# Lecture 12: Fat-tree and Dragonfly Networks

## Abhinav Bhatele, Department of Computer Science

UNIVERSITY OF
MARYLAND

# Summary of last lecture

- Key requirements of HPC networks

  - extremely low latency, high bandwidth, scalable

  - low network diameter, high bisection bandwidth

- Torus networks (less common now)

  - Network diameter grows as $O(\sqrt[3]{N})$ where N is the number of nodes

- Different types of routing algorithms:

  - Shortest path vs. non-minimal

  - Static vs. dynamic

# Fat-tree network

- Most popular network topology
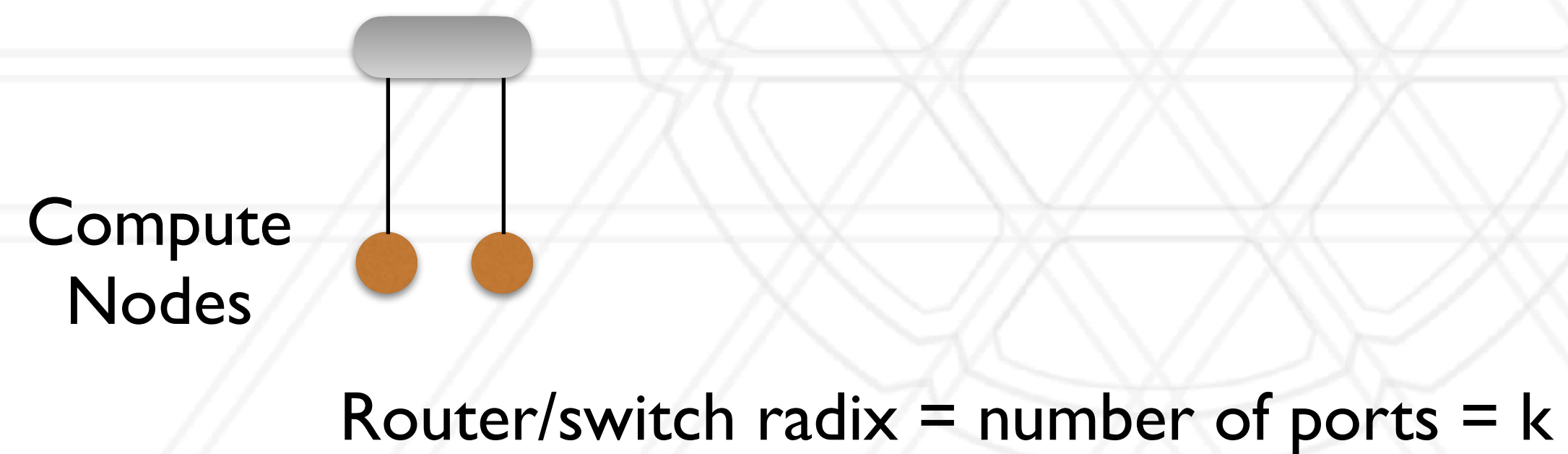
  - Low network diameter, high bandwidth

DEPARTMENT OF
COMPUTER SCIENCE

# Fat-tree network

- Most popular network topology
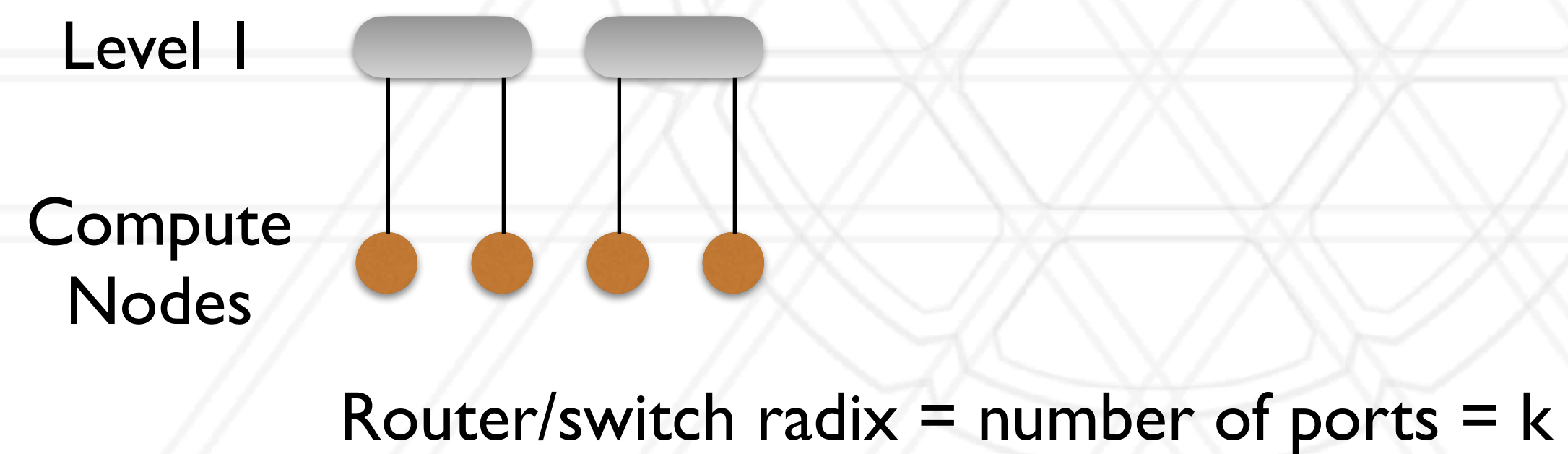
  - Low network diameter, high bandwidth

Compute
Nodes

DEPARTMENT OF
COMPUTER SCIENCE

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth

Compute
Nodes

Router/switch radix = number of ports = k

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth

Compute
Nodes

Router/switch radix = number of ports = k

DEPARTMENT OF
COMPUTER SCIENCE

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth

Level 1

Compute
Nodes

Router/switch radix = number of ports = k

DEPARTMENT OF
COMPUTER SCIENCE

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth

Level 2

Level 1

Compute
Nodes

Router/switch radix = number of ports = k

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth
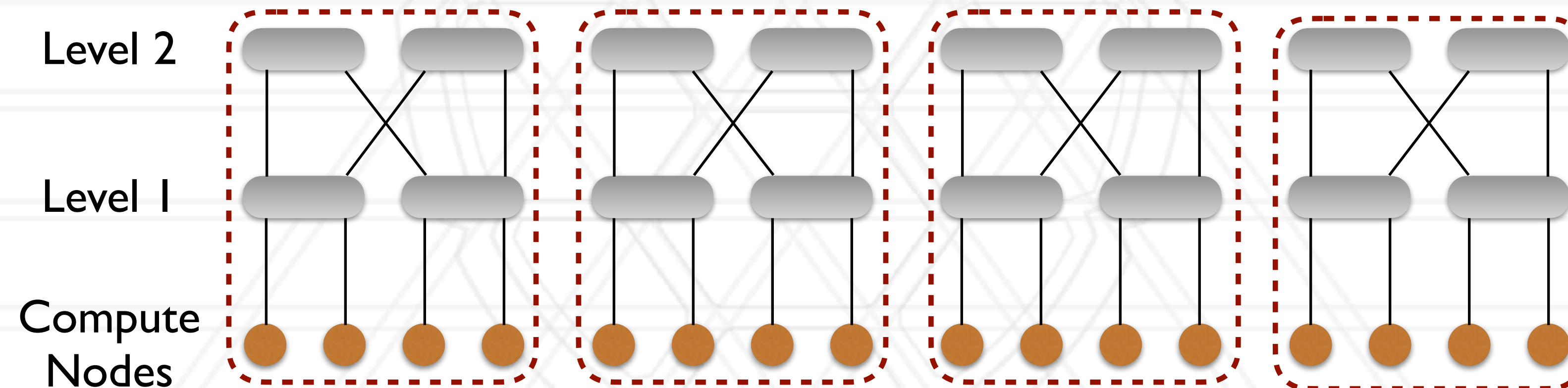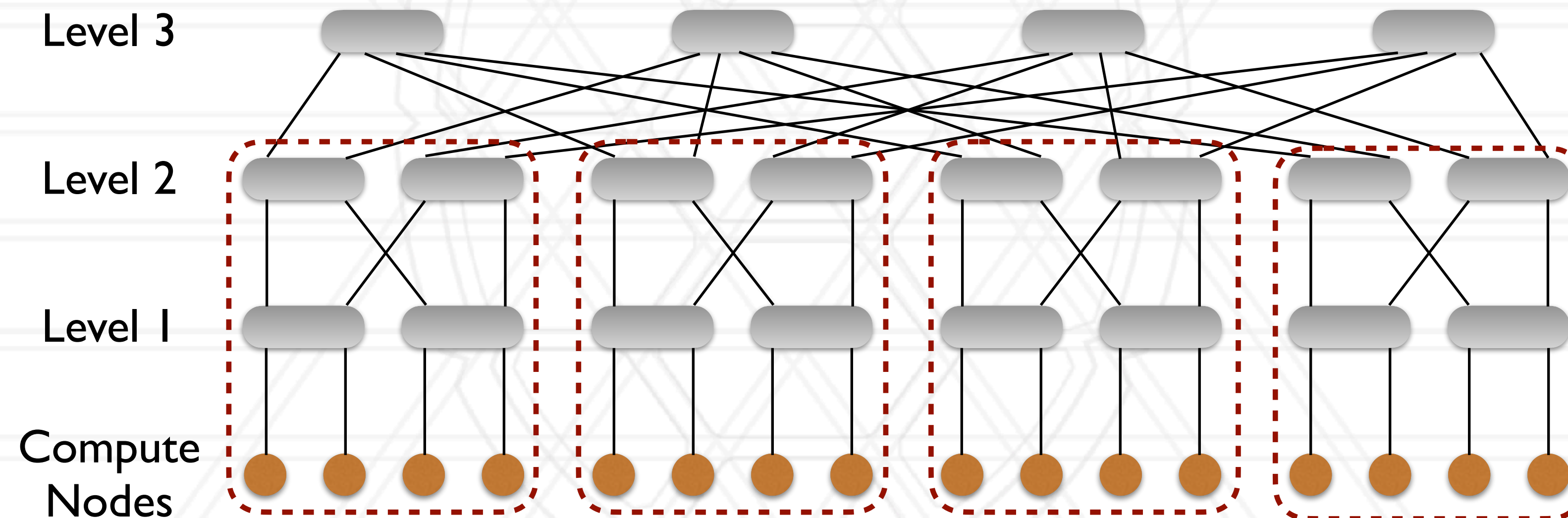


Level 2

Level 1

Compute
Nodes

Router/switch radix = number of ports = k

Pod = group of switches = k/2 switches

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth



Router/switch radix = number of ports = k

Pod = group of switches = k/2 switches

DEPARTMENT OF
COMPUTER SCIENCE

# Fat-tree network

- Most popular network topology
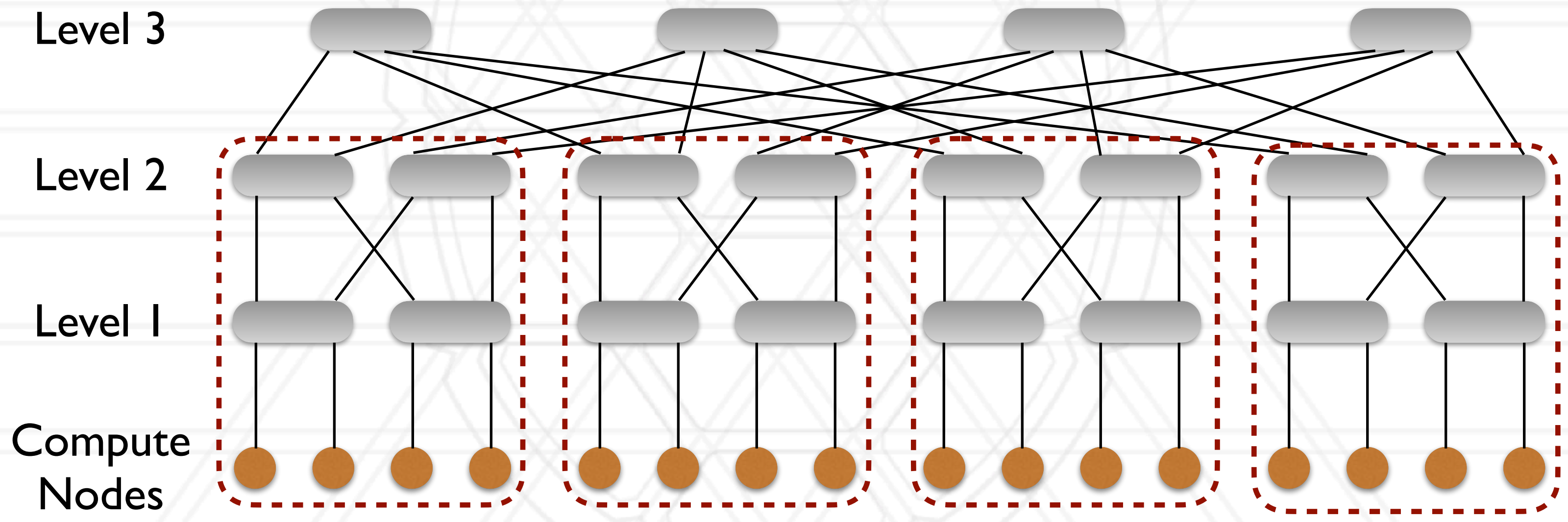
  - Low network diameter, high bandwidth



Router/switch radix = number of ports = k

Pod = group of switches = k/2 switches

# Fat-tree network

- Most popular network topology

  - Low network diameter, high bandwidth



Level 3

Level 2
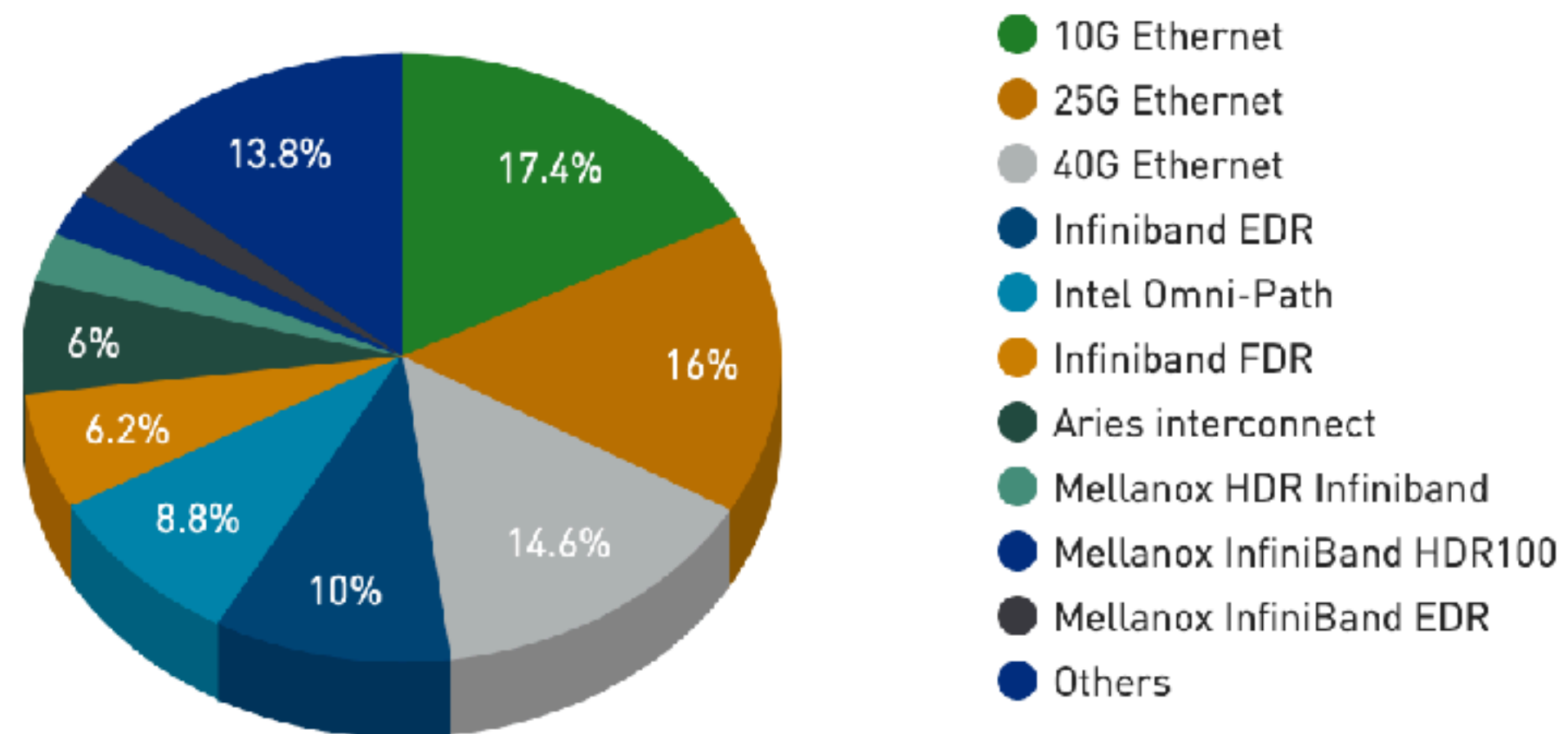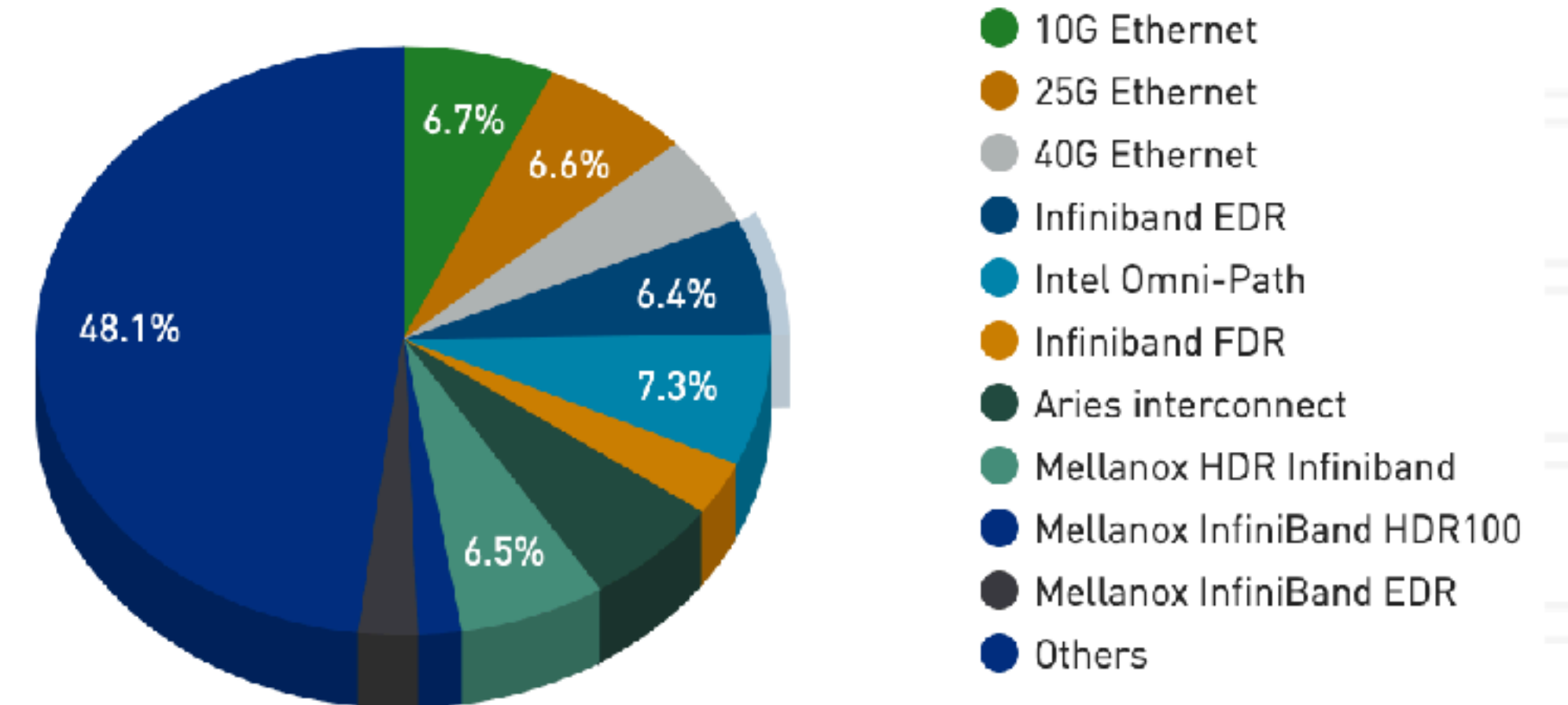
Level 1

Compute
Nodes

Router/switch radix = number of ports = k

Pod = group of switches = k/2 switches          Max. number of pods = k

# Fat-tree networks on the top500 list

- Infiniband EDR/FDR/HDR, Intel Omni-Path



**Interconnect System Share**

| | |
|---|---|
| ● | 10G Ethernet |
| ● | 25G Ethernet |
| ● | 40G Ethernet |
| ● | Infiniband EDR |
| ● | Intel Omni-Path |
| ● | Infiniband FDR |
| ● | Aries interconnect |
| ● | Mellanox HDR Infiniband |
| ● | Mellanox InfiniBand HDR100 |
| ● | Mellanox InfiniBand EDR |
| ● | Others |

17.4%  16%  14.6%  10%  8.8%  6.2%  6%  13.8%

**Interconnect Performance Share**

6.7%  6.6%  6.4%  7.3%  48.1%  6.5%

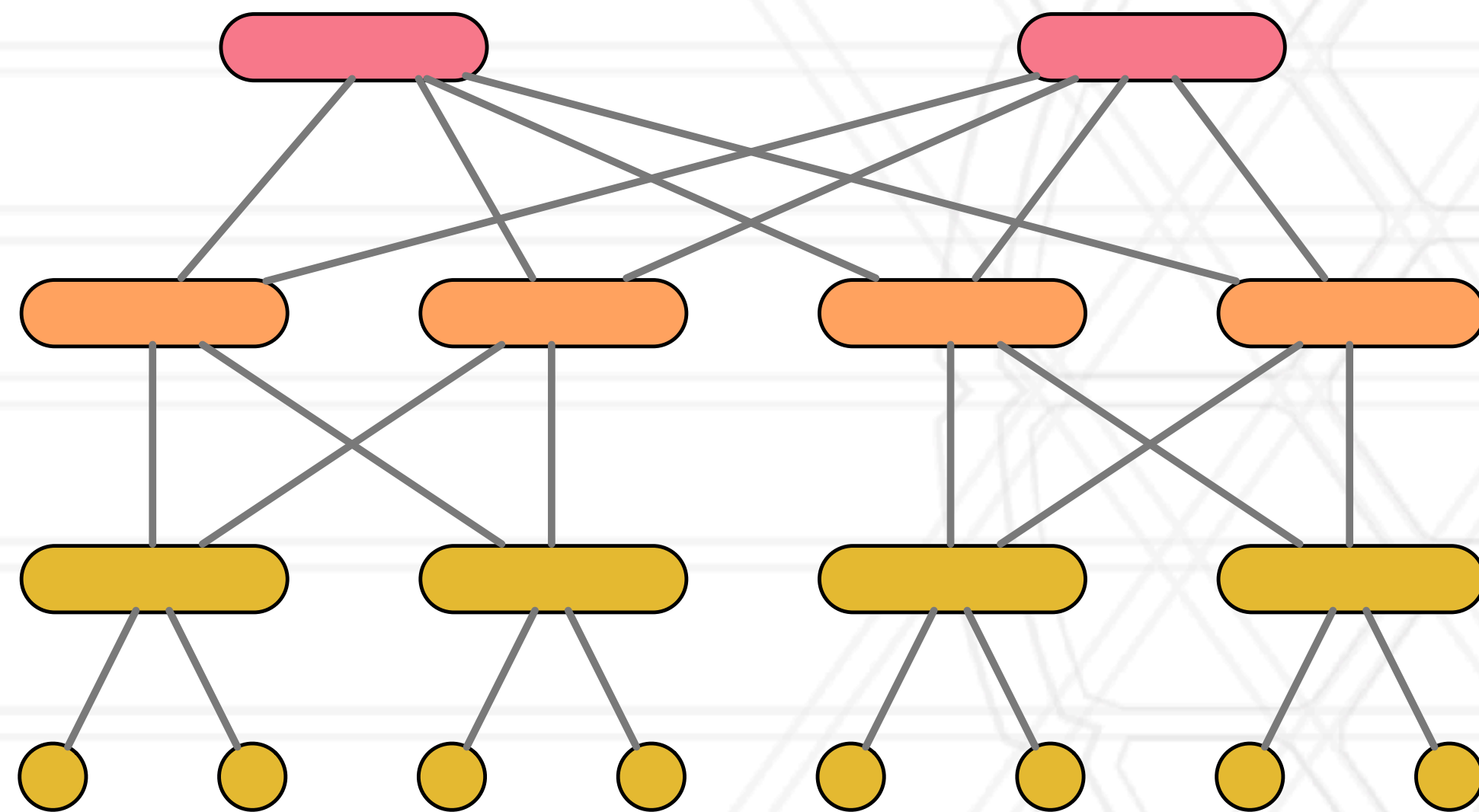https://www.top500.org/statistics/list, November 2020
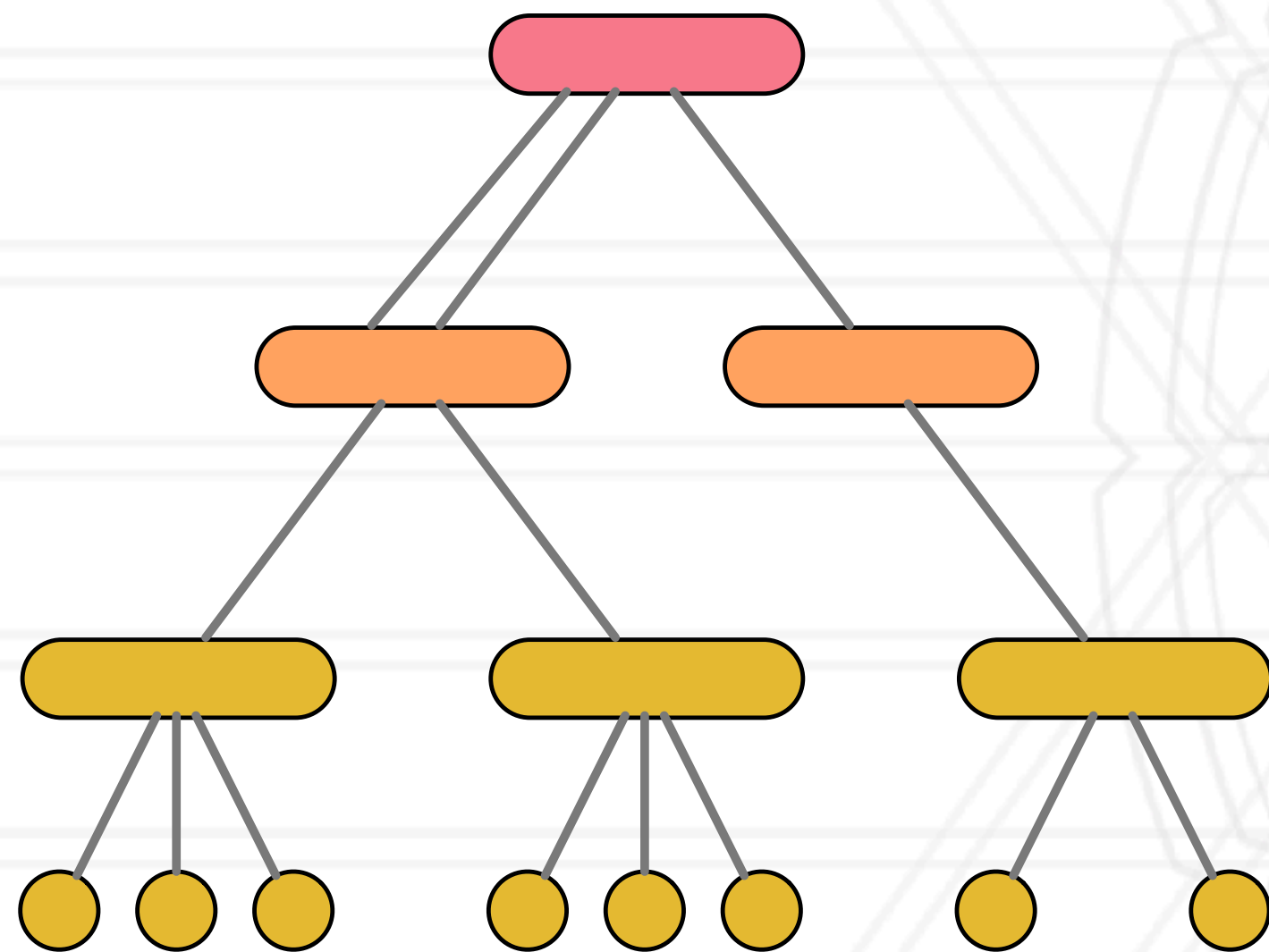
# Routing on a fat-tree

- Until recently, most fat-tree installations used static routing

  - Destination-mod-k (D-mod-k) routing

- Adaptive routing is now starting to be used

DEPARTMENT OF
COMPUTER SCIENCE

# Variations on a full bandwidth fat-tree



Single-rail single-plane fat-tree

DEPARTMENT OF
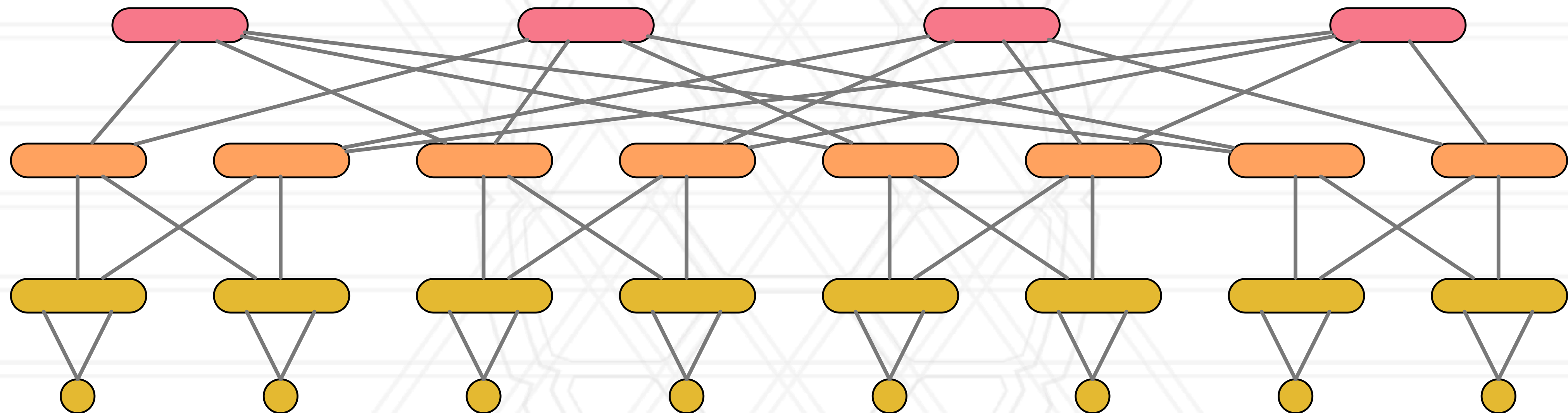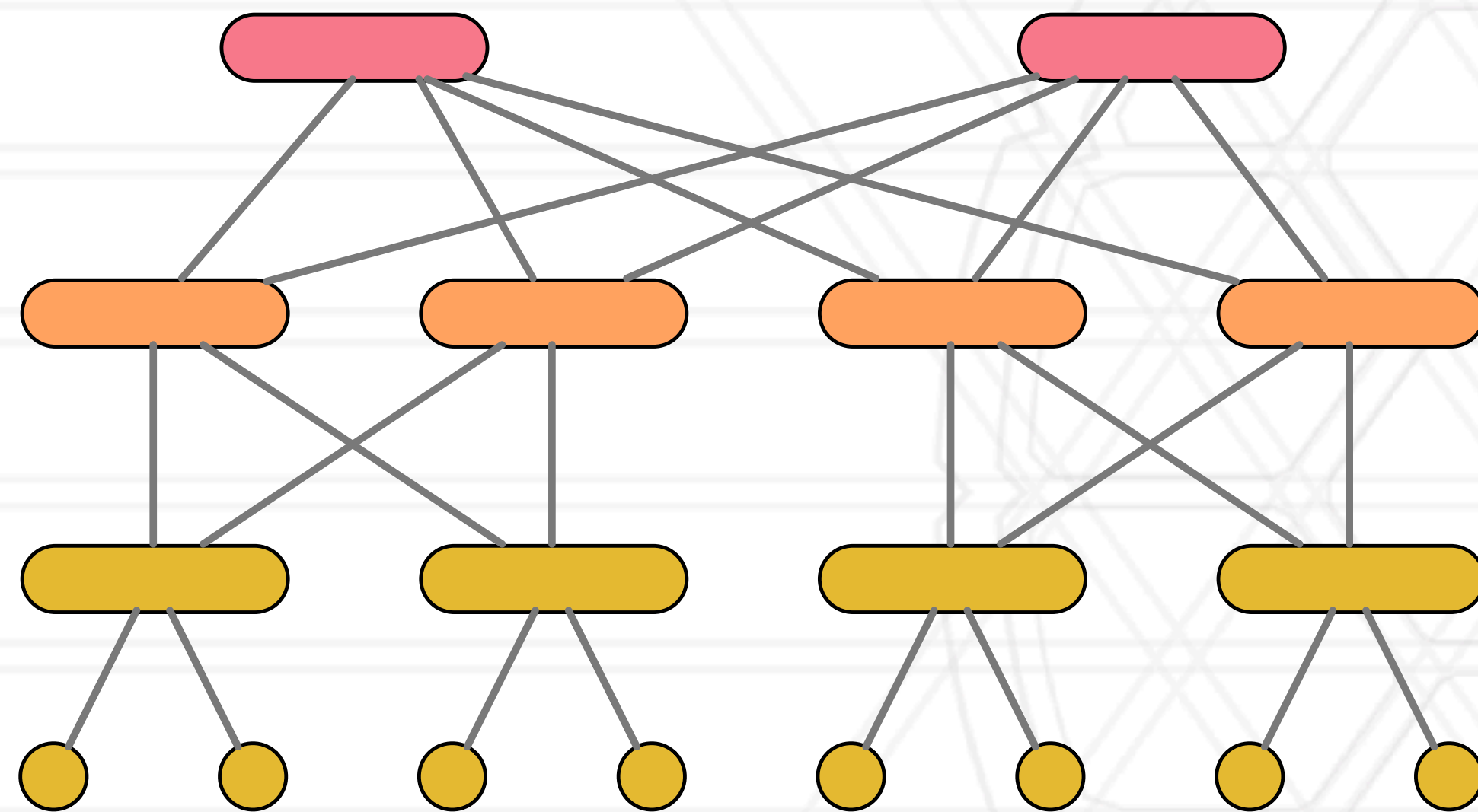COMPUTER SCIENCE

# Variations on a full bandwidth fat-tree



Single-rail single-plane fat-tree (tapered)

# Variations on a full bandwidth fat-tree



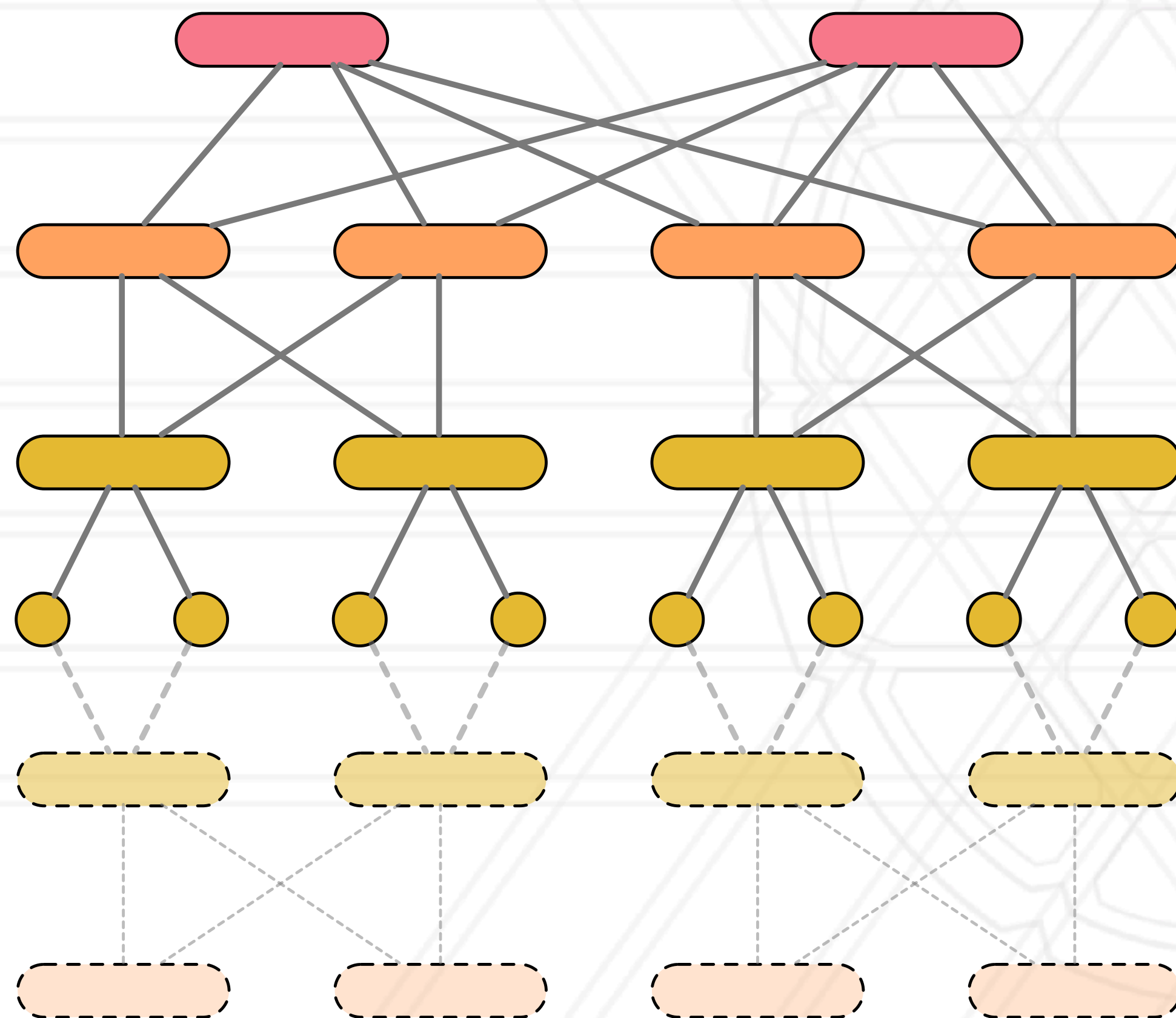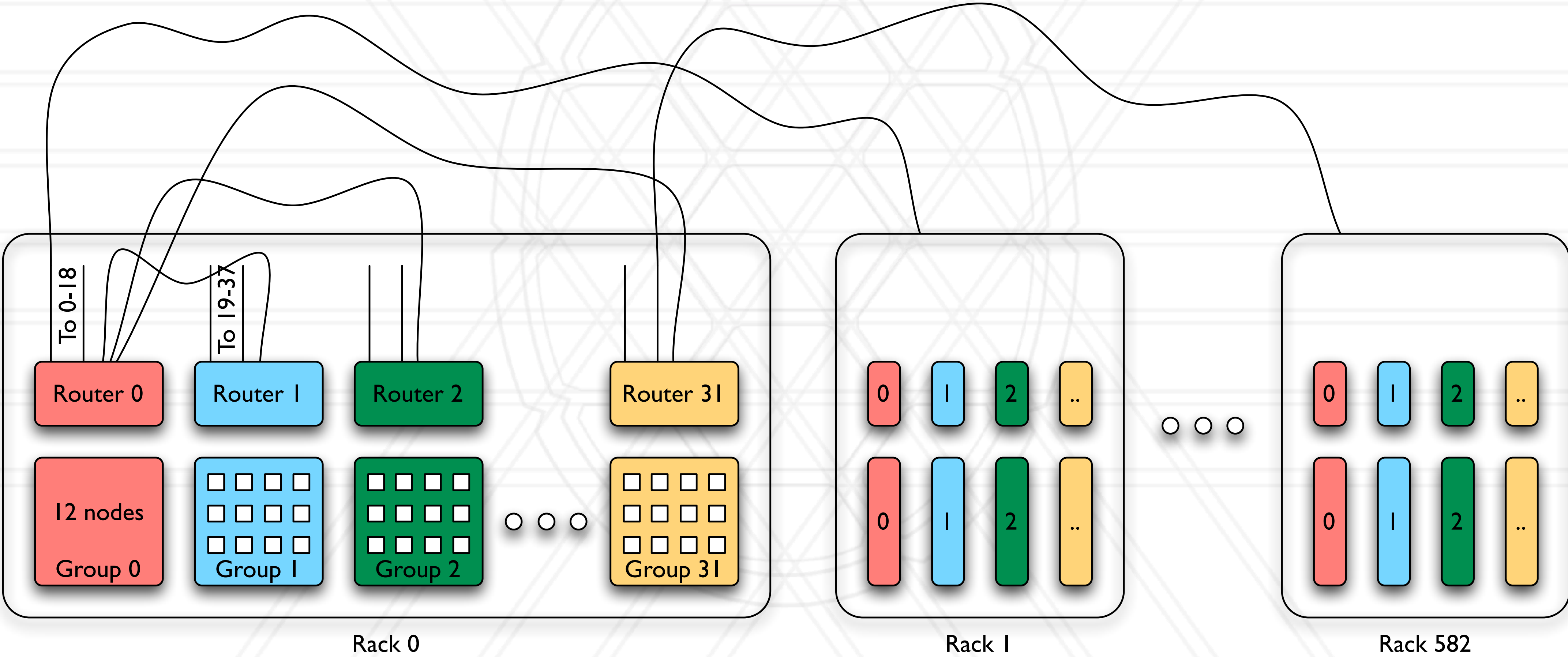Dual-rail single-plane fat-tree

DEPARTMENT OF
COMPUTER SCIENCE

# Variations on a full bandwidth fat-tree



Single-rail single-plane fat-tree

DEPARTMENT OF
COMPUTER SCIENCE

# Variations on a full bandwidth fat-tree
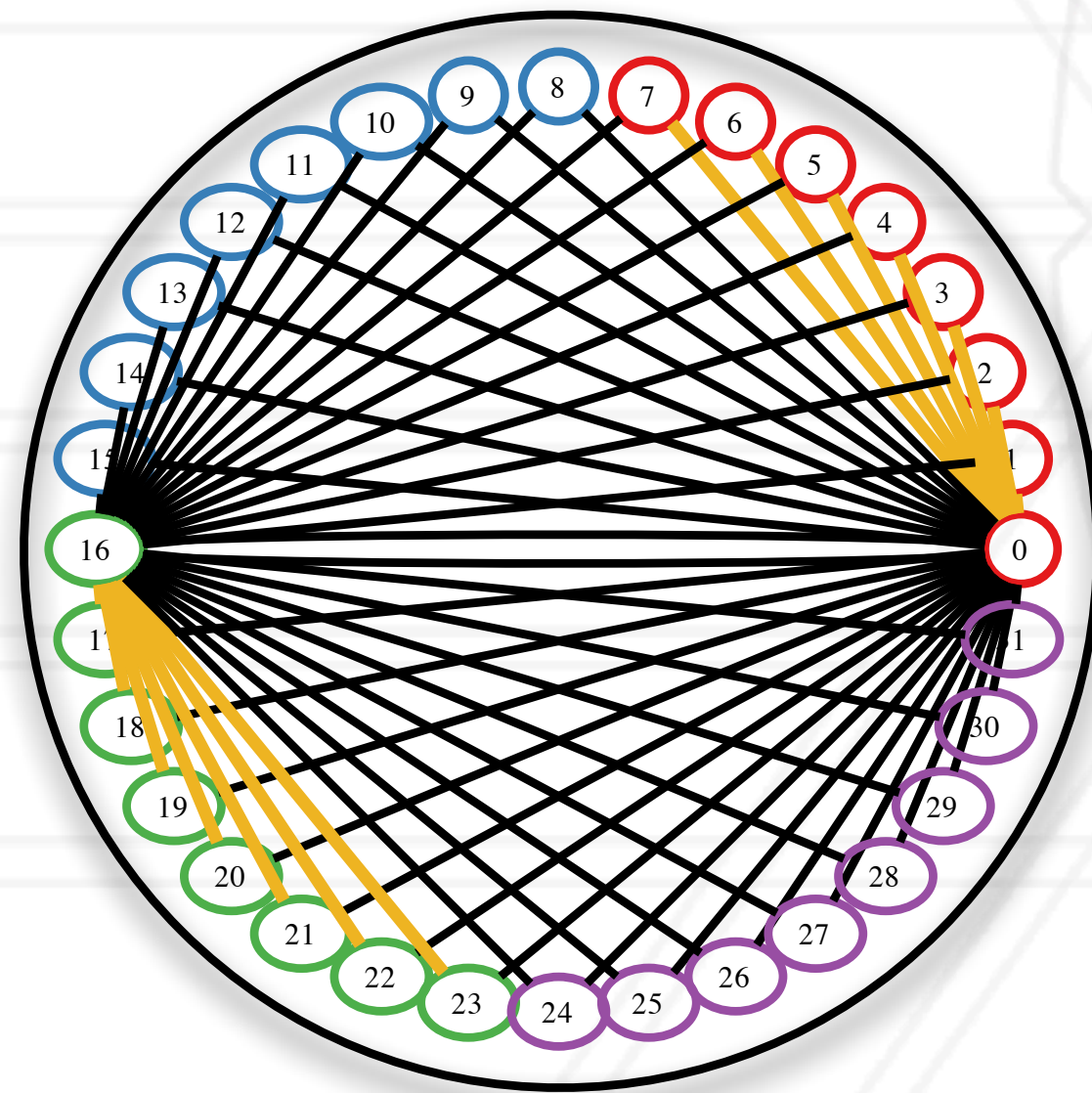


Dual-rail dual-plane fat-tree

DEPARTMENT OF
COMPUTER SCIENCE

# Dragonfly network

# IBM PERCS network

- All-to-all connections within each group



One supernode in the PERCS topology

DEPARTMENT OF
COMPUTER SCIENCE
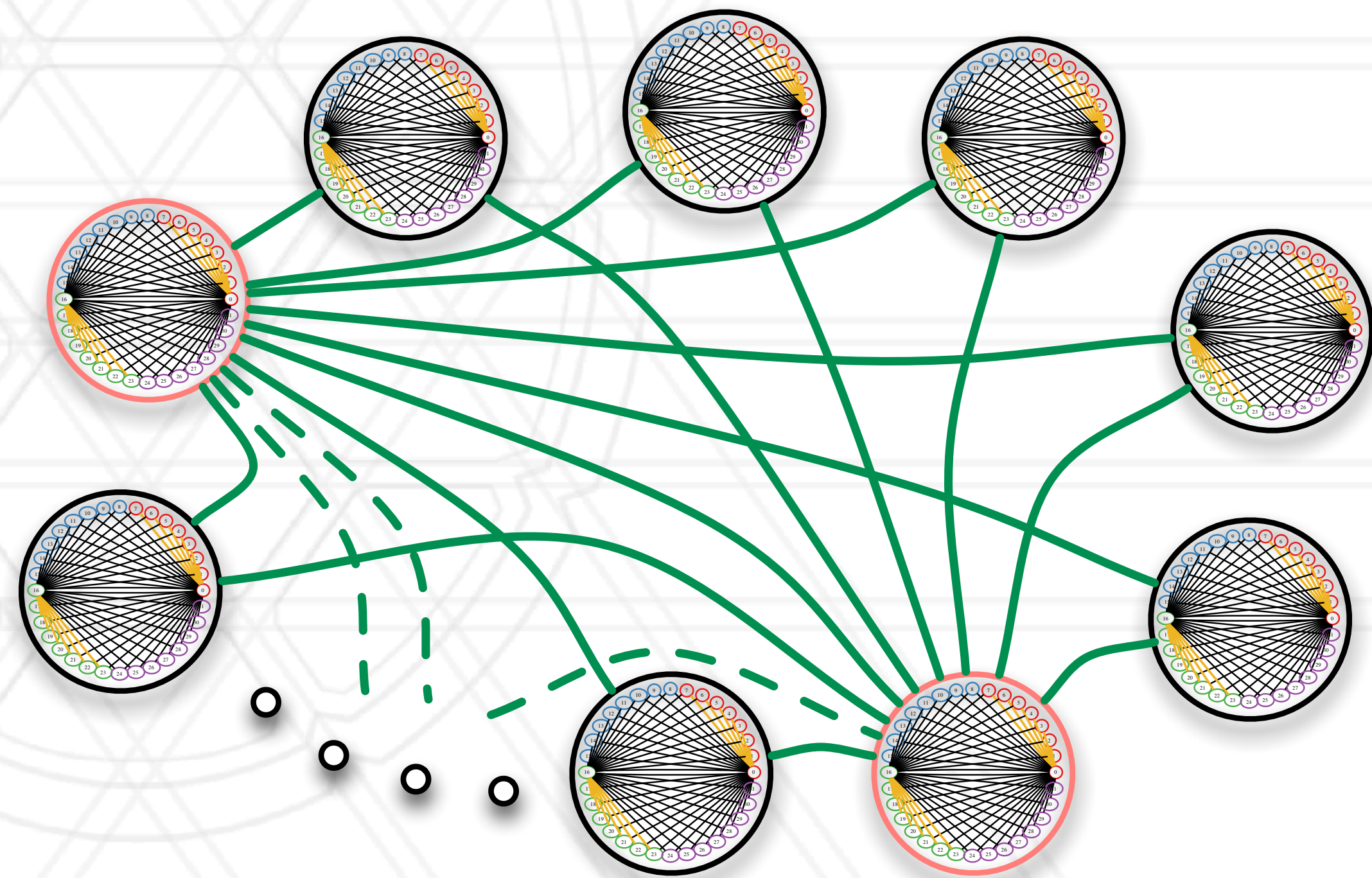
# IBM PERCS network

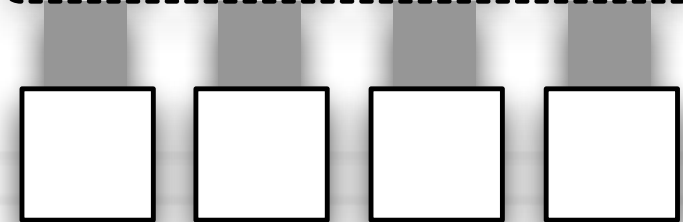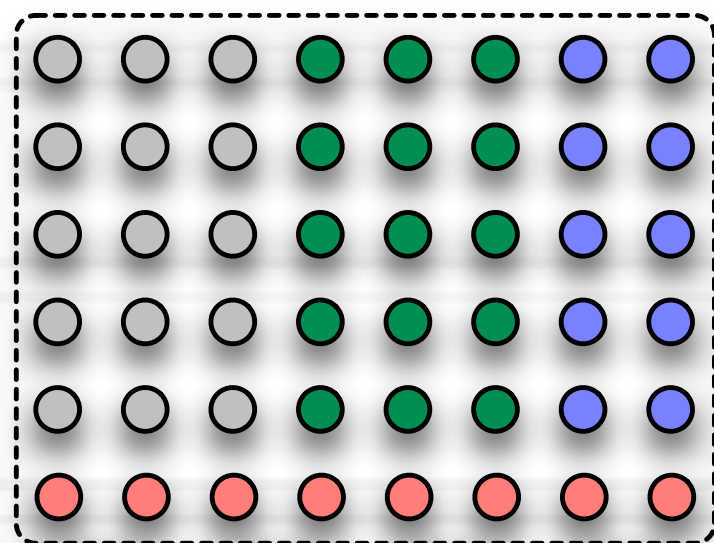- All-to-all connections within each group



One supernode in the PERCS topology

# Cray Aries network

- Row and column all-to-all connections within each group



Aries Router

Compute Nodes

DEPARTMENT OF
COMPUTER SCIENCE

# Cray Aries network

- Row and column all-to-all connections within each group



Aries Router

Compute Nodes

A group with 96 Aries routers

Column all-to-all (black) links          Row all-to-all (green) links

DEPARTMENT OF
COMPUTER SCIENCE

# Cray Aries network

- Row and column all-to-all connections within each group

Aries Router



Compute Nodes

A group with 96 Aries routers

Column all-to-all (black) links          Row all-to-all (green) links

Two-level dragonfly with multiple groups

Inter-group (blue) links
(not all links are shown)

DEPARTMENT OF
COMPUTER SCIENCE

# Network comparisons

| Network topology | #nodes/router | #links/router | Maximum system size (#nodes) |
|---|---|---|---|
| All-to-all (A2A) dragonfly | k/4 | k/2 (**L**), k/4 (**G**) | $(k/2 + 1)^2 \times (k/4 + 1) \times k/4$ |
| Row-column (RC) dragonfly | k/6 | 2k/3 (**L**), k/6 (**G**) | $(k/6 + 1)^5 \times (k/6 + 1) \times k/6$ |
| Express mesh (3D, gap=1) | k/4 | 3k/4 | $(k/4 + 1)^3 \times k/4$ |
| Fat-tree (three-level) | k/2 | k/2 | $k/2 \times k/2 \times k$ |

SIGSIM-PADS '19, June 3-5 2019, Chicago, IL, USA

Parallel Simulation



Abhinav Bhatele (CMSC714)

# Questions?



Abhinav Bhatele

5218 Brendan Iribe Center (IRB) / College Park, MD 20742

phone: 301.405.4507 / e-mail: bhatele@cs.umd.edu