# Topology Switching for Data Center Networks

**Kevin Webb, Alex Snoeren, Ken Yocum**

**UC San Diego Computer Science**

**March 29, 2011**
**Hot-ICE 2011**
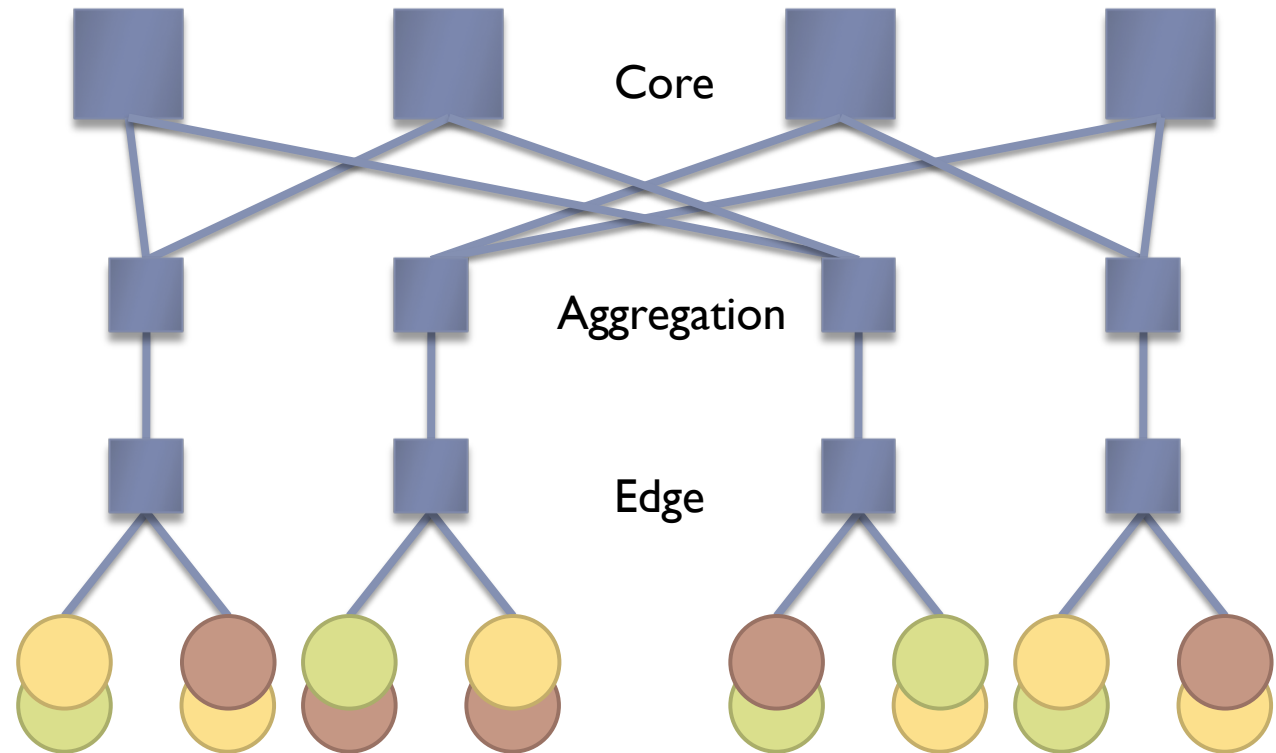
**UCSDCSE**
Computer Science and Engineering
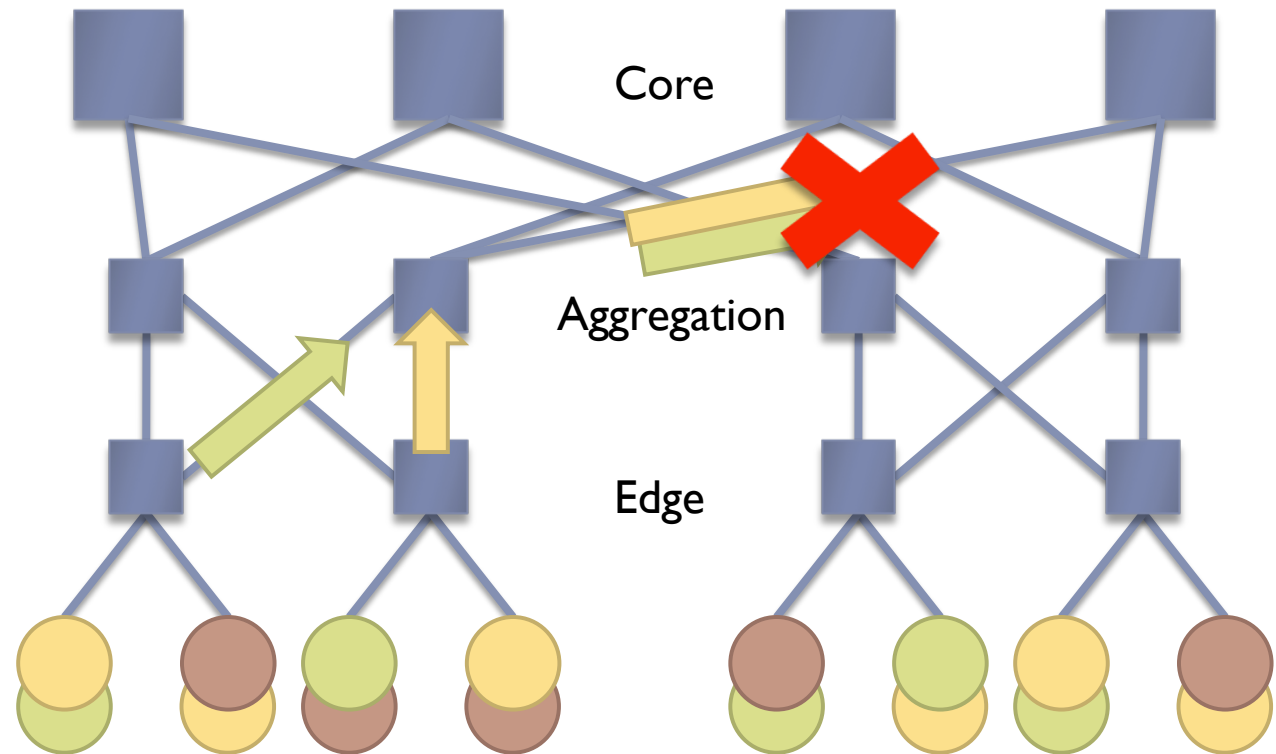
# Data Center Networks

▸ Hosting myriad of applications:

  ▸ Big data: MapReduce

  ▸ Web services

  ▸ HPC: MPI

  ▸ DB, Storage

  ▸ Many others!
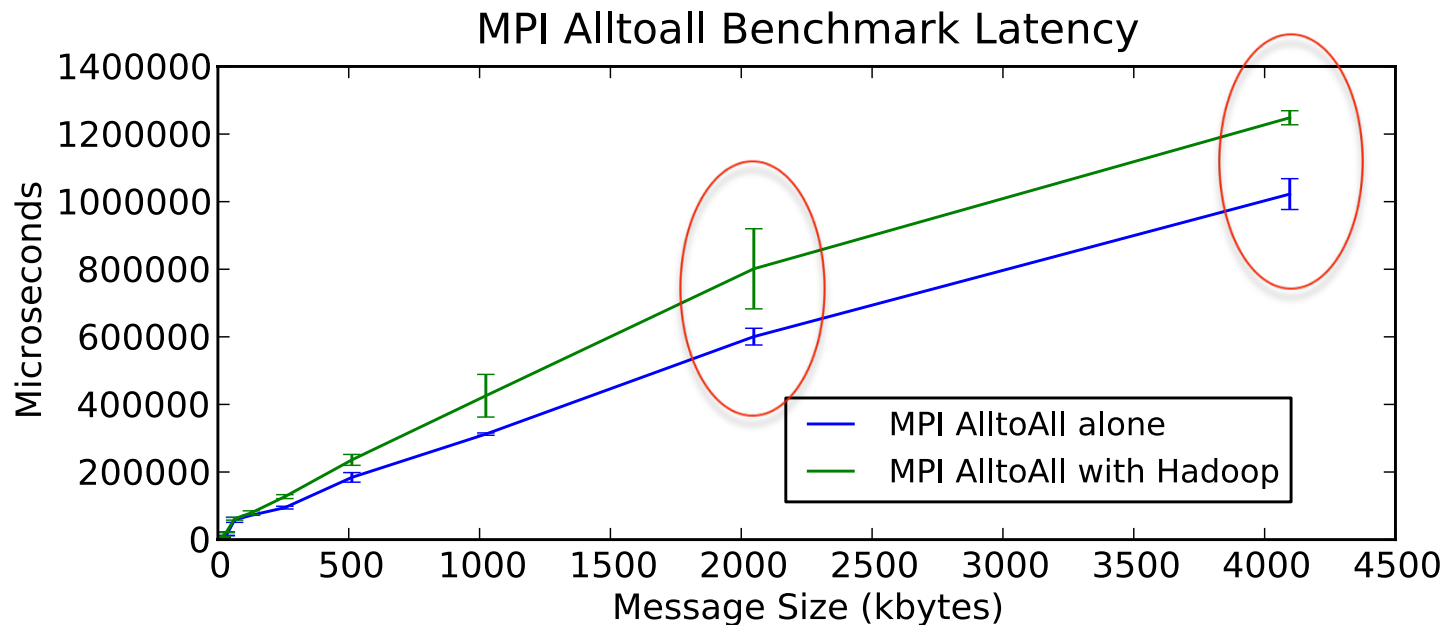
# Data Center Networks

- ## DC engineers adding links
  - Applications need other important characteristics!

# Inter-application Interference



MPI Alltoall Benchmark Latency

▸ Experiment

   ▸ All to all MPI and Hadoop data processing

   ▸ Openflow ECMP network

   ▸ Interference > 20% latency increase

UCSD**CSE**
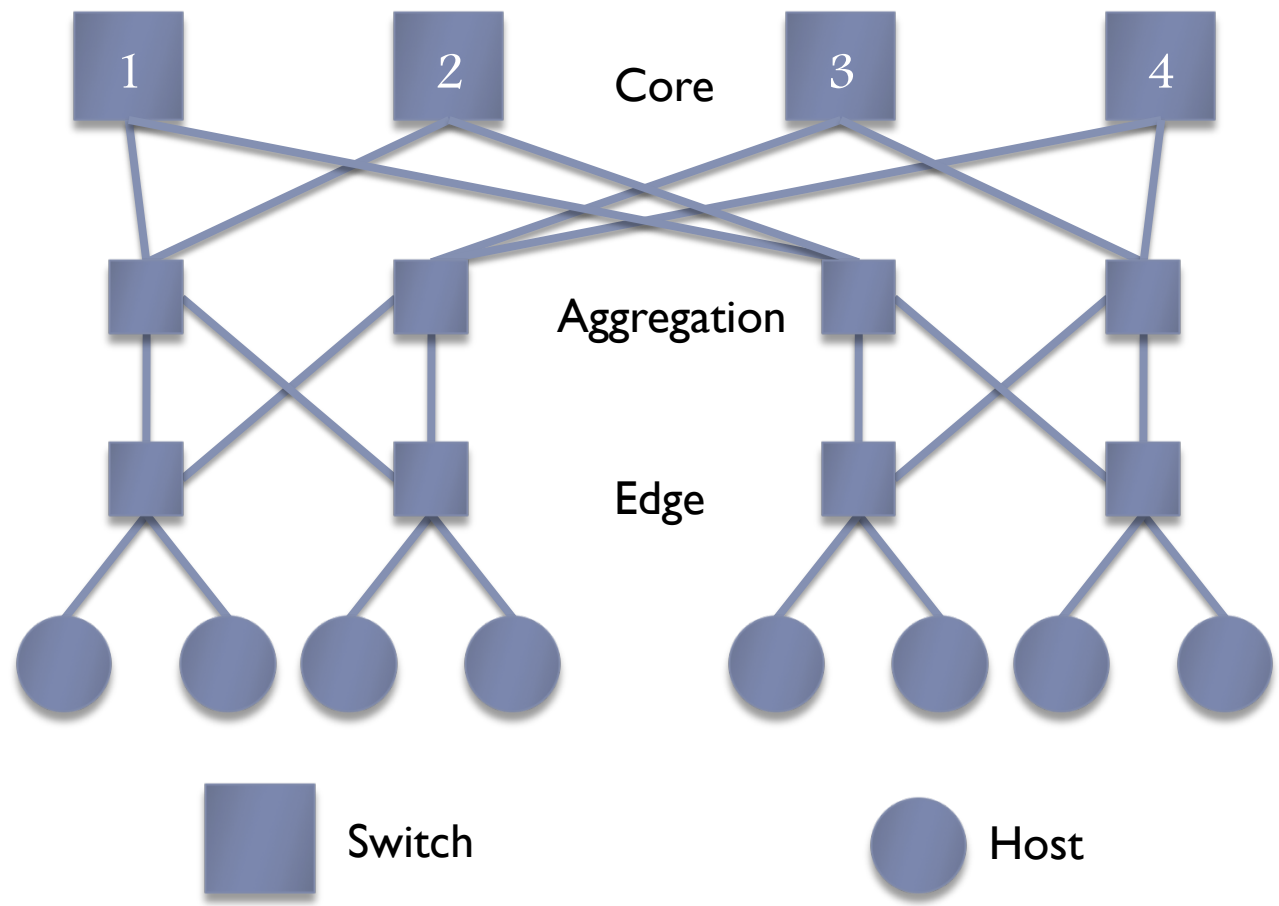Computer Science and Engineering

# Topology Switching Overview

▸ **Applications request specific characteristics**

  ▸ Bandwidth, Redundancy, Latency, Isolation, others…

▸ **Idea: Create routes based on applications' needs**

  ▸ Per application instance
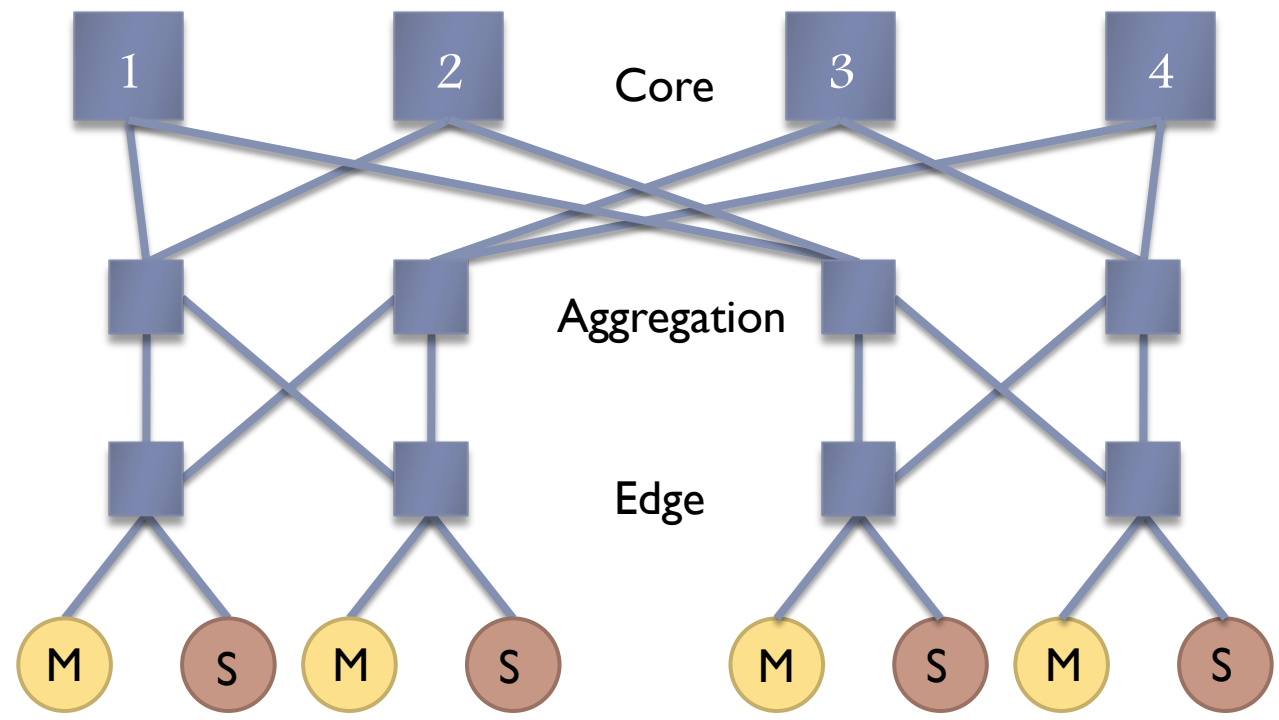
  ▸ Per application phase: Hadoop shuffle vs. HDFS writes

# Example

# Example
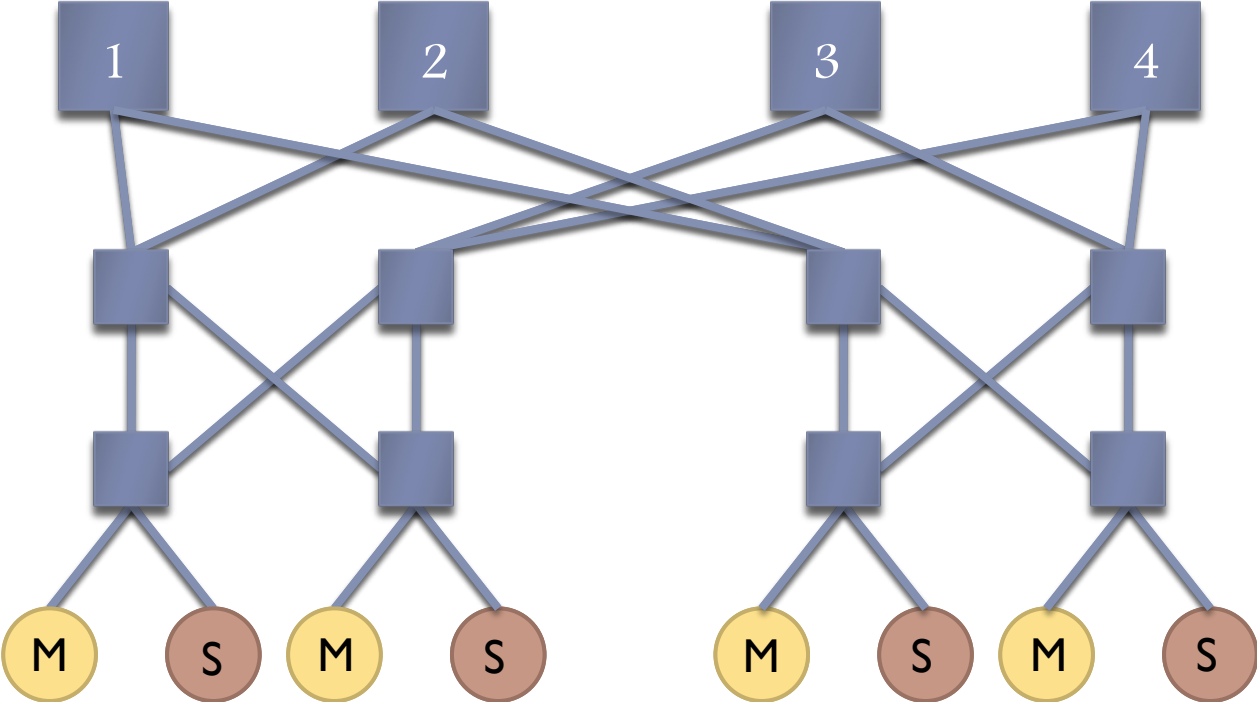
# Example



M     Map Reduce Task (Needs throughput)

S     Scientific Task (MPI - Needs isolation for consistency)
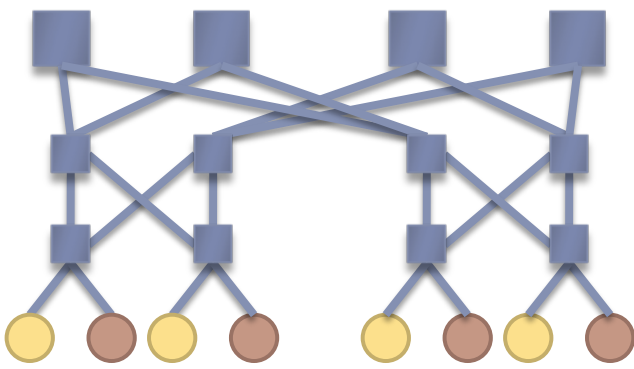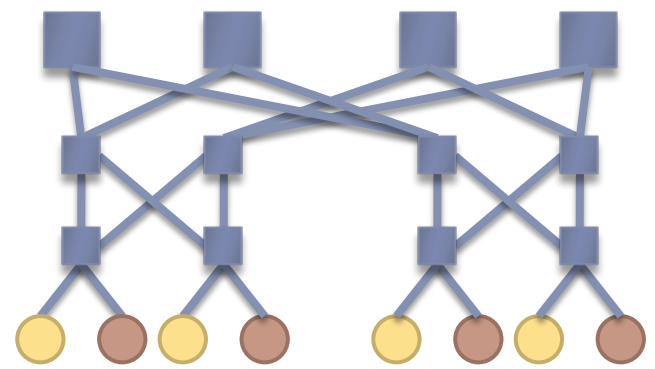
M   S   M   S     M   S   M   S

# Example

# Example



Scientific Network
Exclusive
Free from interference

Map Reduce Network
Multiple paths - high capacity

# Challenges

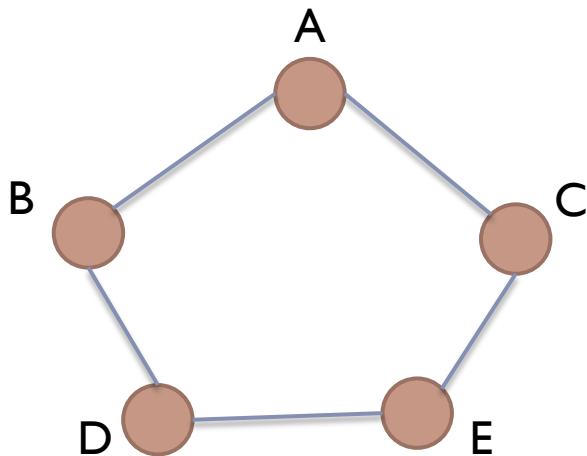- Reconfigurable network infrastructure
  - Frequent allocations
  - Flexible routing rules
  - Openflow

- Allocation algorithms
  - Throughput, Reliability, Isolation, etc.
  - Evaluation metrics

- Cooperative online allocation of network resources
  - Limit conflict between allocations
  - Can't take too long

UCSD**CSE**
Computer Science and Engineering

# Abstraction

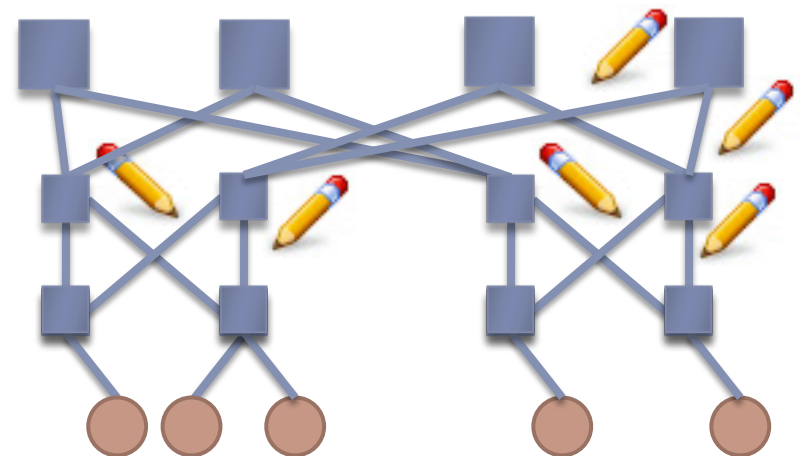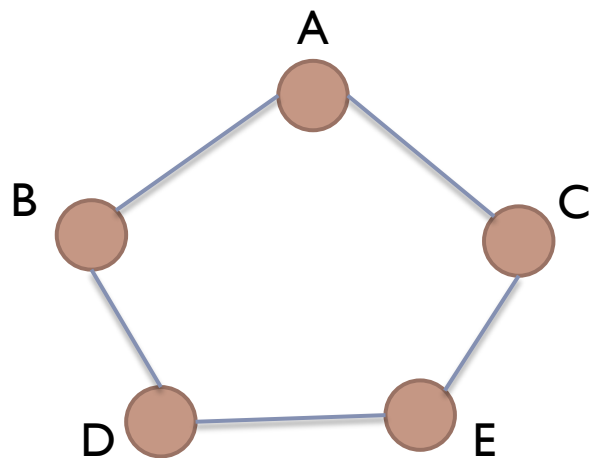▸ **Applications submit *routing tasks*:**

  ▸ Set of communicating end hosts
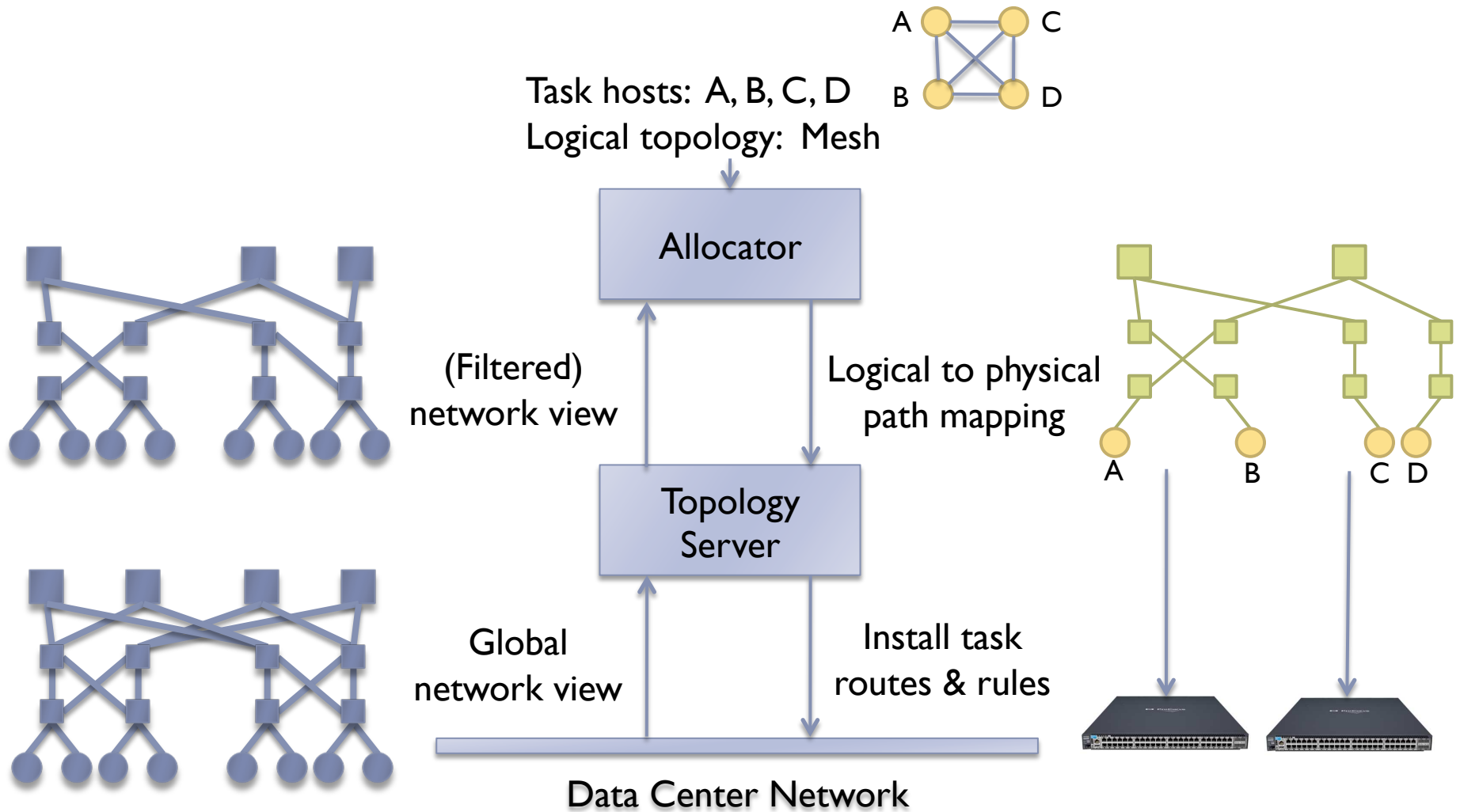
  ▸ Logical topology: mesh, ring, tree, custom

# Abstraction

- Routing tasks utilize an **allocator**:
  - Quantifiable metric
    - Guides allocation, indicates success
  - Allocation algorithm
    - Chooses physical paths
  - Graph annotation & filtering
    - Record allocation results to reduce conflict

# Tasks
# Routes

# Architecture



Task hosts: A, B, C, D
Logical topology: Mesh

Allocator

(Filtered) network view

Logical to physical path mapping

Topology Server

Global network view

Install task routes & rules

Data Center Network

A    B    C   D

# Three Allocators

▸ Bandwidth
  ▸ Finds least loaded paths to maximize capacity

▸ Resiliency
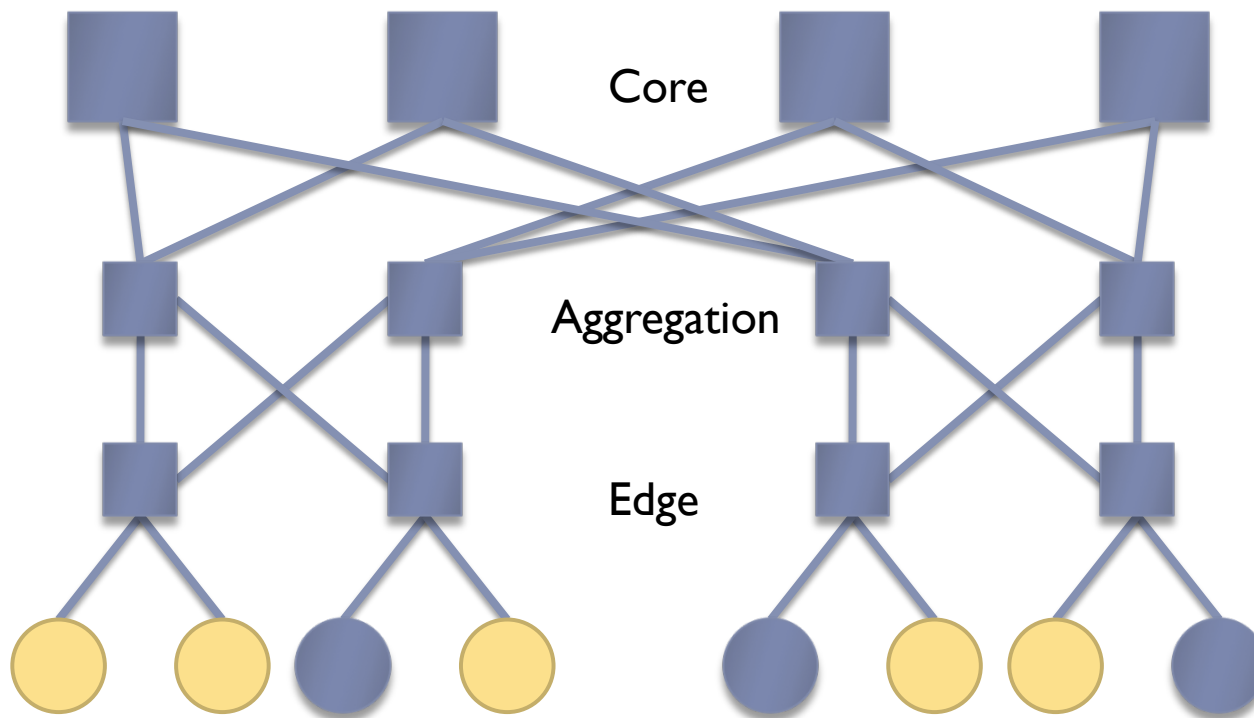  ▸ Allocates $N$ disjoint paths between every host pair

▸ $K$-Isolation
  ▸ Chooses paths with at most $k$ other tasks
  ▸ Reduces inter-task interference, more consistent

UCSDCSE
Computer Science and Engineering
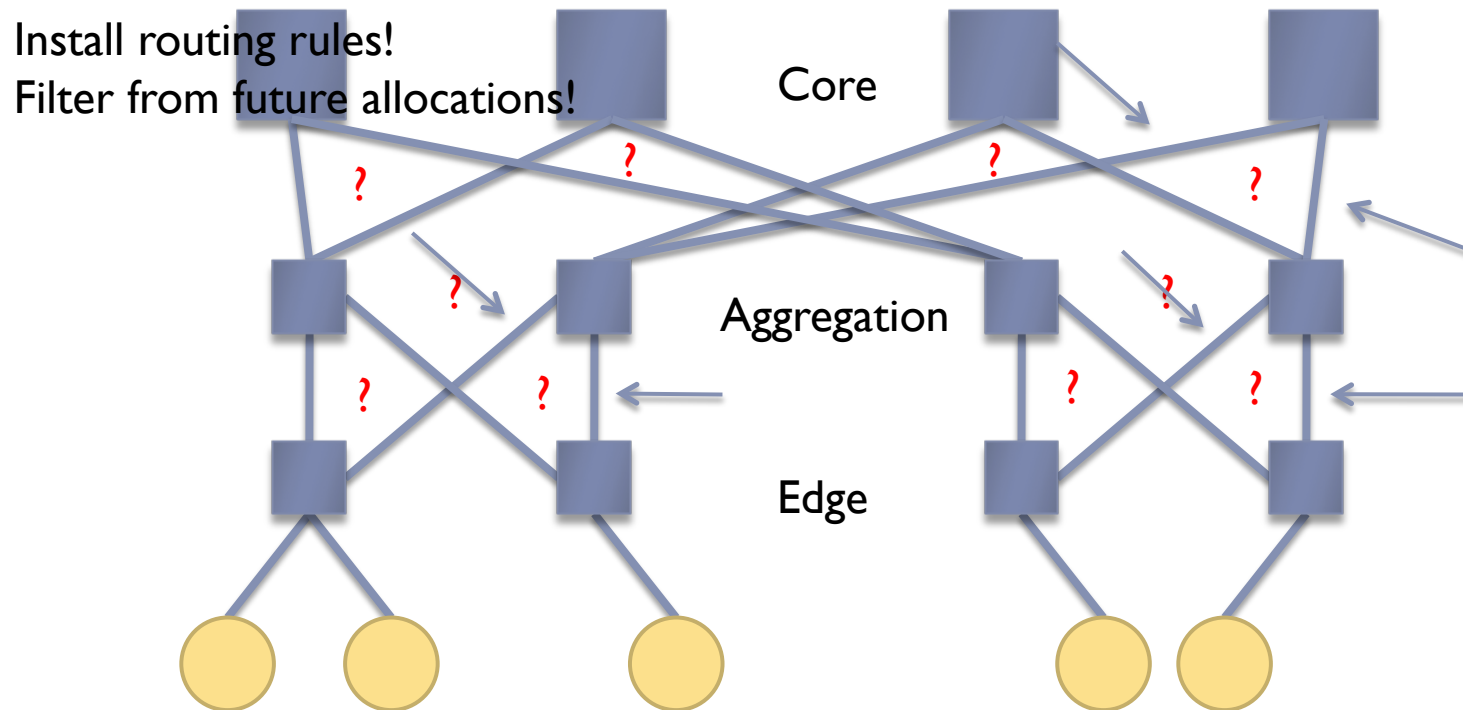
# *K*-Isolation

▸ Goal: Don't share links with more than *K* other tasks

# *K*-Isolation

▸ Goal: Don't share links with any other tasks

Install routing rules!
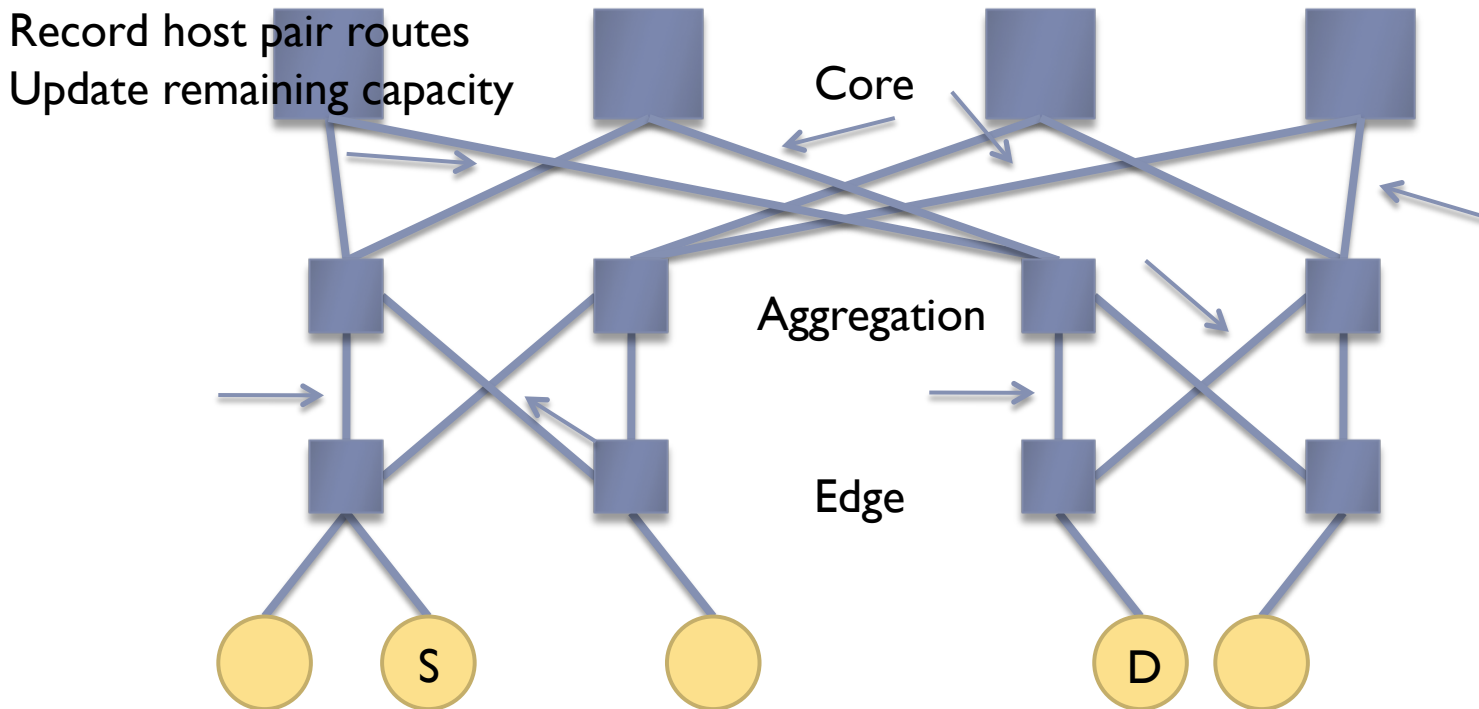Filter from future allocations!

Core

Aggregation

Edge

# Resiliency

▸ Goal: Find *N* disjoint routes between end host pairs



Core

Aggregation

Edge

# Resiliency

▸ Goal: Find **2** disjoint routes between end host pairs



Record host pair routes
Update remaining capacity

Core

Aggregation

Edge

S

D

# Resiliency
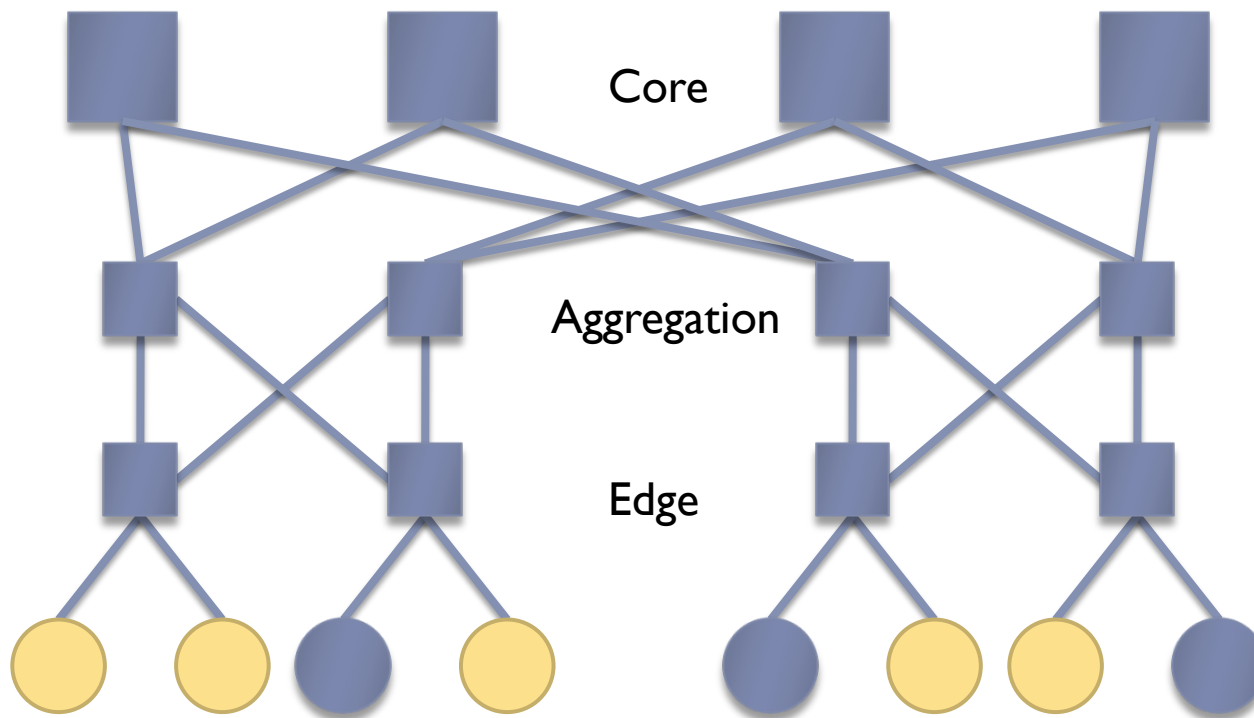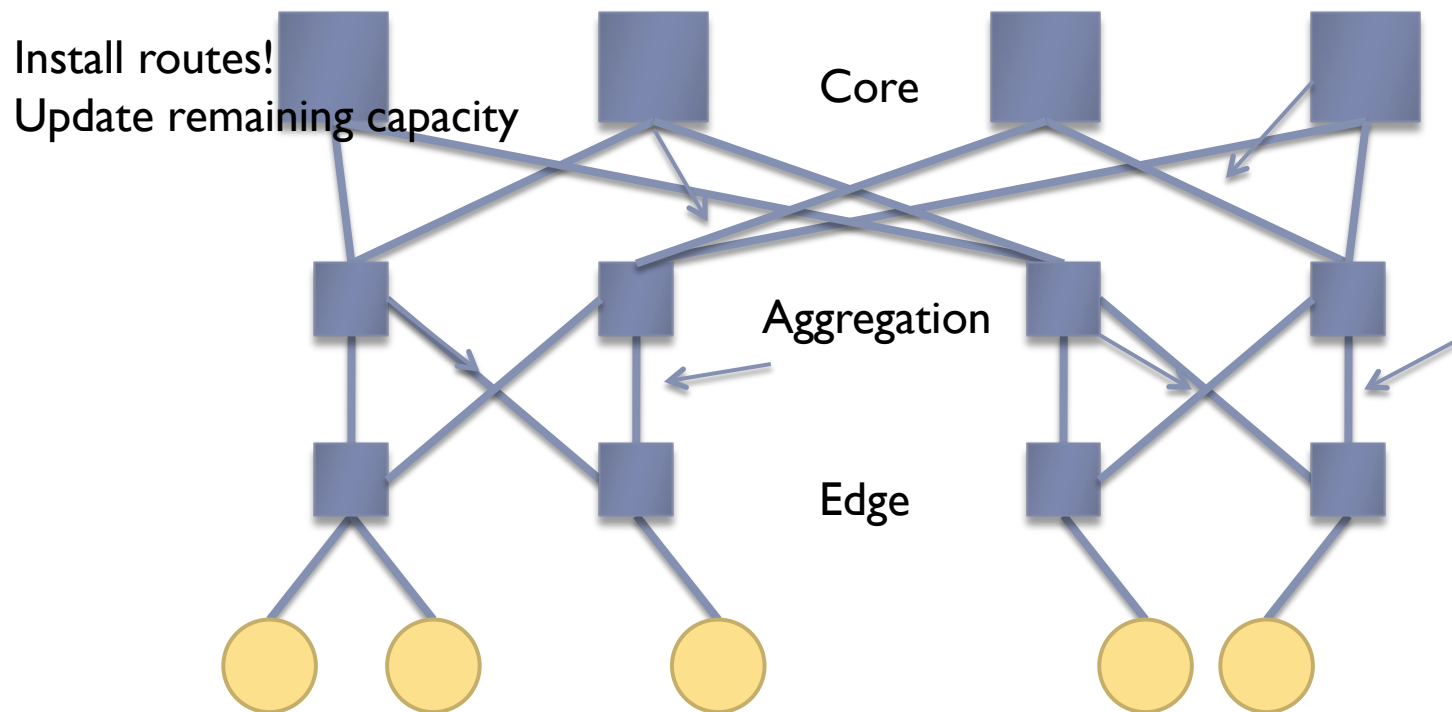
▸ Repeat for remaining source/destination pairs…

# Bandwidth

▸ Goal: Maximize bisection bandwidth between hosts

# Bandwidth

▸ Build max spanning tree over remaining capacity



Install routes!
Update remaining capacity

Core

Aggregation

Edge

# Simulations

▸ Does it work?

▸ Comparison against state of the art:
  ▸ "Optimal" equal-cost paths for Resiliency/Bandwidth tasks
  ▸ Spanning tree VLAN for isolation tasks

▸ Two distinct workloads:
  ▸ Balanced - 6 tasks: 2 isolation, 2 resilience, 2 bandwidth
  ▸ Stressed - 16 tasks: 7 isolation, 5 resilience, 4 bandwidth

# Results

Isolation metric: [0, 1]

Value indicates average path isolation.

1 – Completely isolated

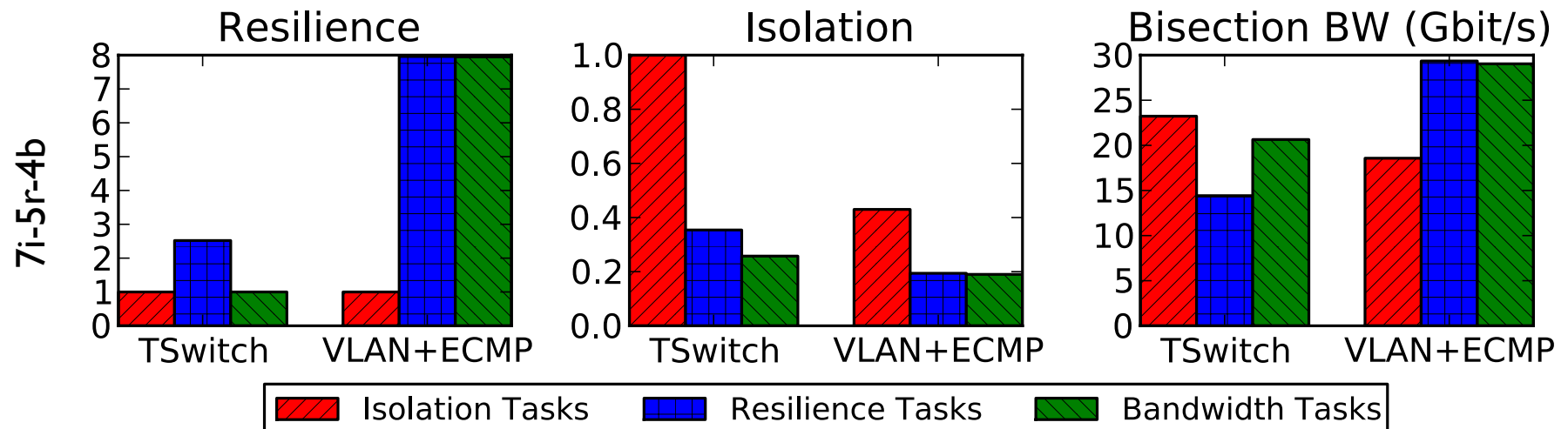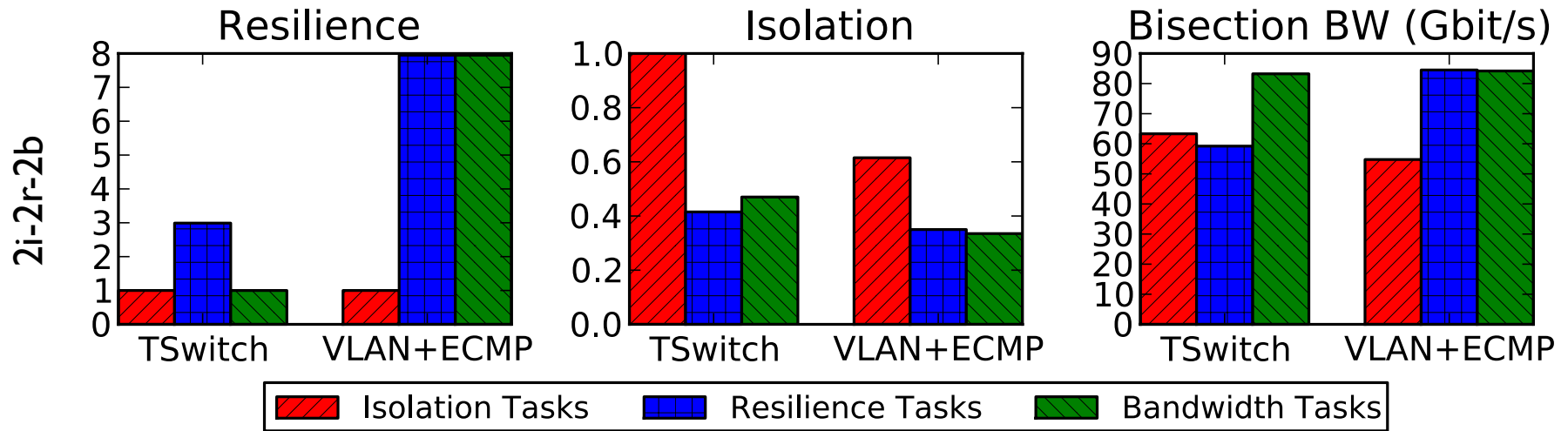Resiliency

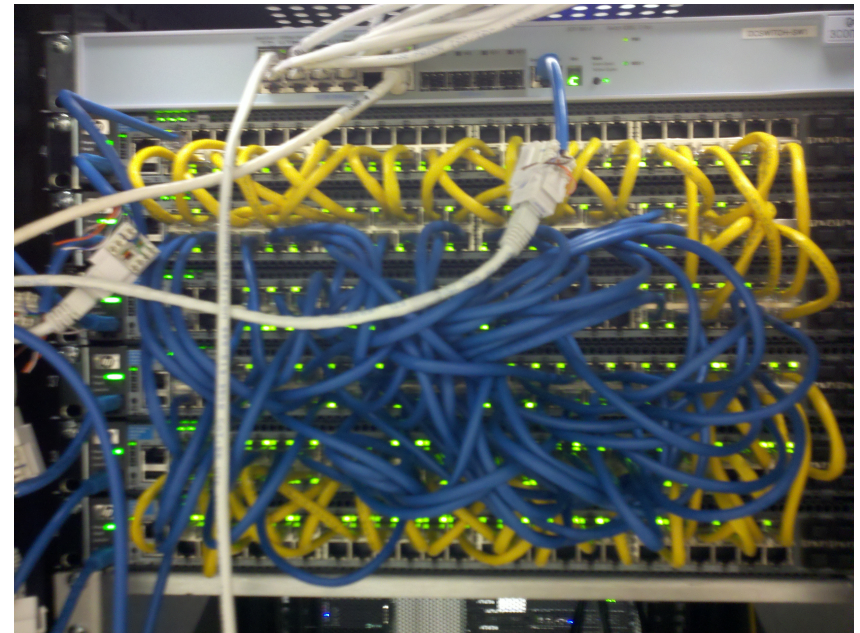Number of disjoint paths between hosts.

# Results

# Current Status

▸ **Simulations promising!**

▸ **Ongoing work:**
- ▸ Quantifying interference
- ▸ Refining allocation strategies

▸ **Building architecture**
- ▸ Openflow-enabled switches
- ▸ Routing rule instantiation
  - ▸ Limited TCAM size / speed



UCSD**CSE**
Computer Science and Engineering

# Thanks!

▸ Questions?