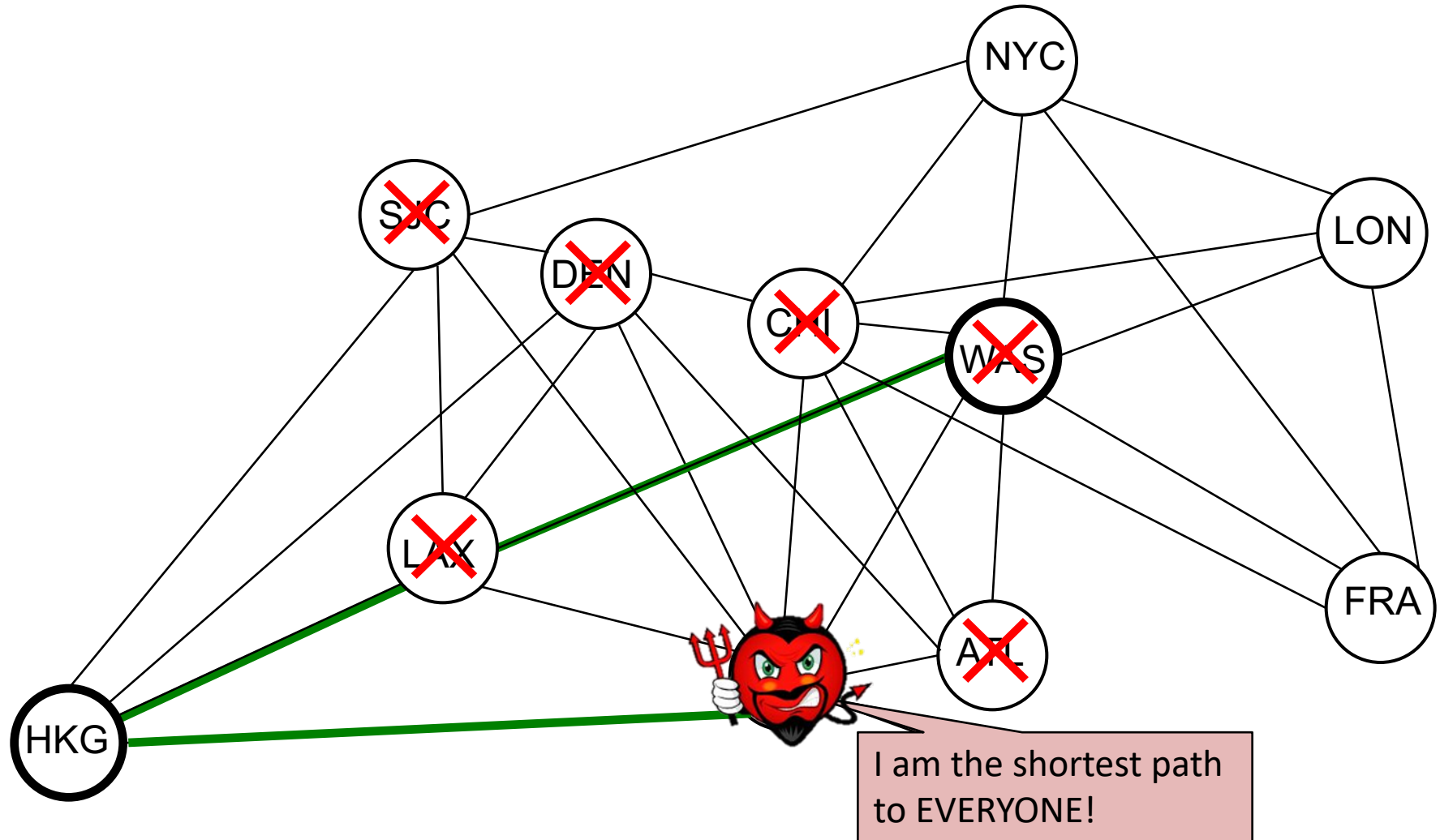  - Source- or flow-based fairness in resource allocation

"Practical Intrusion-Tolerant Networks", D. Obenshain, T. Tantillo, A. Babay, J. Schultz, A. Newell, Md. E. Hoque, Y. Amir, C. Nita-Rotaru,
*IEEE International Conference on Distributed Computing Systems (ICDCS),* 2016
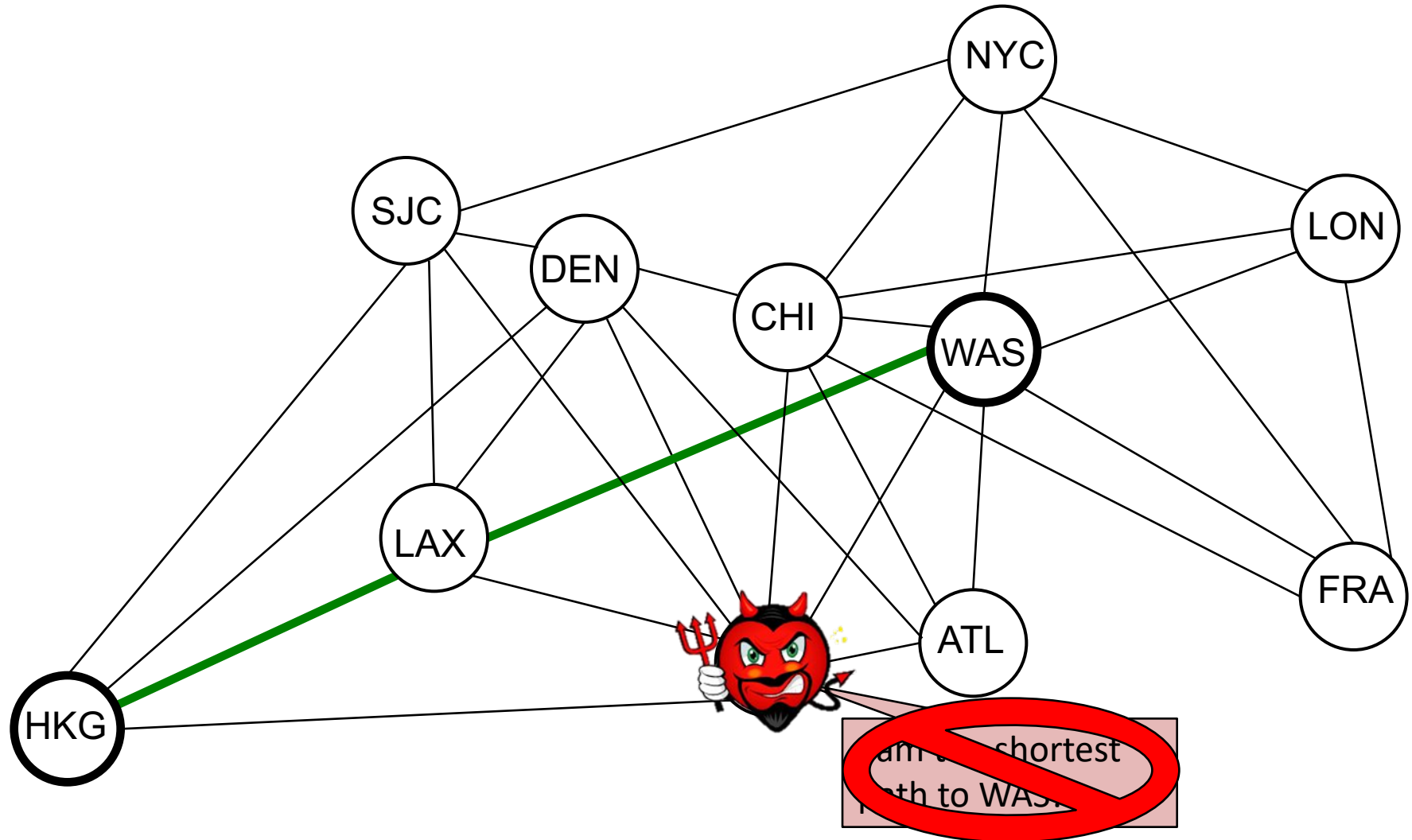
# Regular Secure Routing



Regular secure routing takes the shortest path from source (HKG) to destination (WAS).

# Regular Secure Routing Under Attack
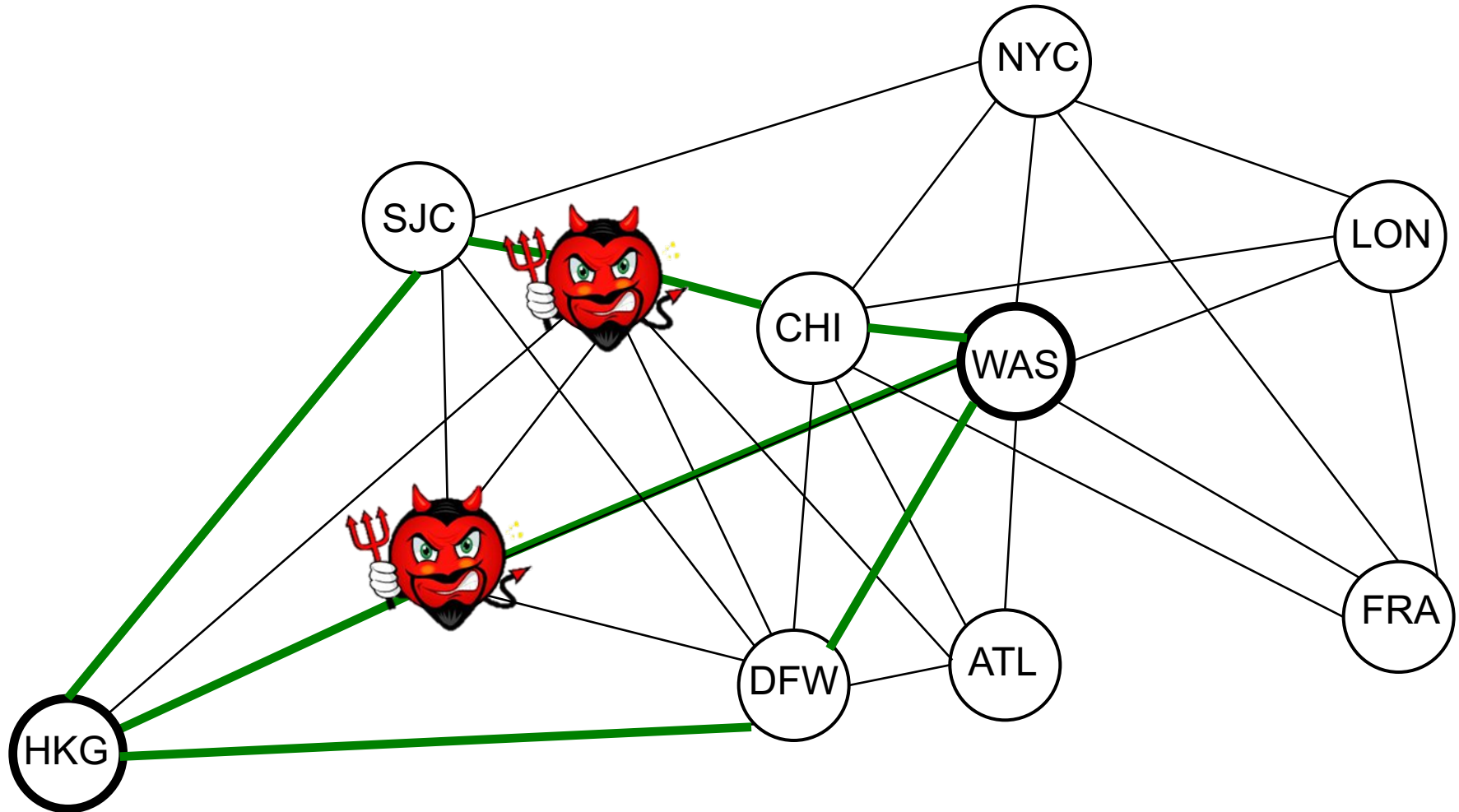


I am the shortest path to EVERYONE!

A compromised node can lie and attract traffic, which can then be dropped.
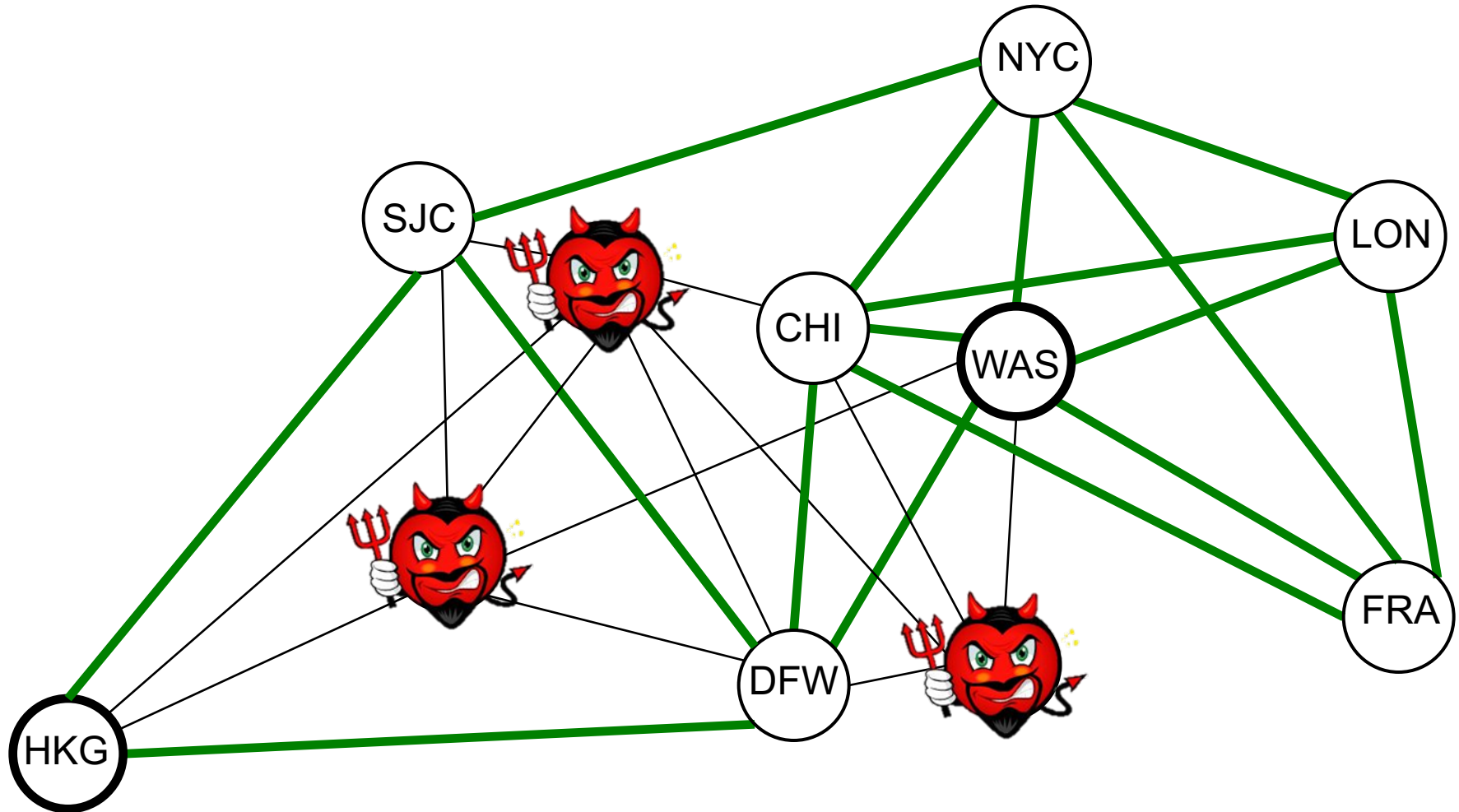
# Maximal Topology with Minimal Weights



- The nodes and edges in the topology are known ahead of time
- No node can advertise weights below the minimal weights – attack defeated
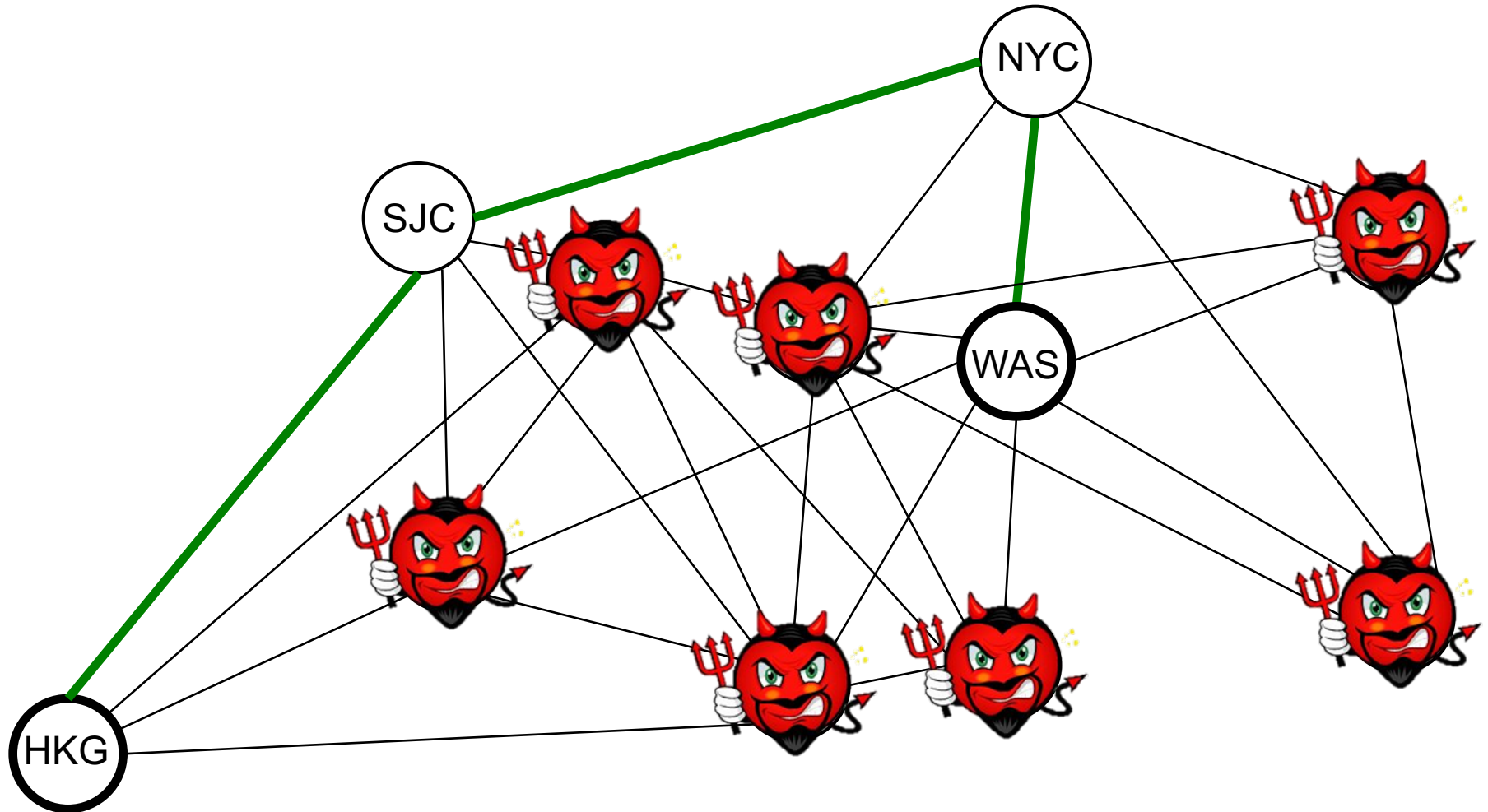
# K Node-Disjoint Paths



K node-disjoint paths defends against K-1 compromised nodes.
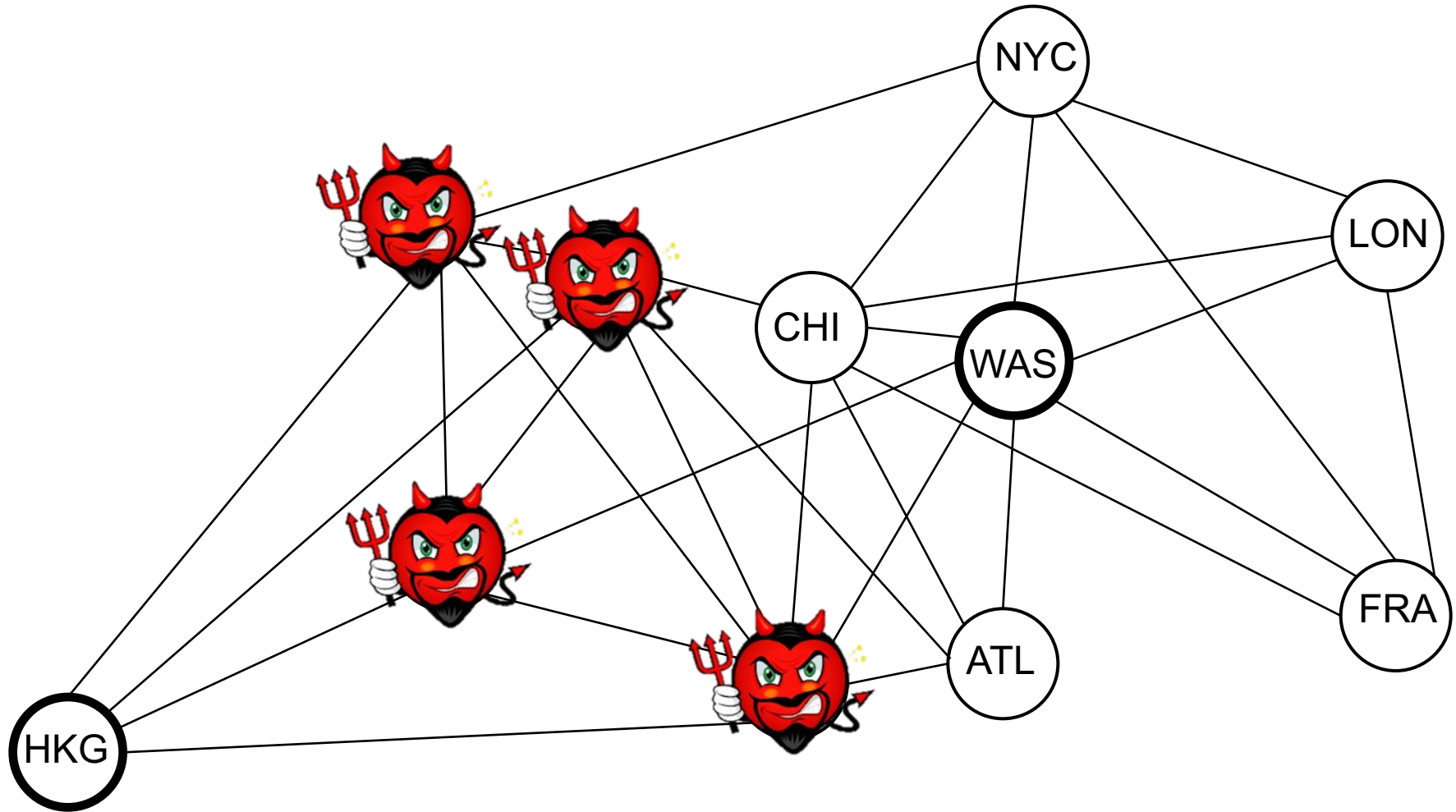
# Constrained Flooding



Flooding across the **overlay** network provides optimal resiliency.
Costs more, but we're willing to pay for the most important messages.

# Constrained Flooding



If even a single good path exists, constrained flooding will pass messages from source to destination in a timely manner.
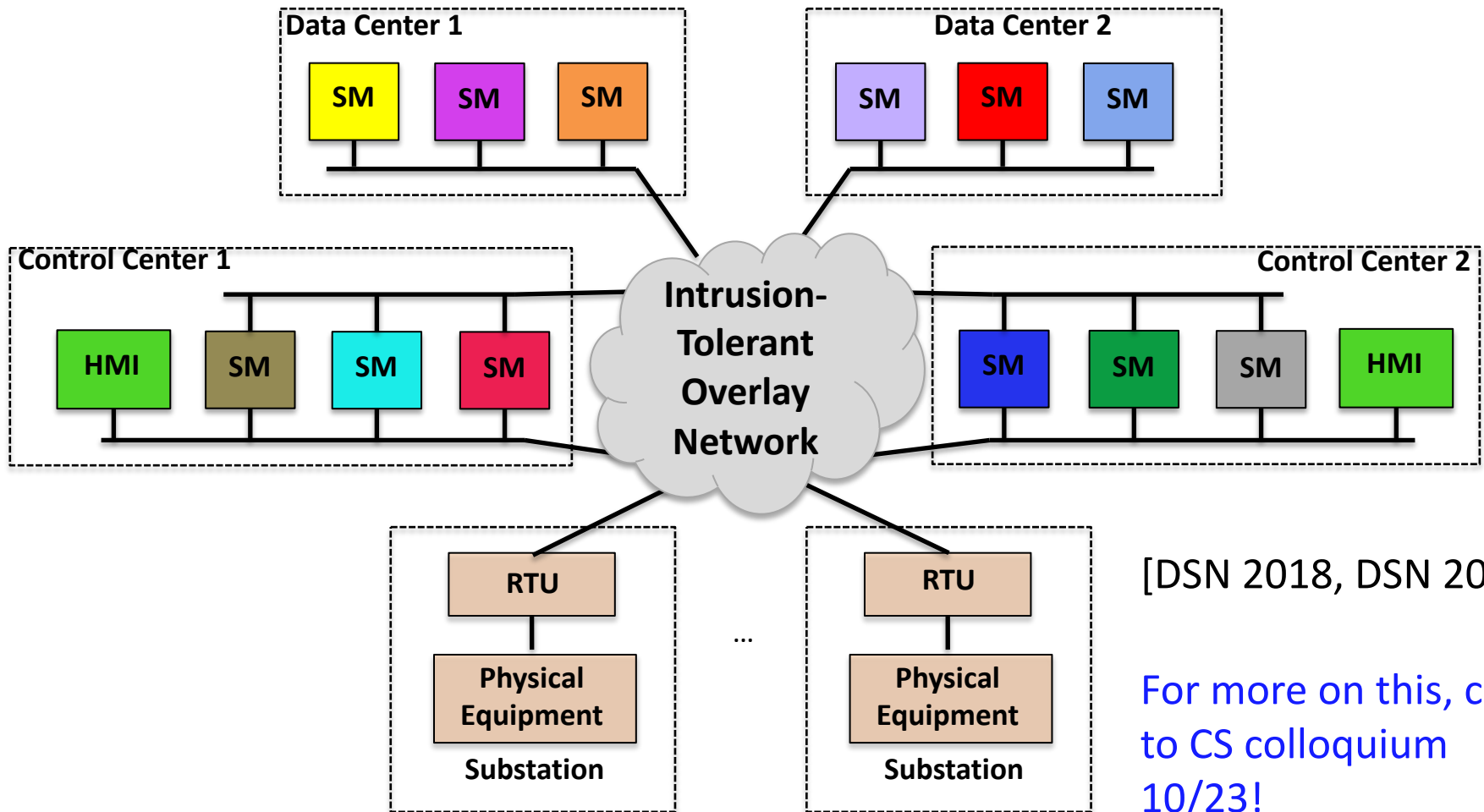
If the compromised nodes cut the network, no protocol can succeed.

# Critical Infrastructure Applications: SCADA for the Power Grid

- **Intrusion-tolerant overlay network** provides the communication foundation for our intrusion-tolerant <u>SCADA system</u> for the power grid

- **Supervisory Control and Data Acquisition (SCADA)** systems monitor and control critical infrastructure services

- SCADA system failures and downtime can cause **catastrophic consequences** (equipment damage, blackouts, human casualties)

- **Perimeter defenses** are **not sufficient** against determined attackers
  - Stuxnet, Dragonfly/Energetic Bear, Black energy (Ukraine 2015), Crashoverride (Ukraine 2016)

# Intrusion-Tolerant SCADA for the Power Grid



[DSN 2018, DSN 2019]

For more on this, come to CS colloquium 10/23!

# Outline

- A New Generation of Internet Services
- The Structured Overlay Network Vision
  - Resilient network architecture
  - Overlay node software architecture with global state and unlimited programmability
  - Flow-based processing
- First Steps and Benefits
  - Responsive overlay routing with a resilient network architecture
  - Hop-by-hop reliability with flow-based processing and unlimited programmability
- The Quest for QoS
  - Almost-reliable real-time protocol for VoIP
  - Almost-reliable real-time protocol for Live TV
- Going even Faster
  - Remote manipulation, remote robotic surgery, collaborative virtual reality
  - Dissemination graphs with targeted redundancy
- Resilient Communication in a Hostile World
  - Intrusion-tolerant networking via structured overlays
  - Critical infrastructure applications
- **Future Directions**

# Unlimited Network Programmability at Scale

- **New generation of Internet services**
  - Low-latency interactivity                    [ICDCS 2017 – Best paper]
  - High-performance reliability
  - Flow processing, transformation, analytics
  - Resilience, security, access control       [ICDCS 2016, DSN 2018]
- **Unlimited programmability at scale**
  - Structured Overlays: put general-purpose application-level processing into the middle of the network
  - Software Defined Networking: enables line speed classification and redirection
  - Combine to enable sophisticated new Internet services at scale