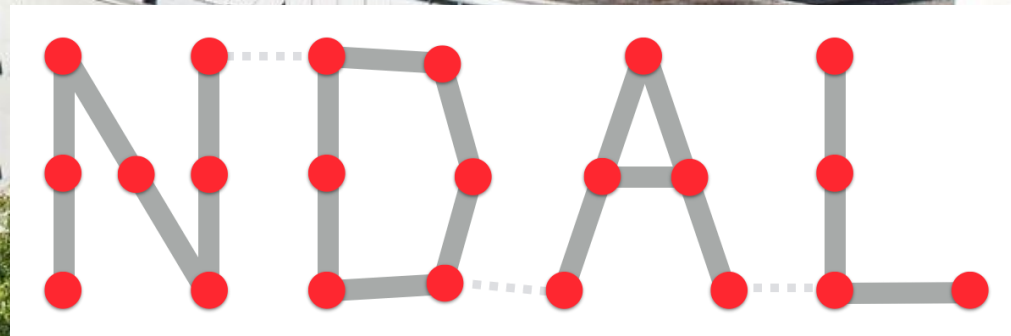
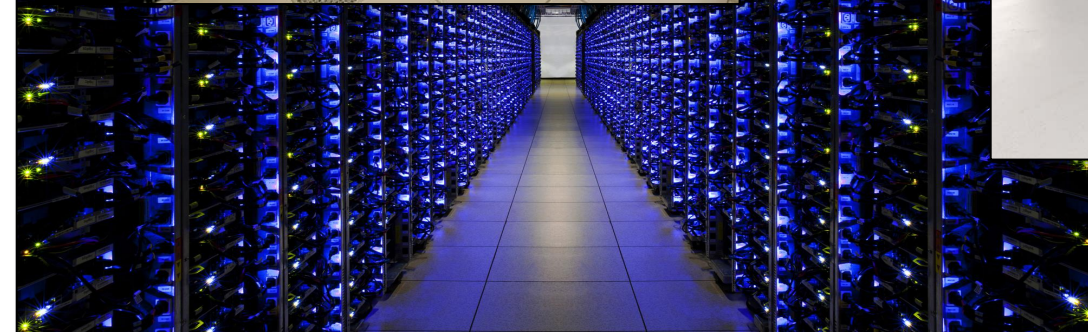
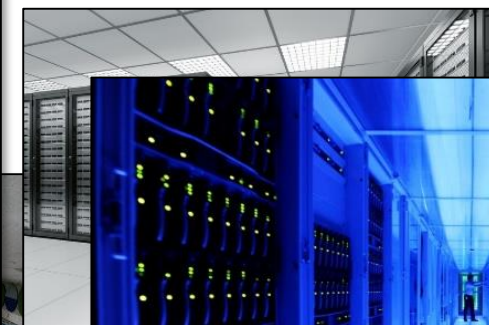


# FatPaths: Routing in Supercomputers and Data Centers when Shortest Paths Fall Short

MACIEJ BESTA, MARCEL SCHNEIDER, MAREK KONIECZNY, KAROLINA CYNK  
ERIK HENRIKSSON, SALVATORE DI GIROLAMO, ANKIT SINGLA, TORSTEN HOEFLER

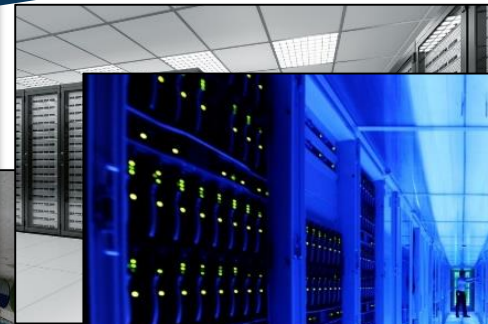




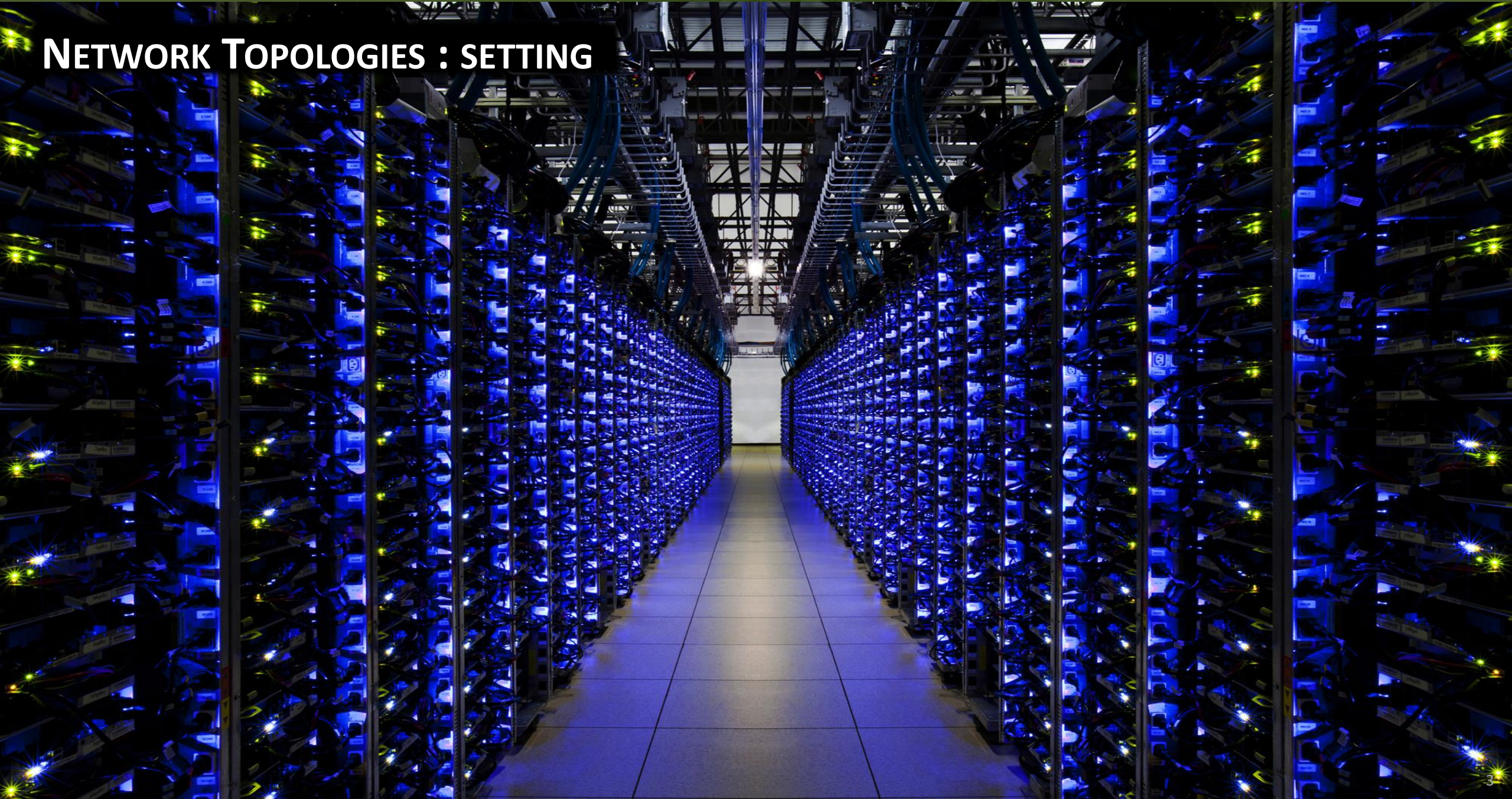




**The interconnect:** a key part of supercomputers and data centers, relevant both for high performance and low cost



# NETWORK TOPOLOGIES : SETTING



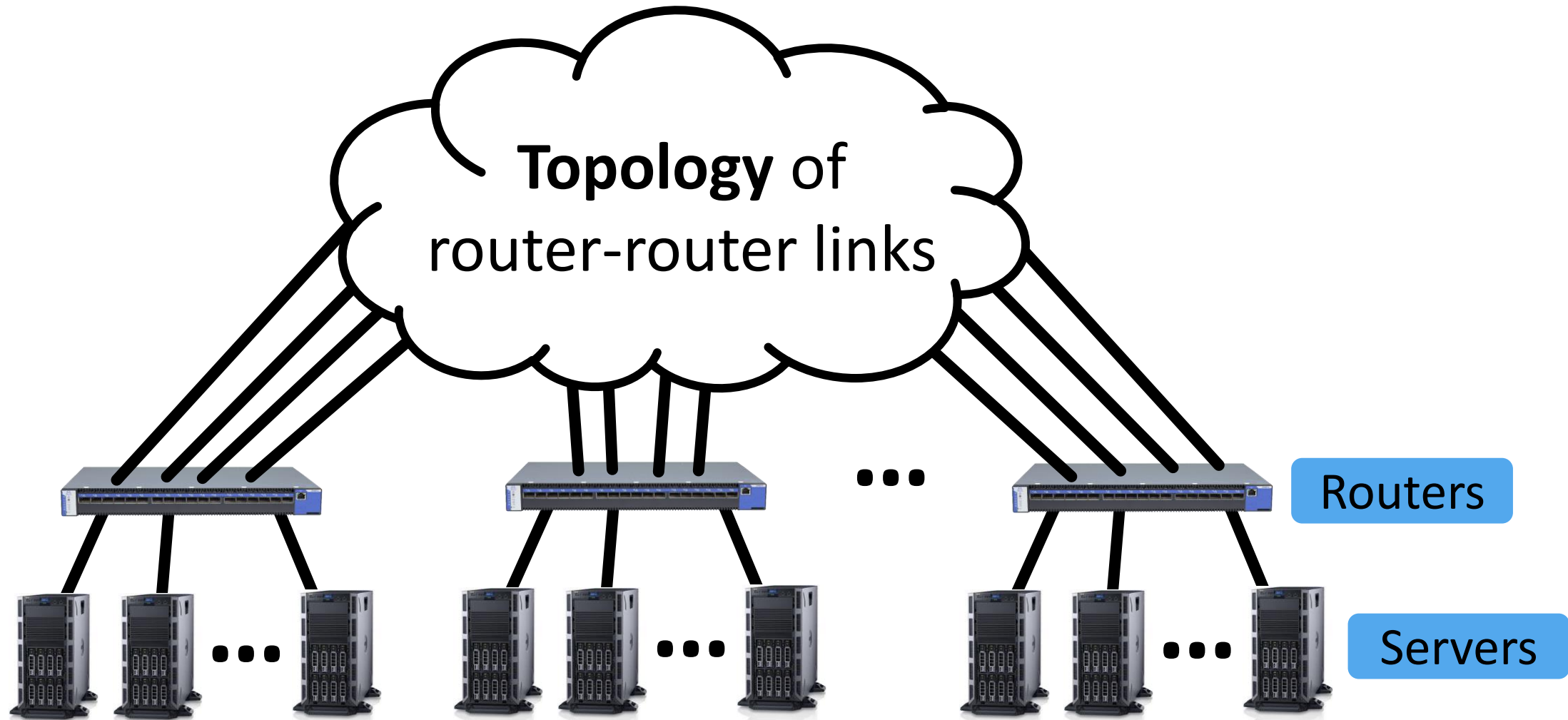
# NETWORK TOPOLOGIES : SETTING



# NETWORK TOPOLOGIES : SETTING



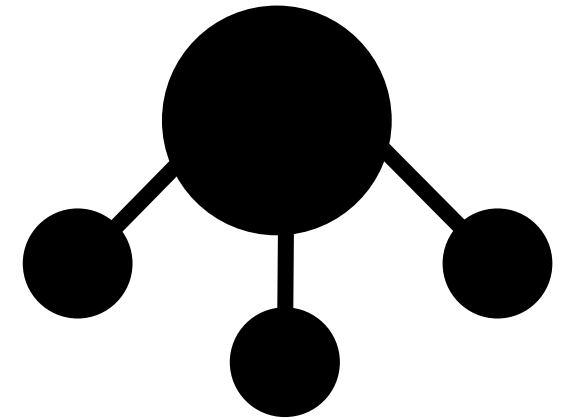
# NETWORK TOPOLOGIES : SETTING





# NETWORK TOPOLOGIES : SETTING

We will often use the following symbols for routers and servers



Topology of  
router-router links

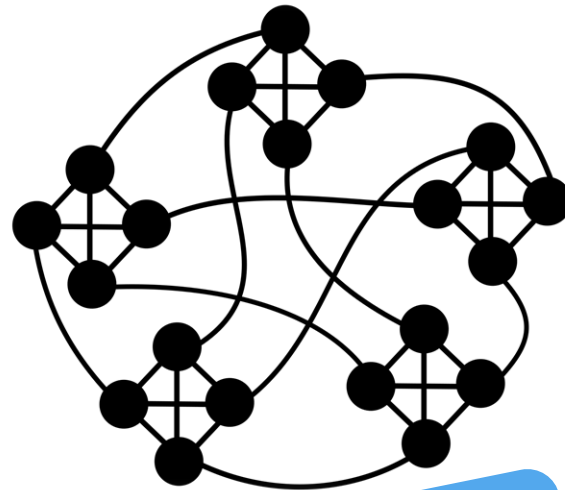


Routers

Servers

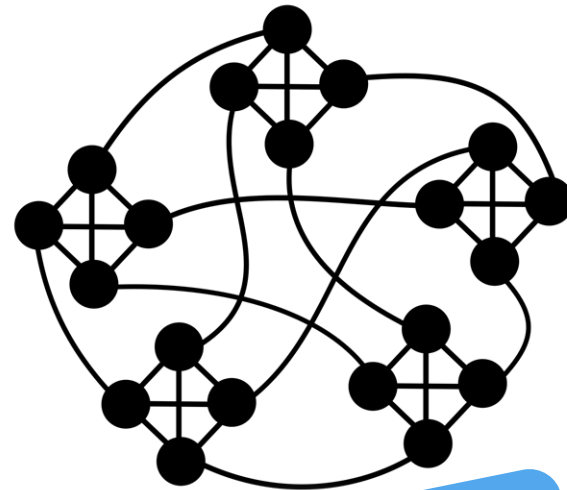
# LOW-DIAMETER NETWORK TOPOLOGIES

# LOW-DIAMETER NETWORK TOPOLOGIES



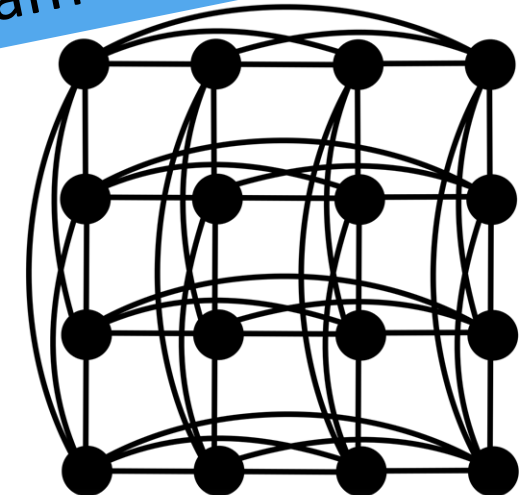
Dragonfly [2008]  
diameter = 3

# LOW-DIAMETER NETWORK TOPOLOGIES



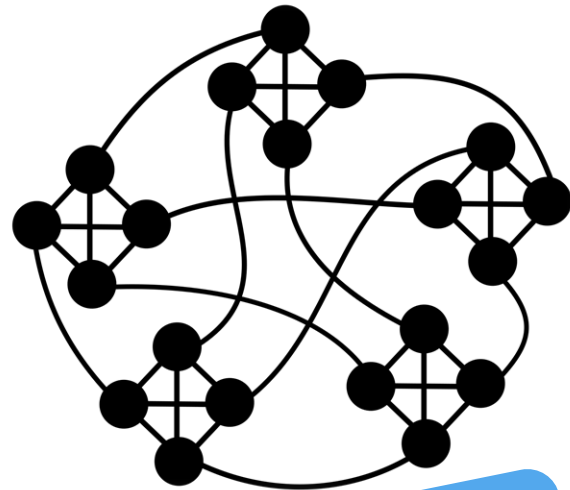
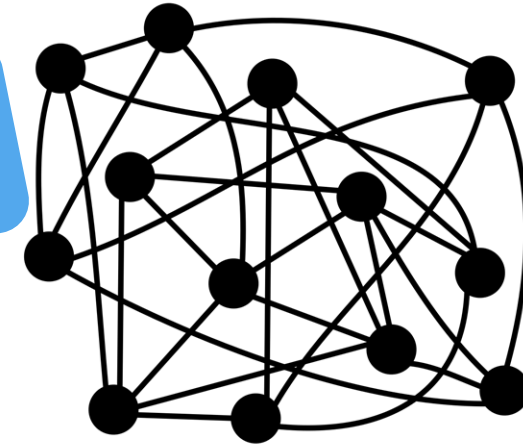
Dragonfly [2008]  
diameter = 3

HyperX [2009]  
diameter = 3



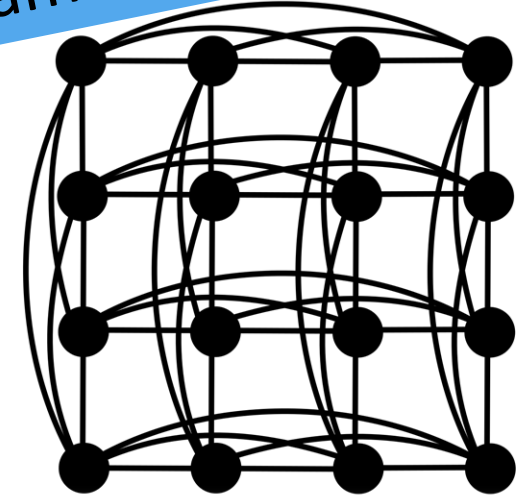
# LOW-DIAMETER NETWORK TOPOLOGIES

Jellyfish [2012]  
diameter  $\leq 3$



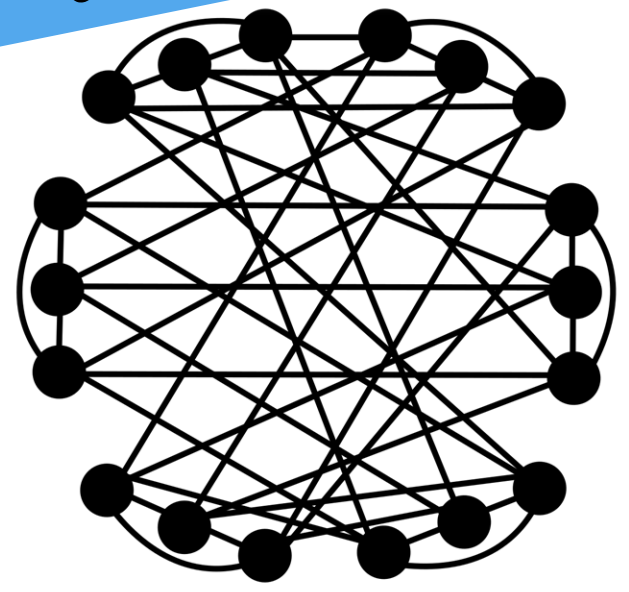
Dragonfly [2008]  
diameter = 3

HyperX [2009]  
diameter = 3

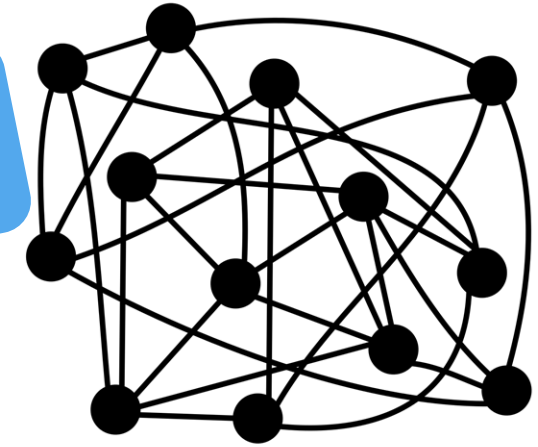


# LOW-DIAMETER NETWORK TOPOLOGIES

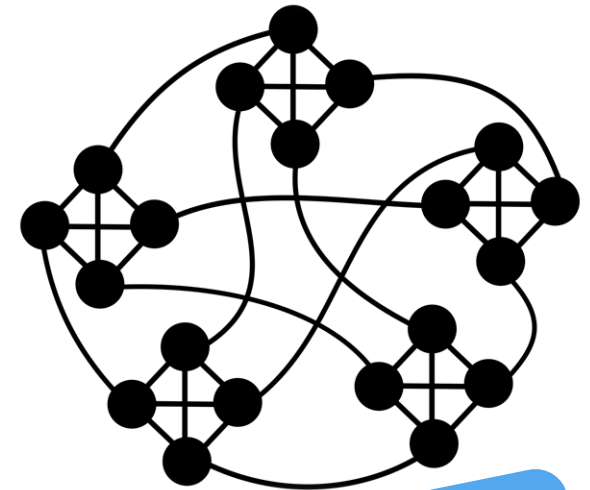
Slim Fly [2014]  
diameter = 2



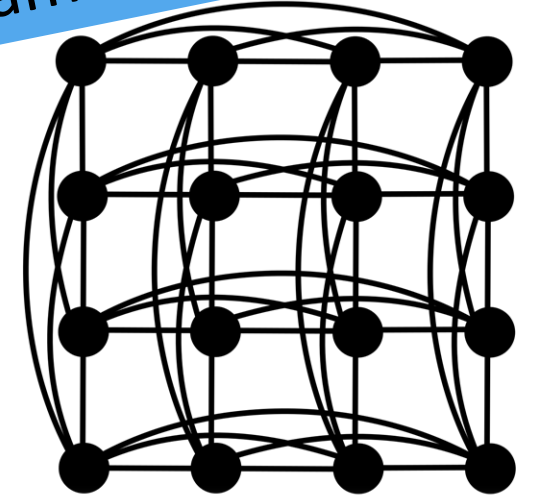
Jellyfish [2012]  
diameter  $\leq 3$



Dragonfly [2008]  
diameter = 3



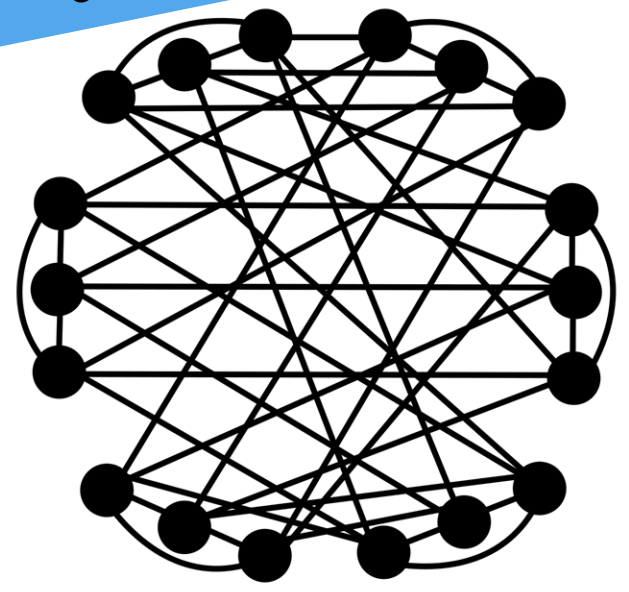
HyperX [2009]  
diameter = 3



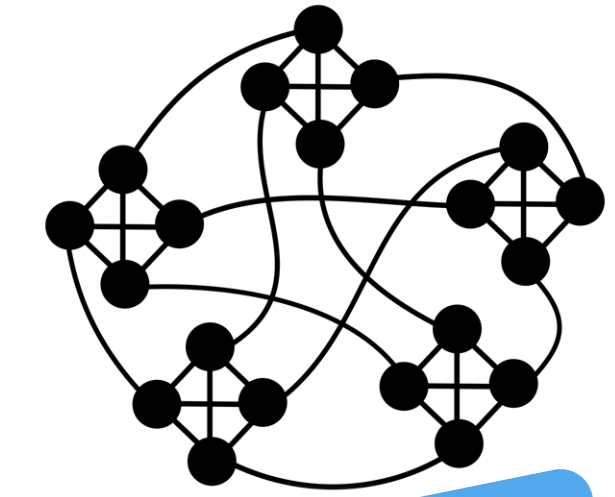
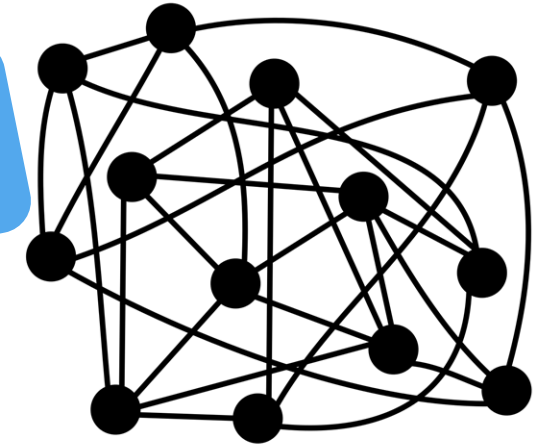
# LOW-DIAMETER NETWORK TOPOLOGIES

Xpander [2016]  
diameter  $\leq 3$

Slim Fly [2014]  
diameter = 2

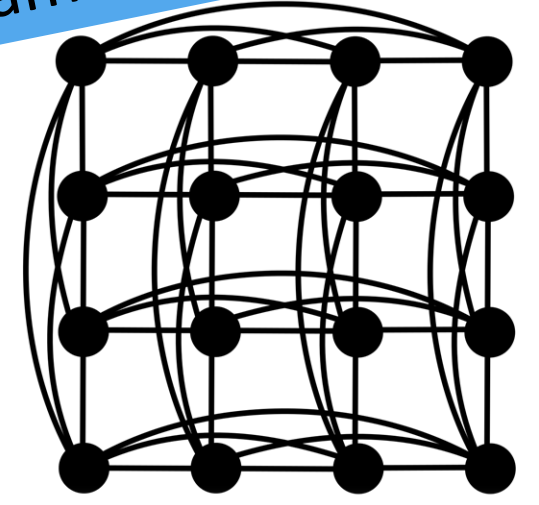


Jellyfish [2012]  
diameter  $\leq 3$



Dragonfly [2008]  
diameter = 3

HyperX [2009]  
diameter = 3

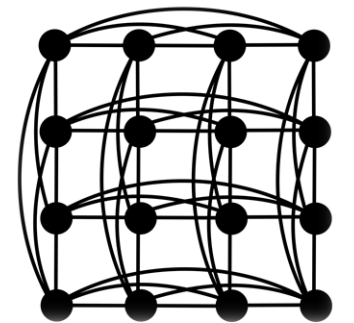
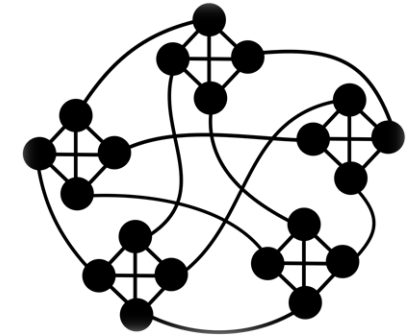
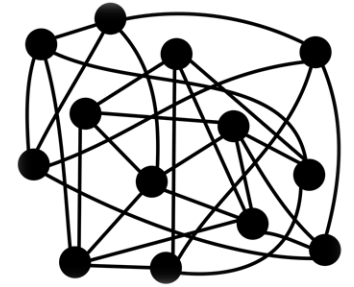
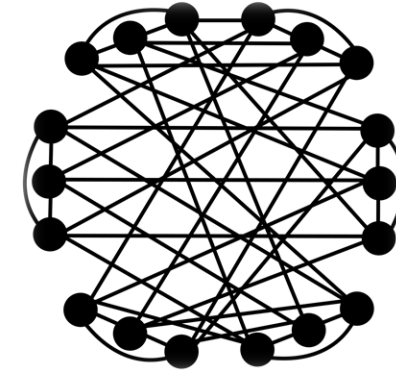


# LOW-DIAMETER NETWORK TOPOLOGIES



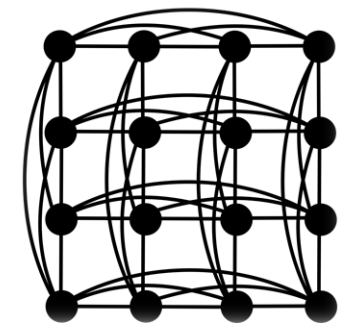
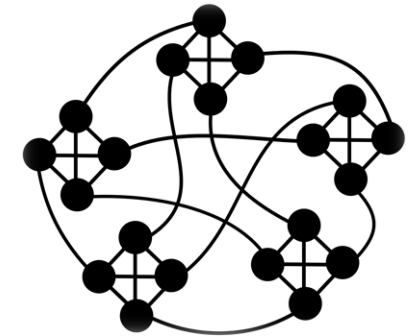
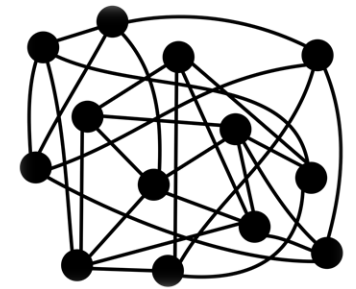
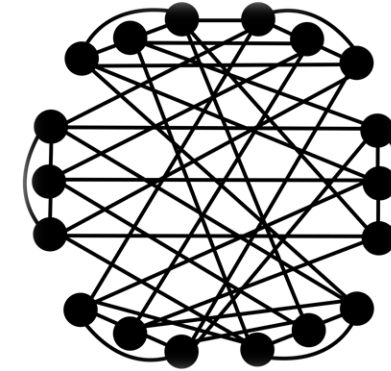
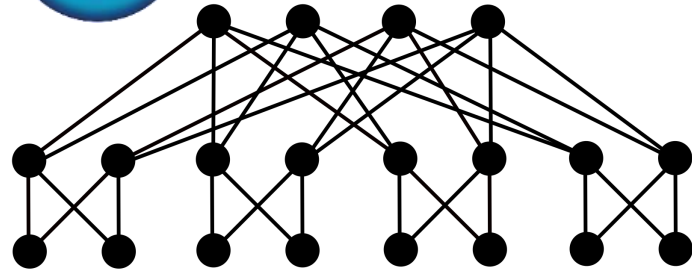
# LOW-DIAMETER NETWORK TOPOLOGIES

Why care? (there is already Clos, Fat tree, etc.).



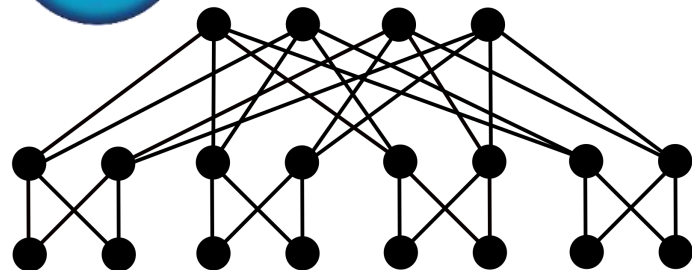
# LOW-DIAMETER NETWORK TOPOLOGIES

Why care? (there is already Clos, Fat tree, etc.).

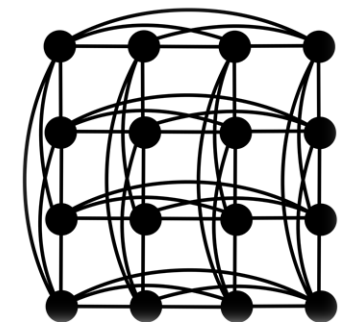
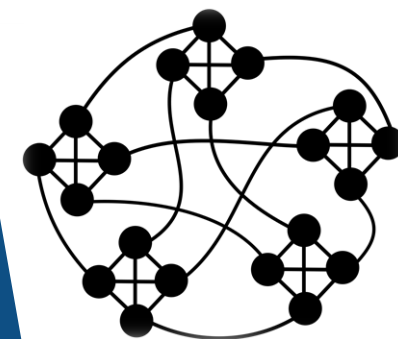
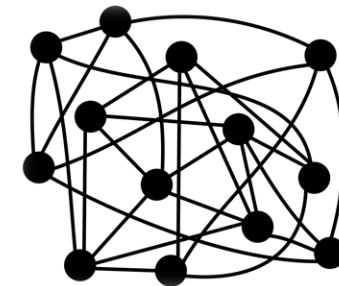
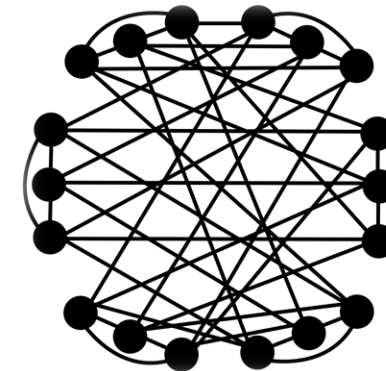


# LOW-DIAMETER NETWORK TOPOLOGIES

Why care? (there is already Clos, Fat tree, etc.).

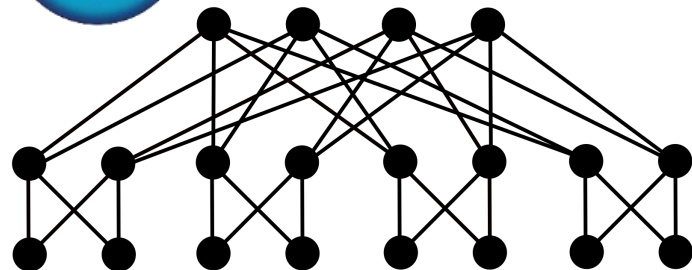


Cost-effective (Slim fly is 50% cheaper than Fat tree for a fixed network size counted in #endpoints)



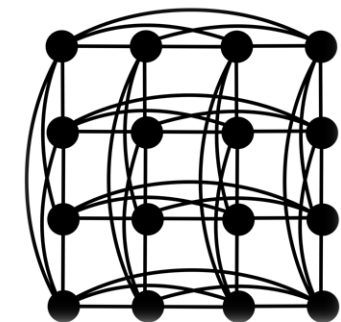
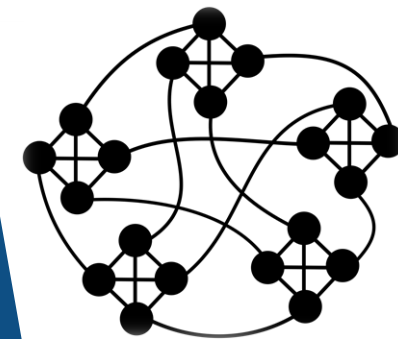
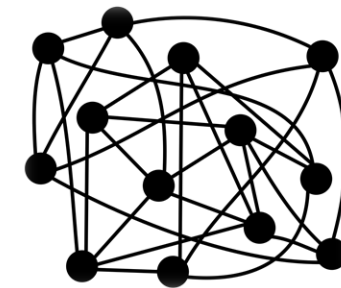
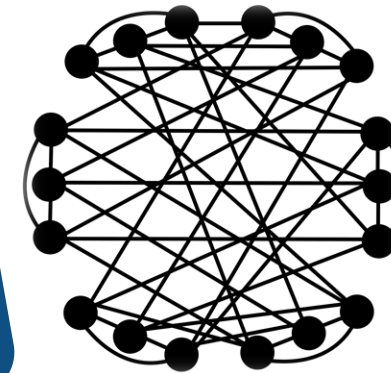
# LOW-DIAMETER NETWORK TOPOLOGIES

Why care? (there is already Clos, Fat tree, etc.).



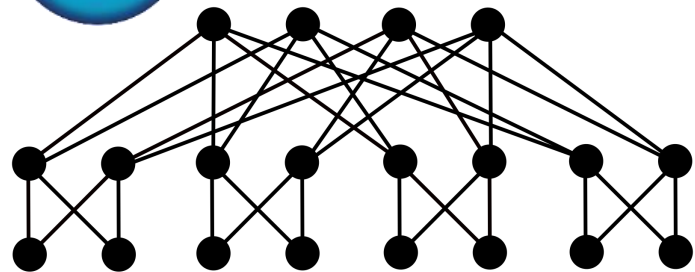
Energy-efficient  
(fewer routers and thus buffers)

Cost-effective (Slim fly is 50% cheaper than Fat tree for a fixed network size counted in #endpoints)



# LOW-DIAMETER NETWORK TOPOLOGIES

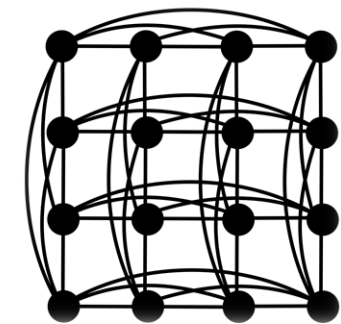
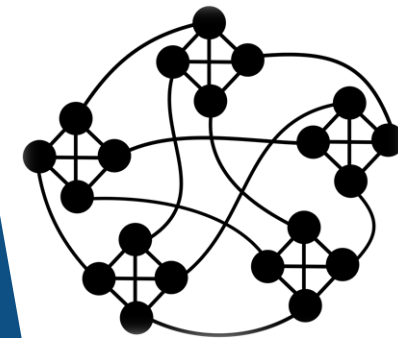
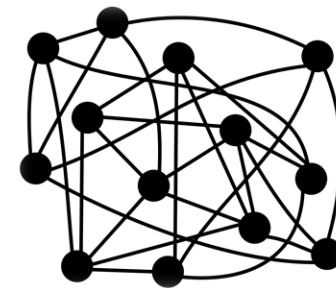
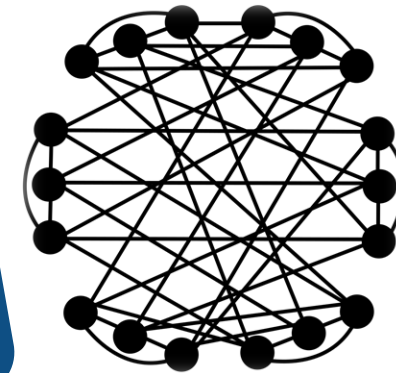
Why care? (there is already Clos, Fat tree, etc.).



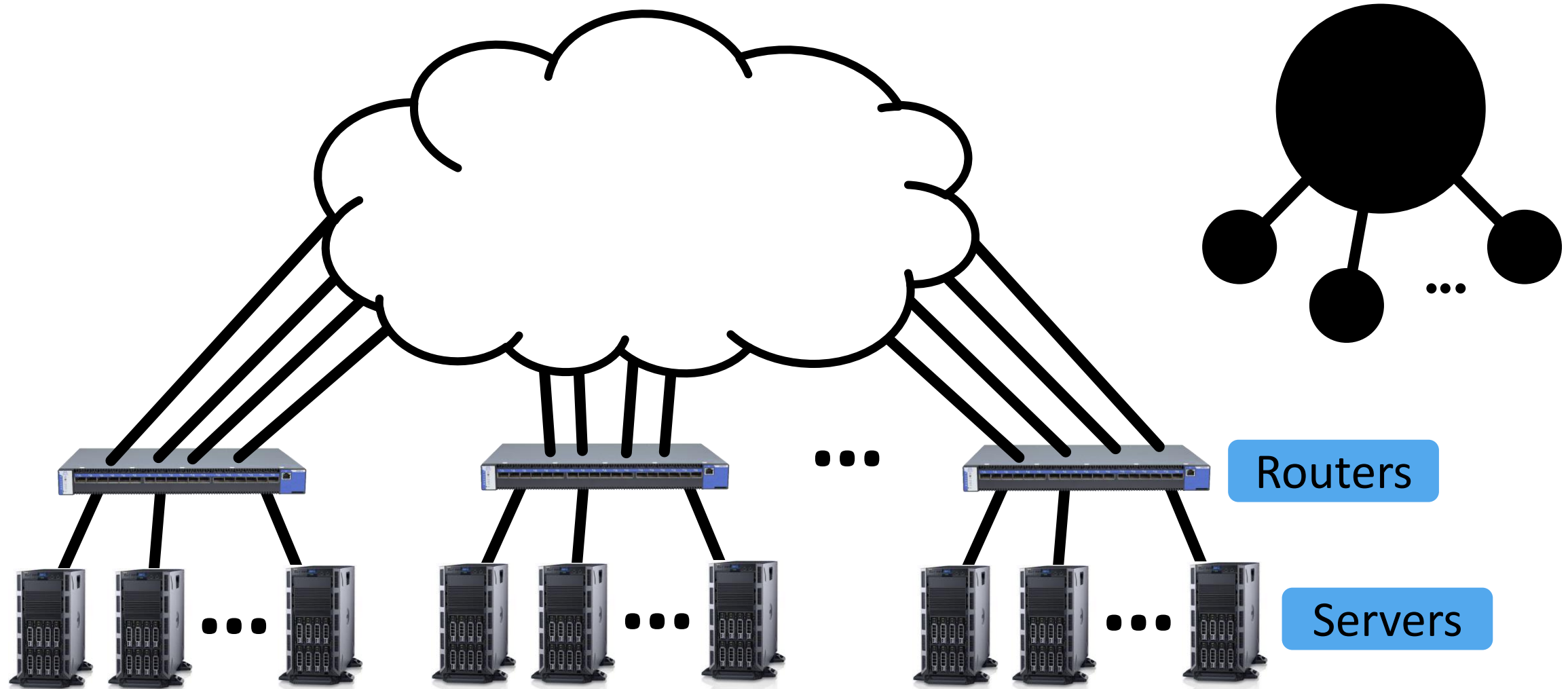
Energy-efficient  
(fewer routers and thus buffers)

Cost-effective (Slim fly is 50% cheaper than Fat tree for a fixed network size counted in #endpoints)

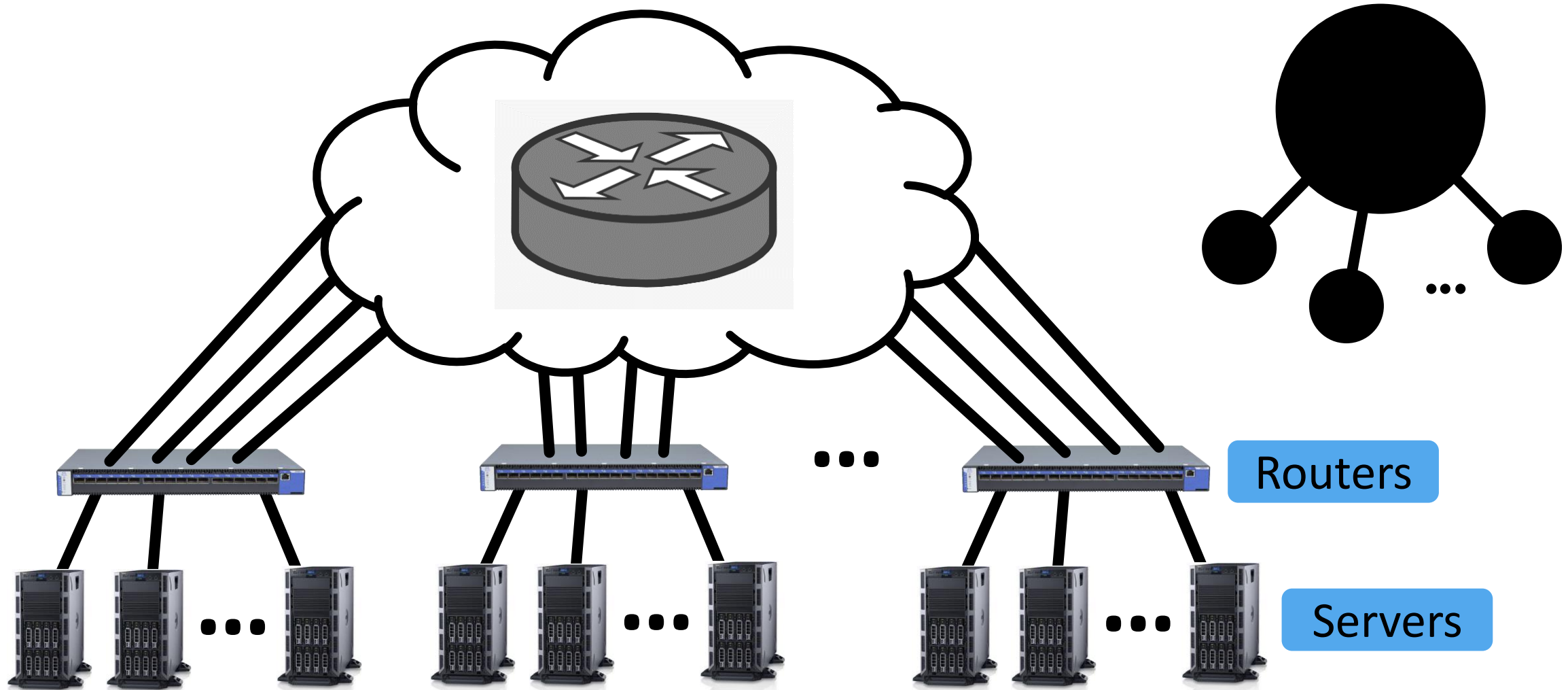
Low latency (from low diameters)



# NETWORK TOPOLOGIES : ROUTING



# NETWORK TOPOLOGIES : ROUTING



# ROUTING IN FAT TREES





## ROUTING IN FAT TREES

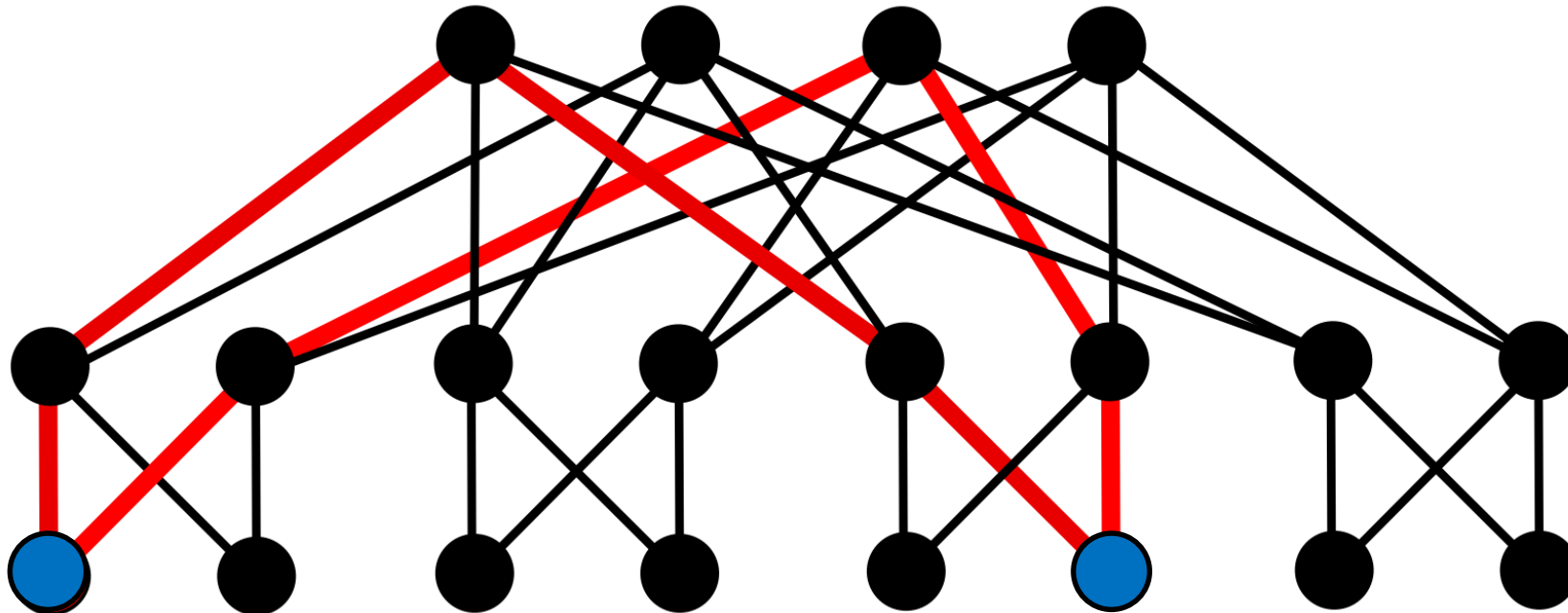


High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints

## ROUTING IN FAT TREES



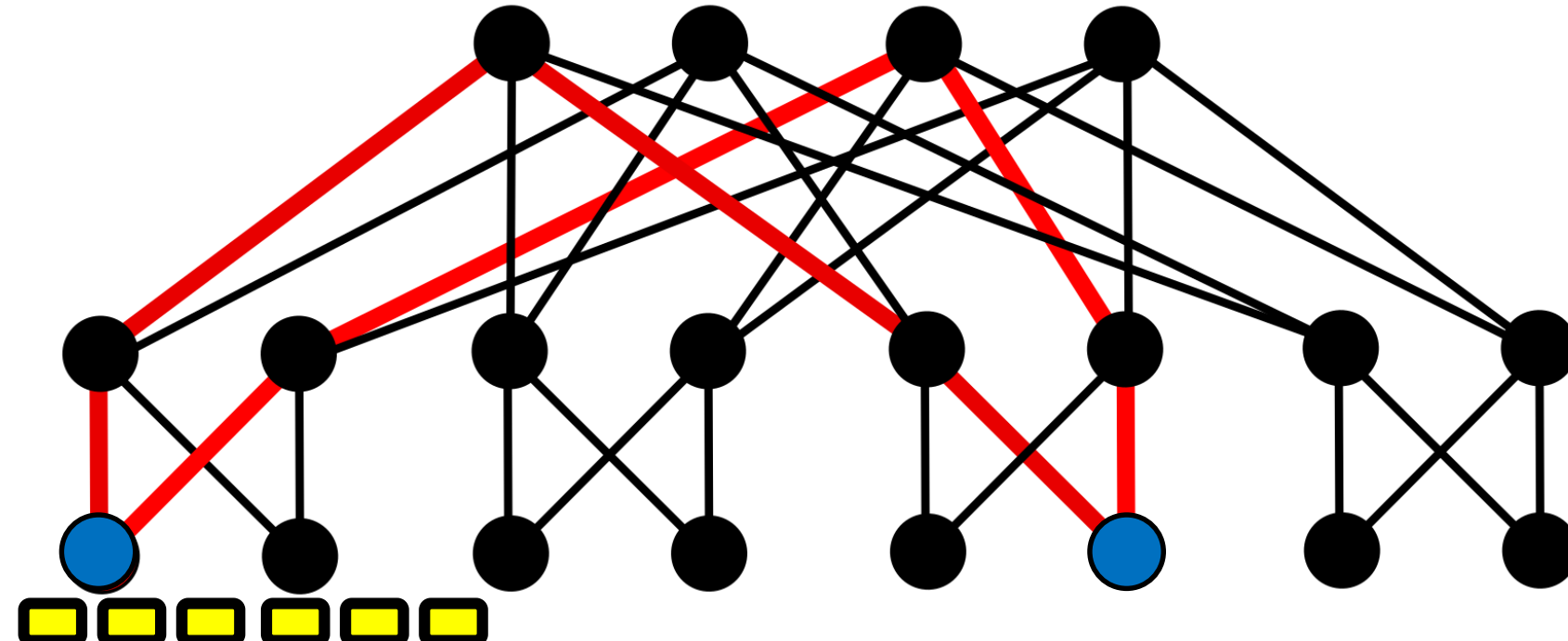
! High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints



# ROUTING IN FAT TREES



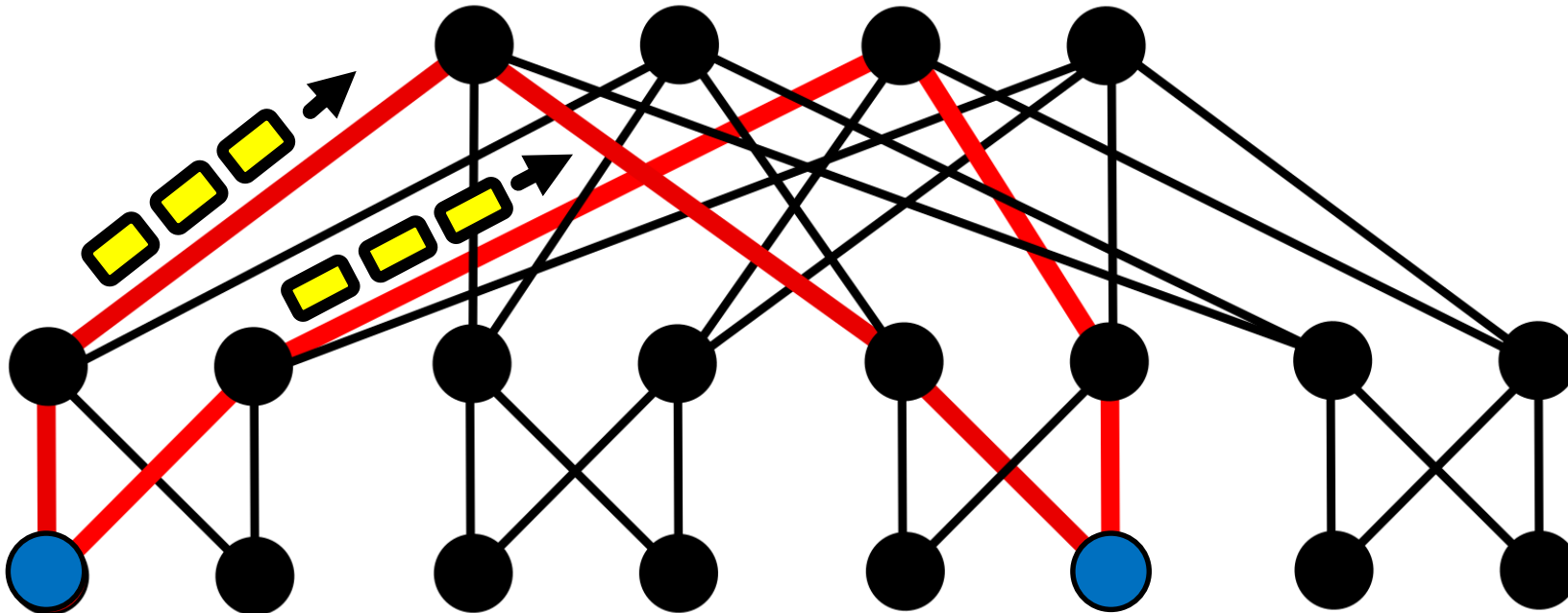
! High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints



## ROUTING IN FAT TREES



High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints

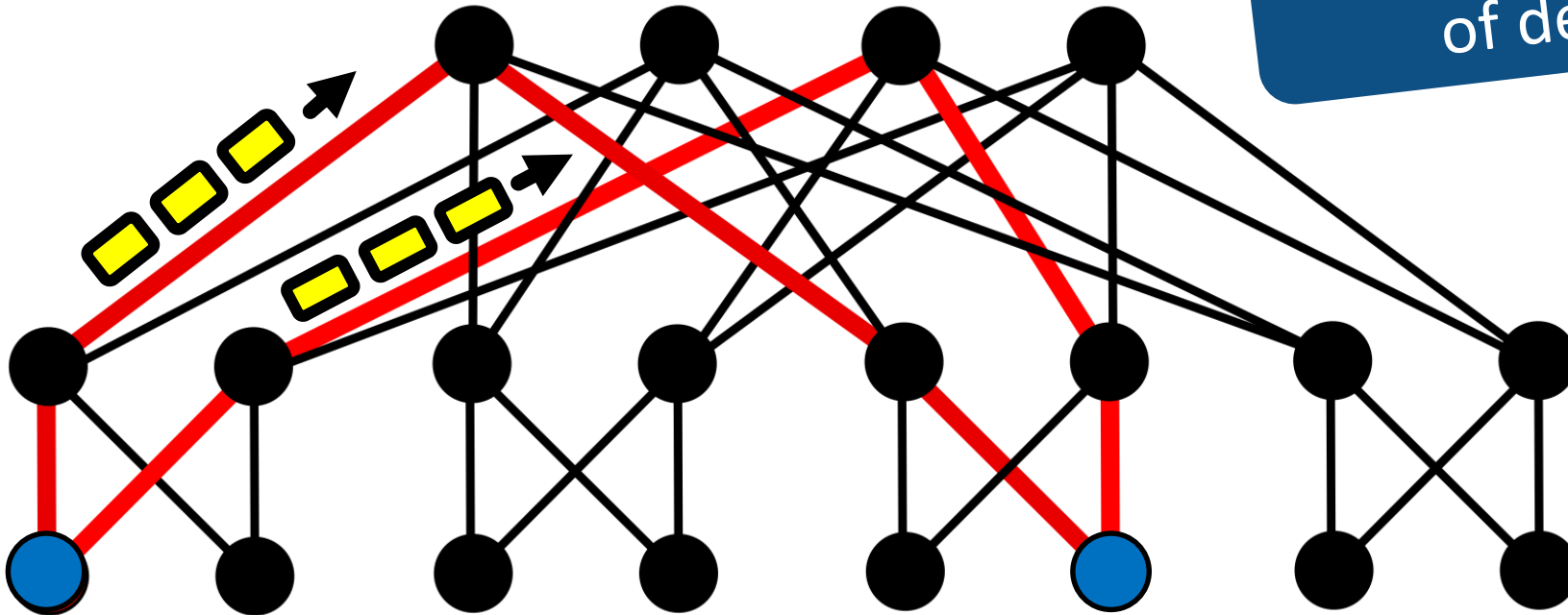


## ROUTING IN FAT TREES



High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints

Established techniques for using multipathing & plethora of designs available



# ROUTING IN FAT TREES



! High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints

Network Working Group  
Request for Comments: 2992  
Category: Informational

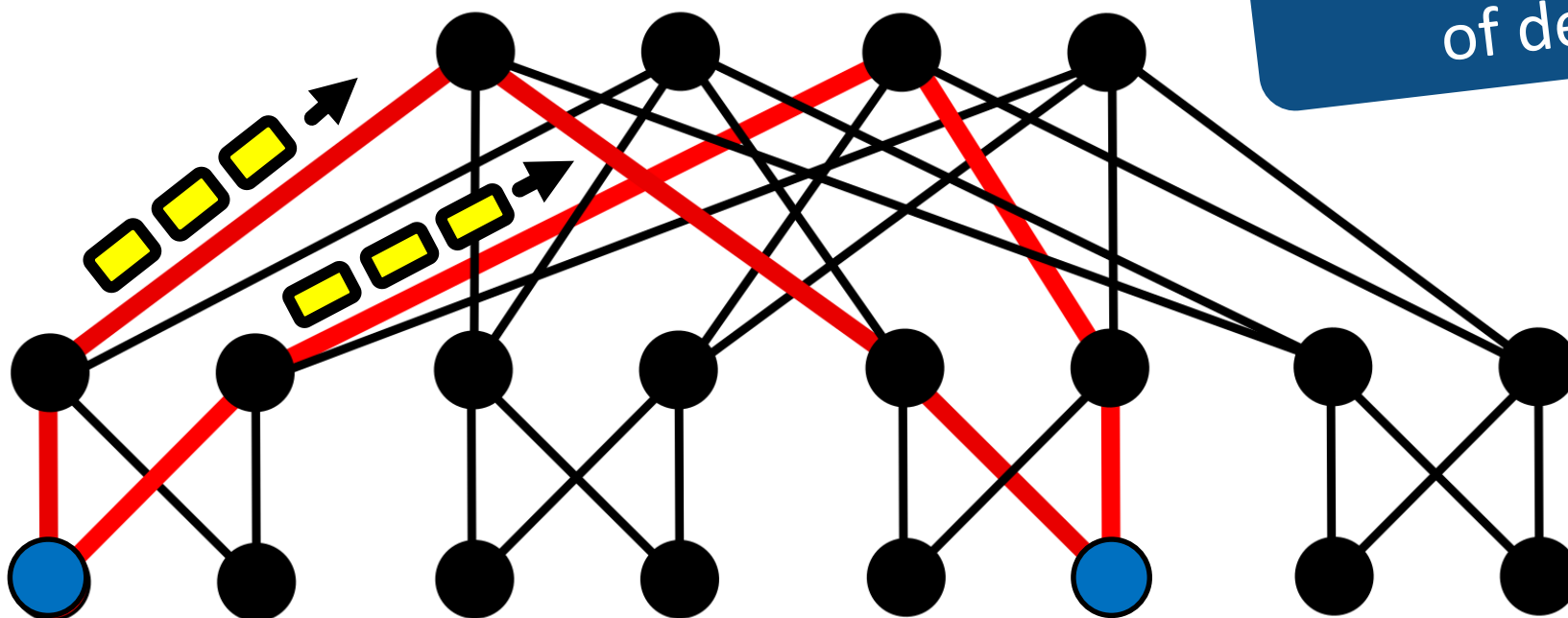
**ECMP**

INFORMATIONAL

C. Hopps  
NextHop Technologies  
November 2000

Analysis of an Equal-Cost Multi-Path Algorithm

Established techniques for using multipathing & plethora of designs available

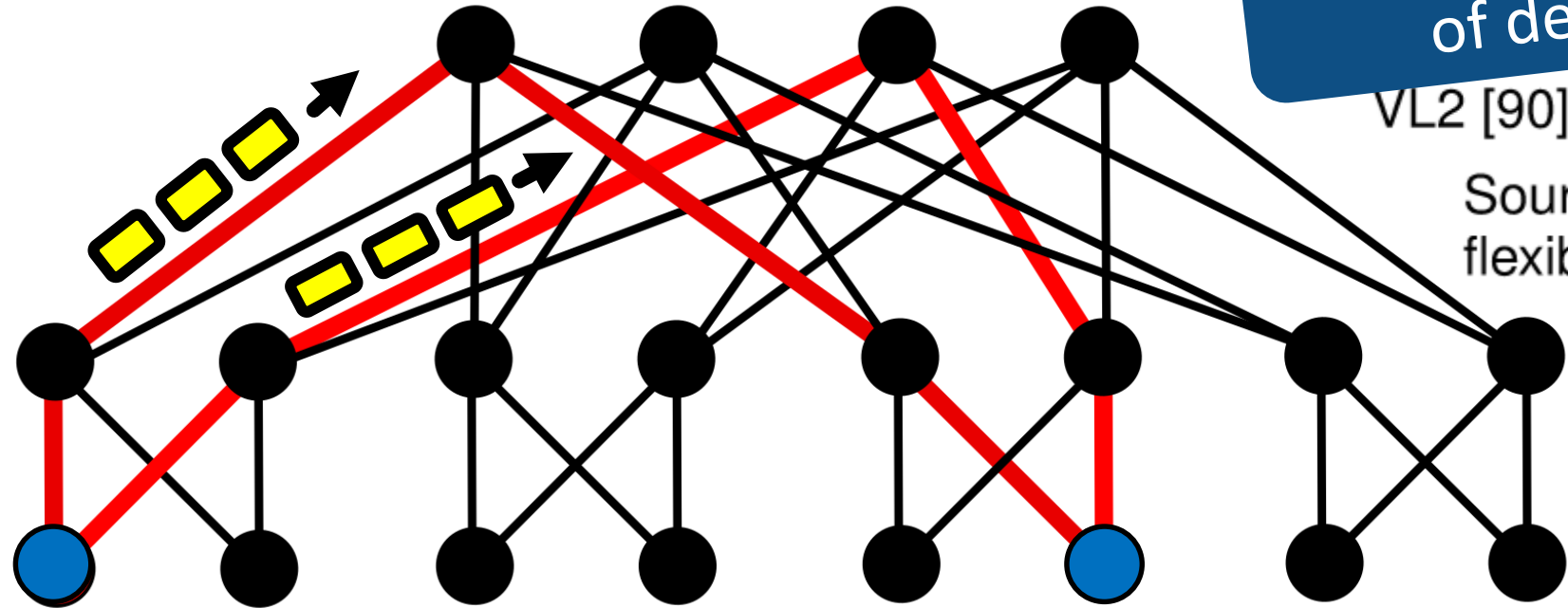


# ROUTING IN FAT TREES



! High-performance routing is facilitated by numerous multiple shortest paths of equal lengths between any endpoints

Established techniques for using multipathing & plethora of designs available



INFORMATIONAL

Network Working Group  
Request for Comments: 2992  
Category: Informational

**ECMP**

C. Hopps  
NextHop Technologies  
November 2000

Analysis of an Equal-Cost Multi-Path Algorithm

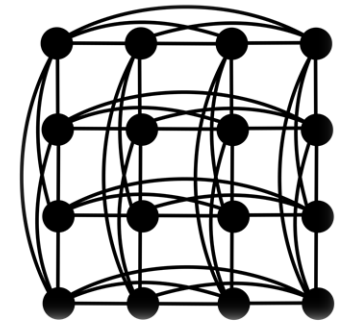
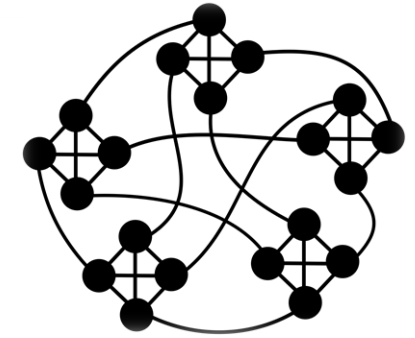
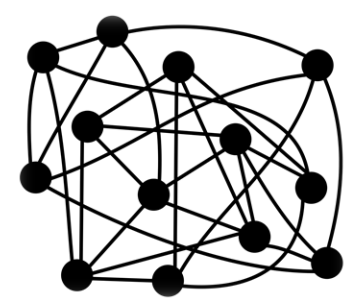
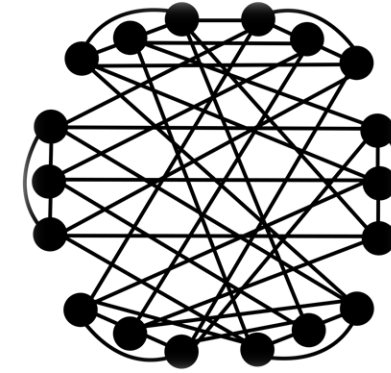
Work by Al-Fares et al. [7]  
 VL2 [90] WCMP for DC [233]  
 Source routing for flexible DC fabric [117] Monsoon [91]  
 PortLand [160] SPAIN [158]  
 ECMP-VLB [123]  
 Work by Linden et al. [215]  
 Work by Suchara et al. [204]

# ROUTING IN LOW-DIAMETER TOPOLOGIES?





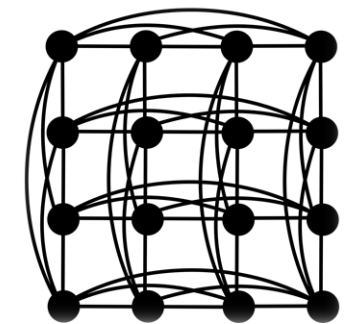
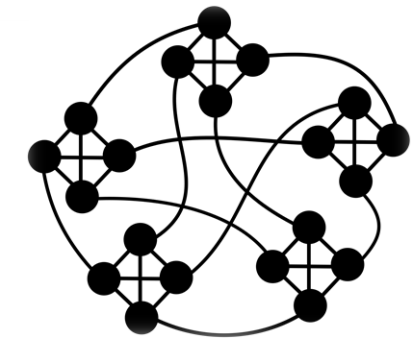
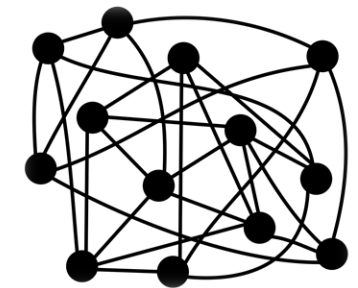
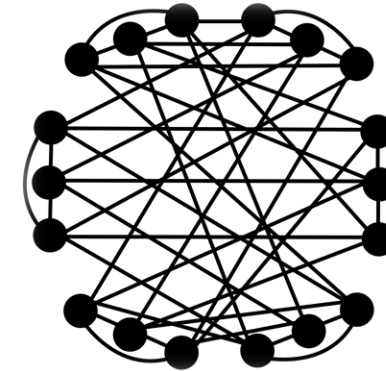
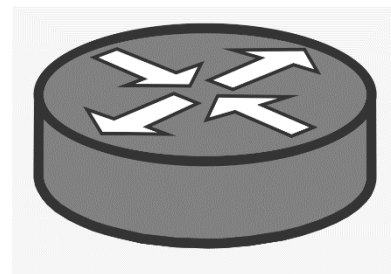
# ROUTING IN LOW-DIAMETER TOPOLOGIES?



# ROUTING IN LOW-DIAMETER TOPOLOGIES?



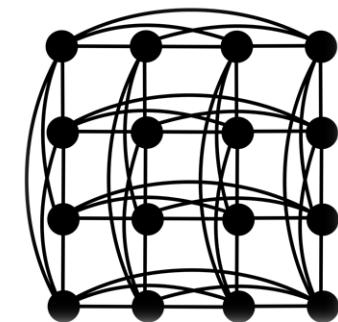
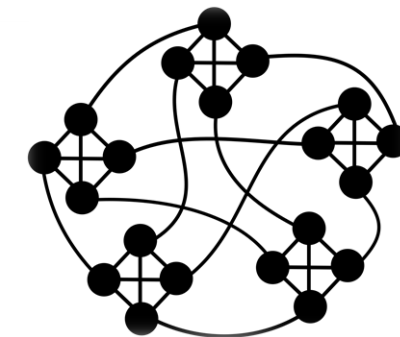
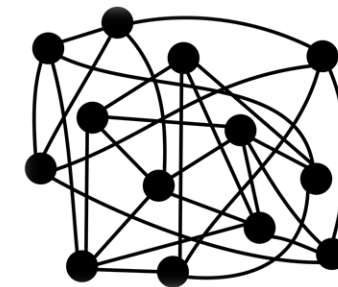
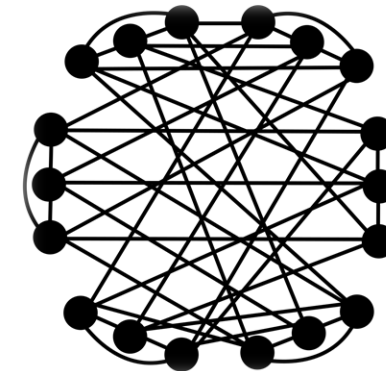
How to route modern low diameter topologies using multipathing?



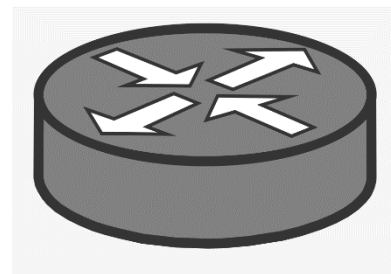
# ROUTING IN LOW-DIAMETER TOPOLOGIES?



How to route modern low diameter topologies using multipathing?



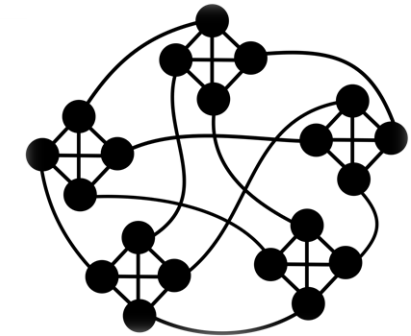
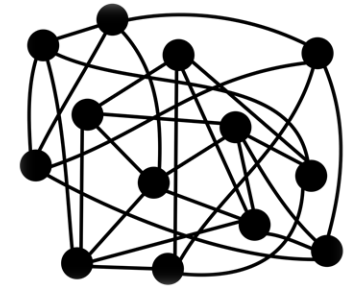
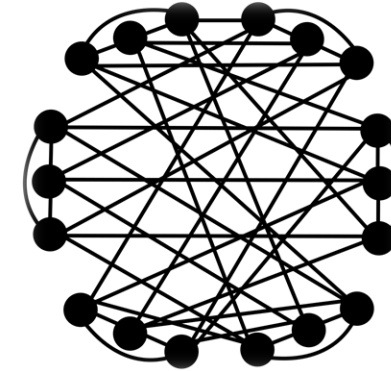
The driving question of this whole work



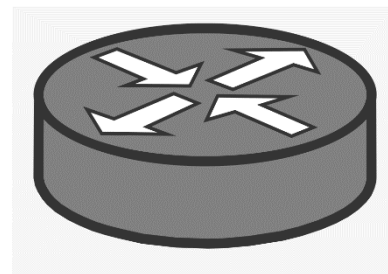
# ROUTING IN LOW-DIAMETER TOPOLOGIES?



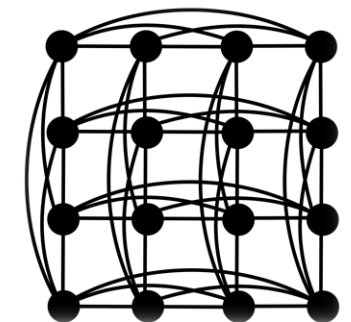
How to route modern low diameter topologies using multipathing?



The driving question of this whole work



...but it touches on many aspects. Let's go!



# MULTIPATH ROUTING: MOTIVATION



## MULTIPATH ROUTING: MOTIVATION



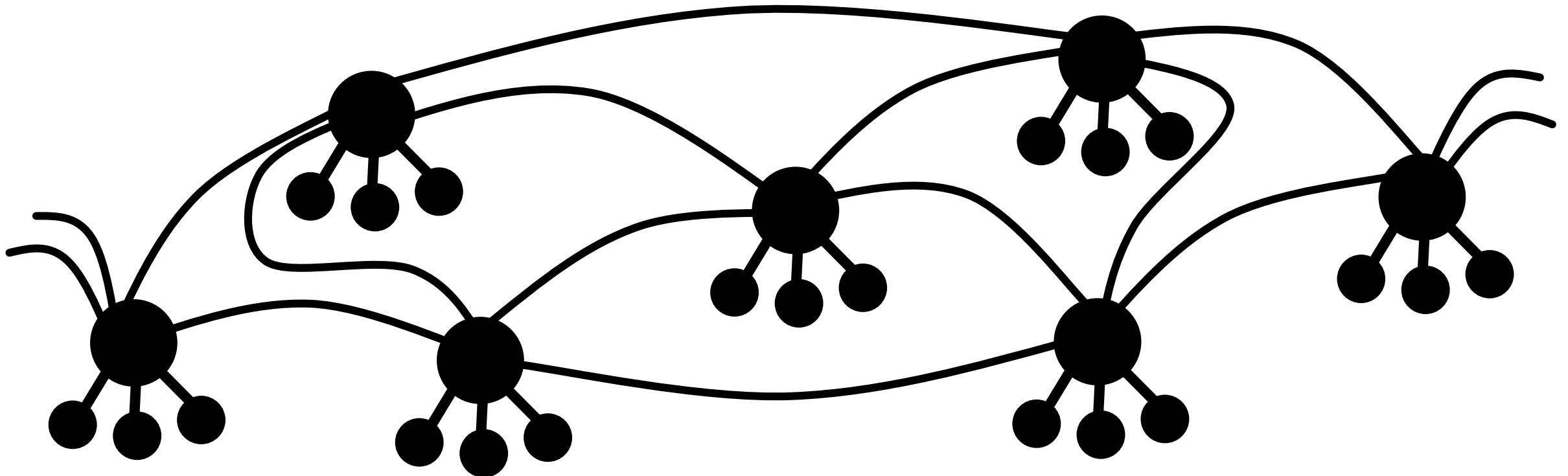
What are the problems that we want to tackle with multipathing?



## MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?



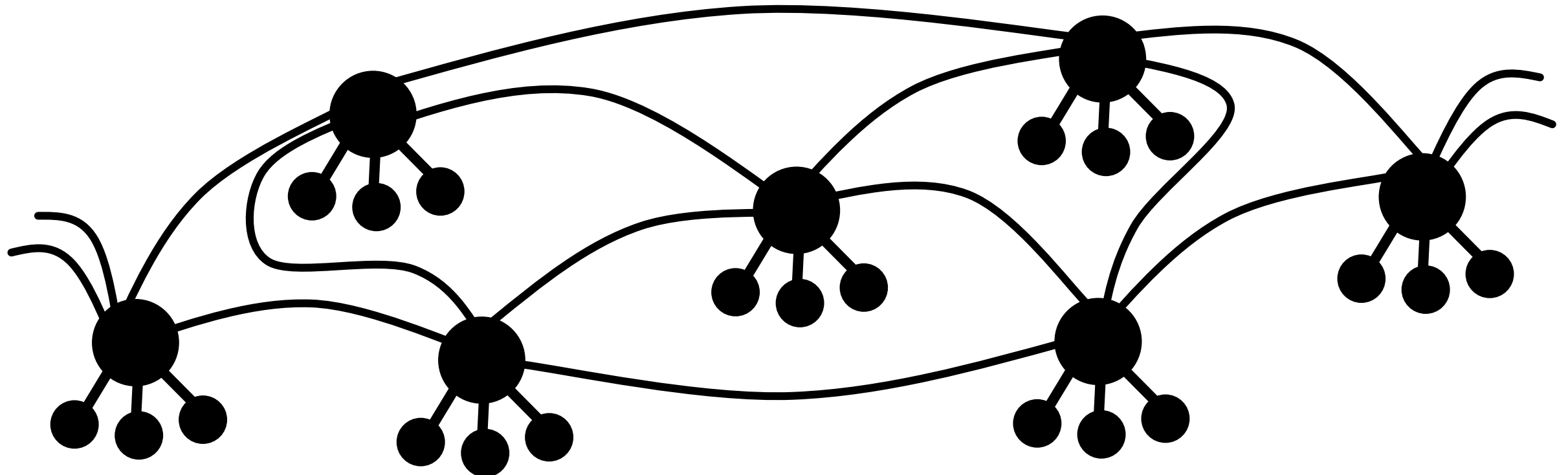
## MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?



Flows collide!



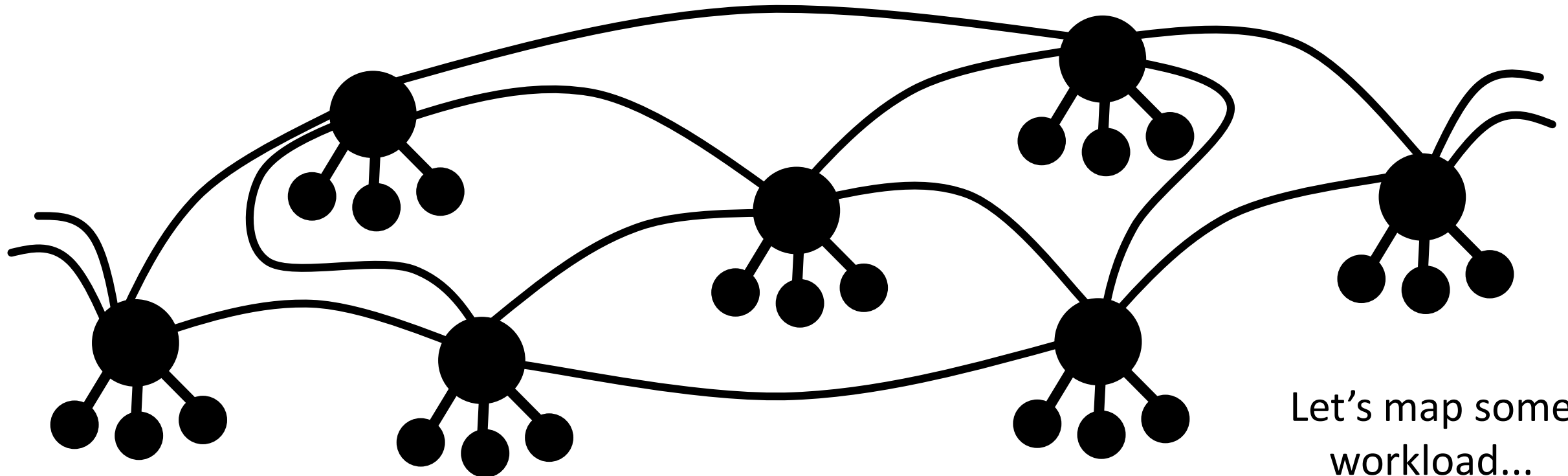


# MULTIPATH ROUTING: MOTIVATION



**Flows collide!**

What are the problems that we want to tackle with multipathing?



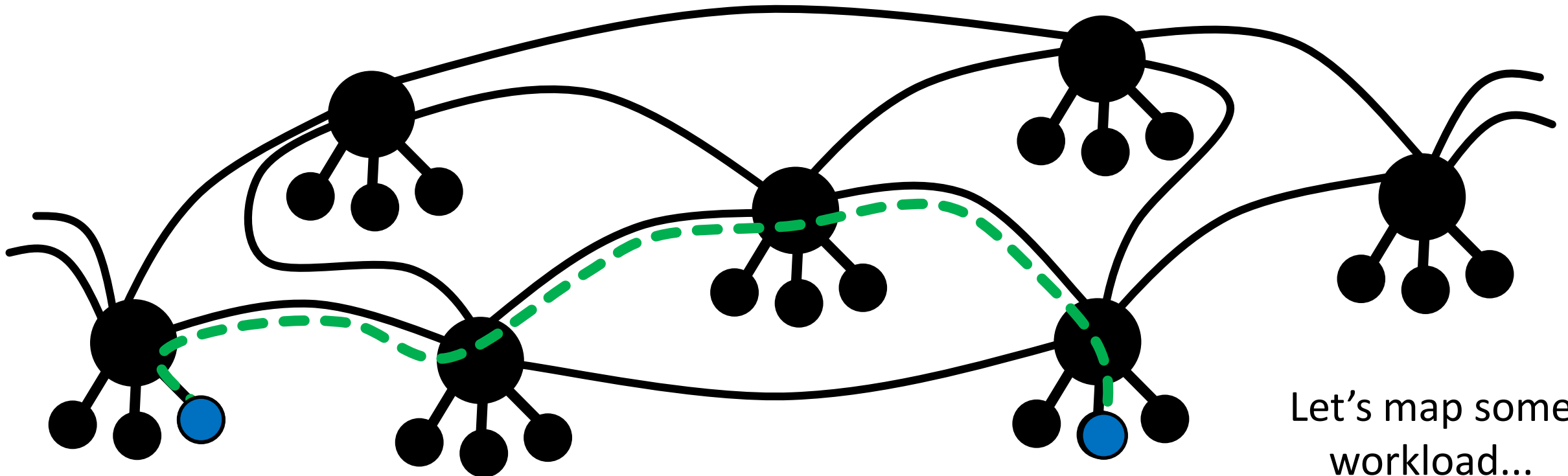
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



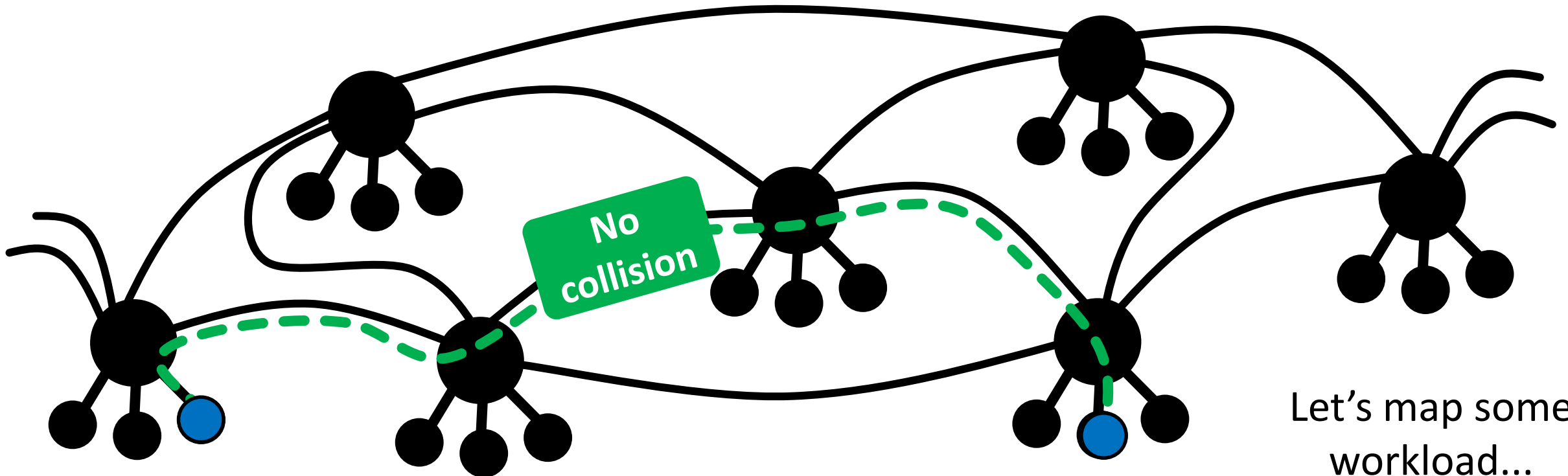
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



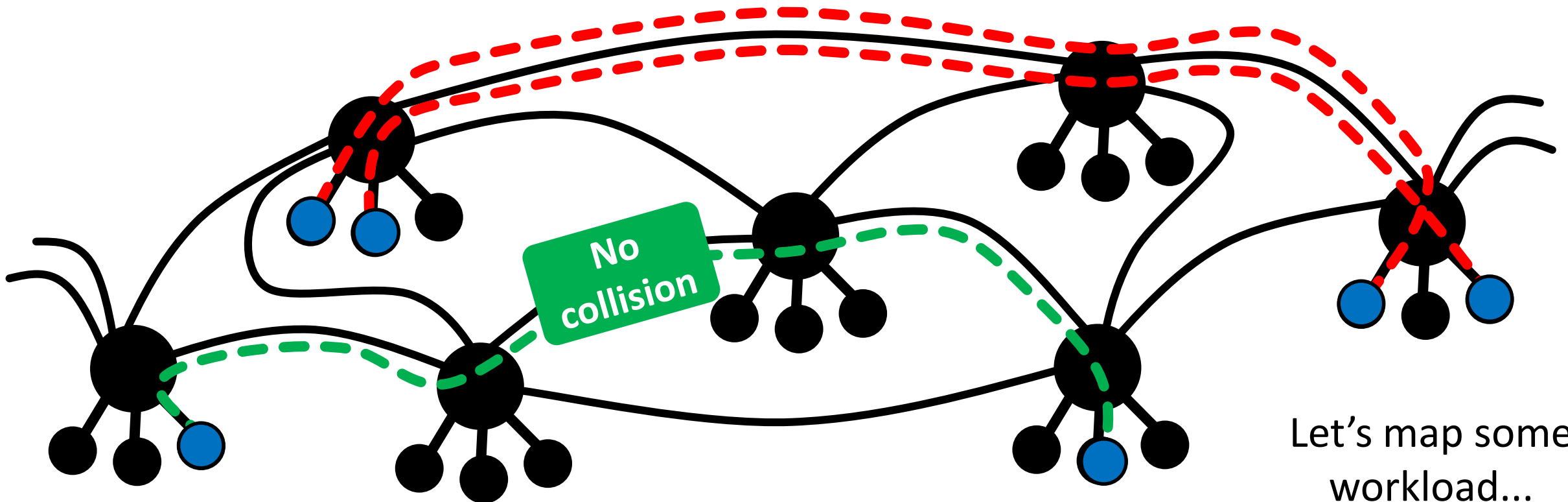
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



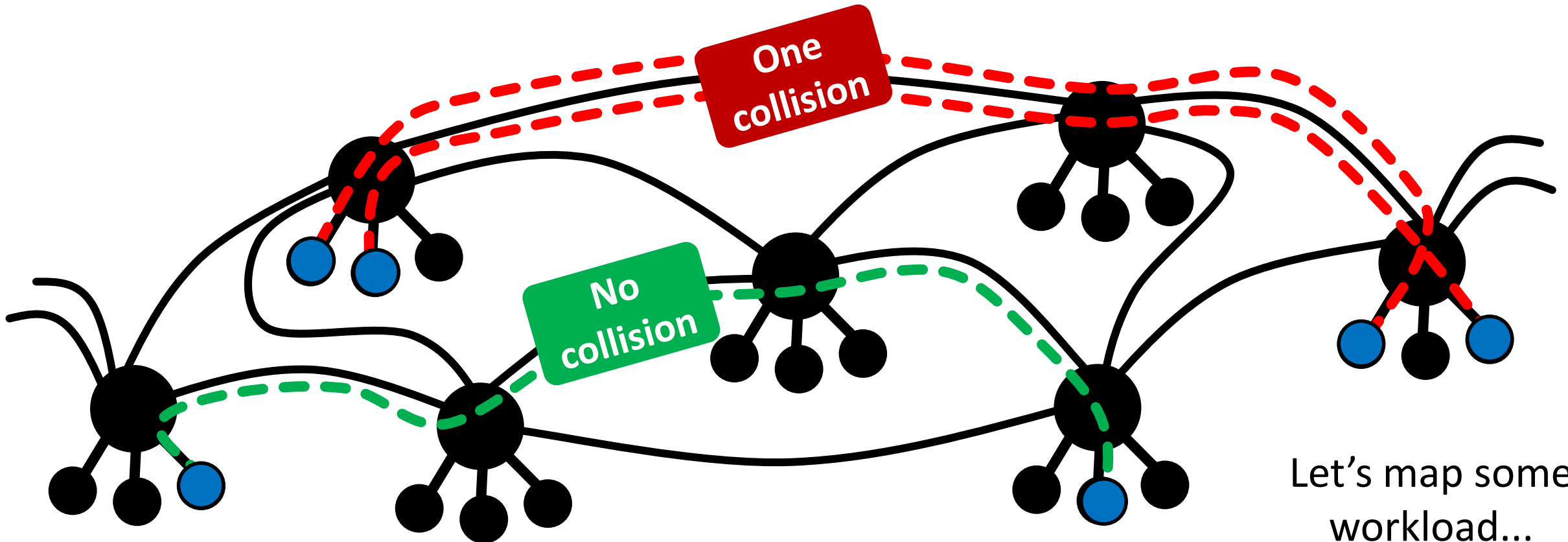
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



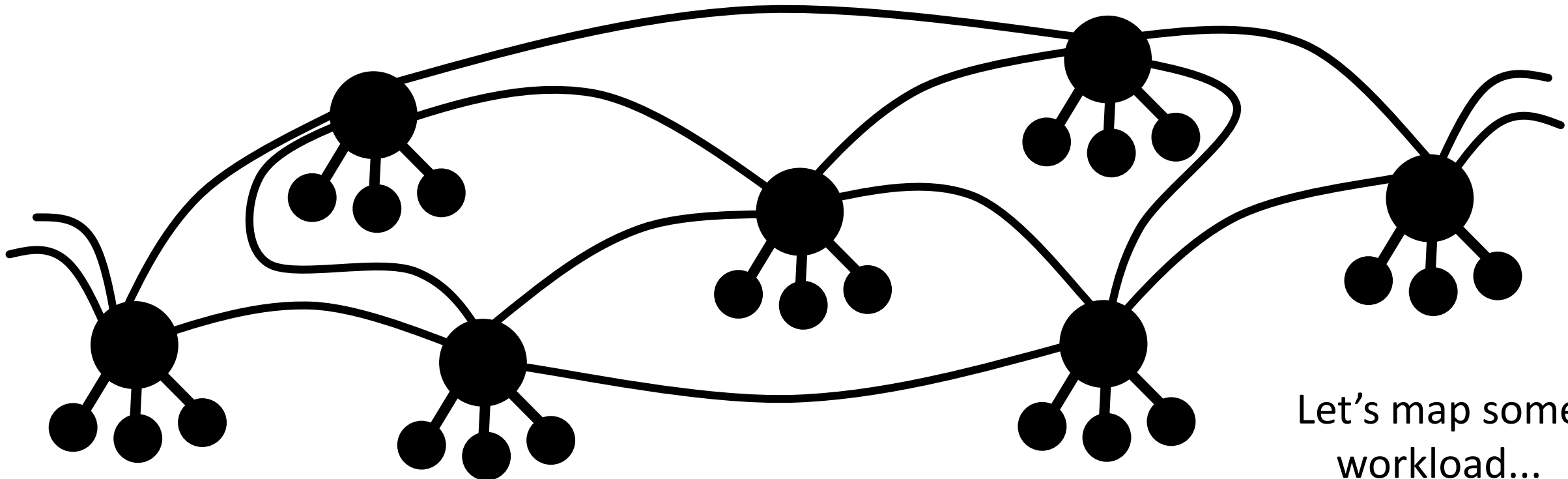
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



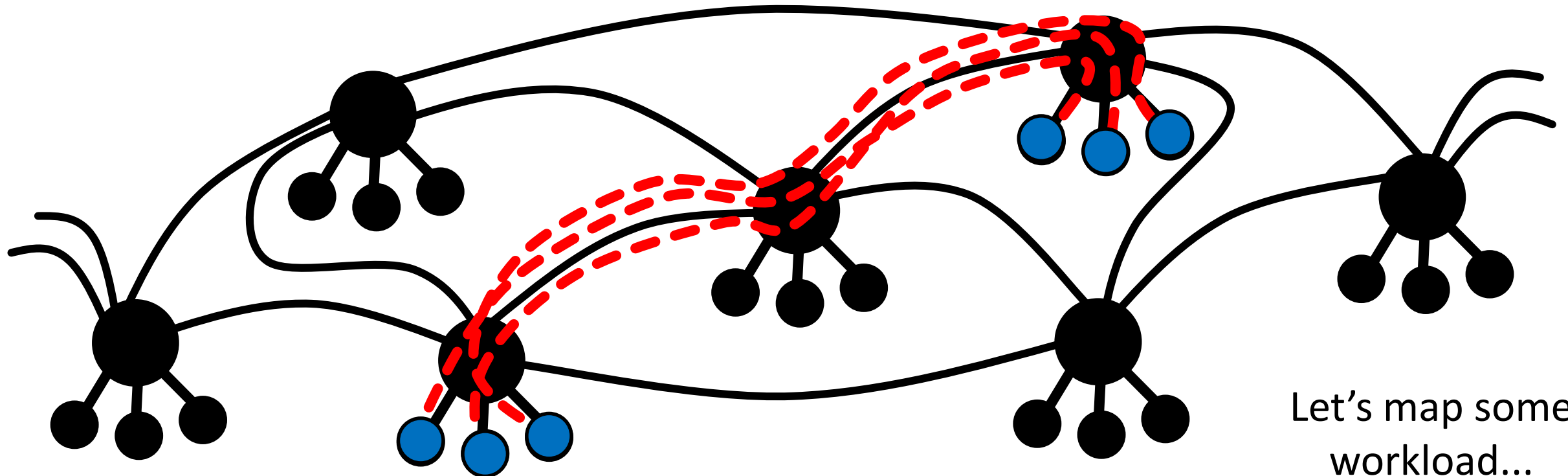
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



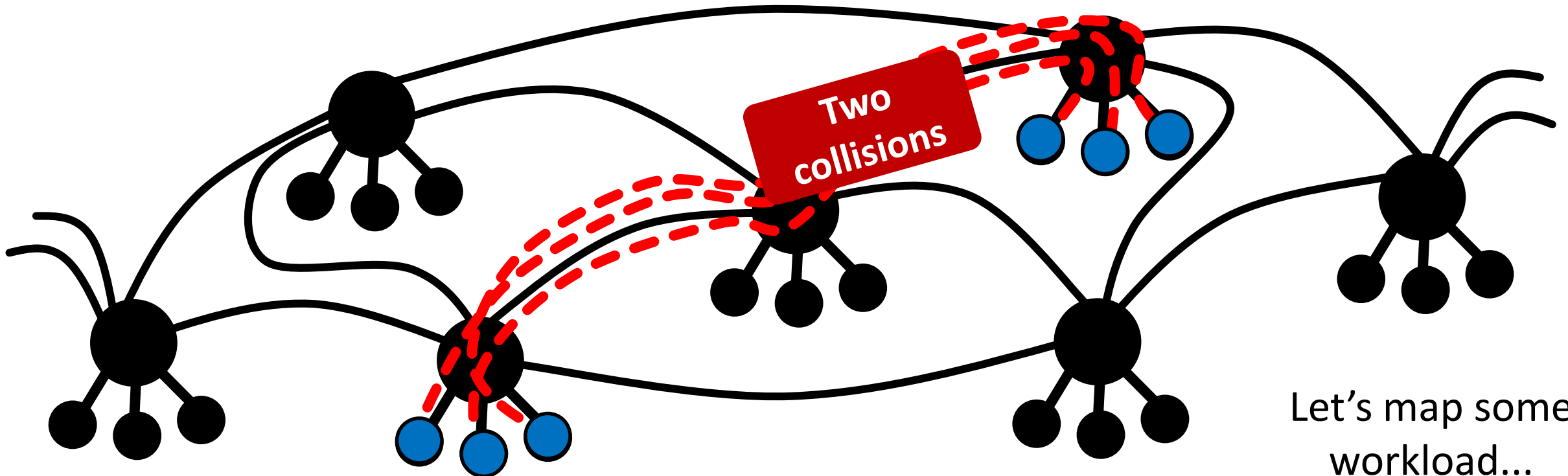
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Let's map some workload...

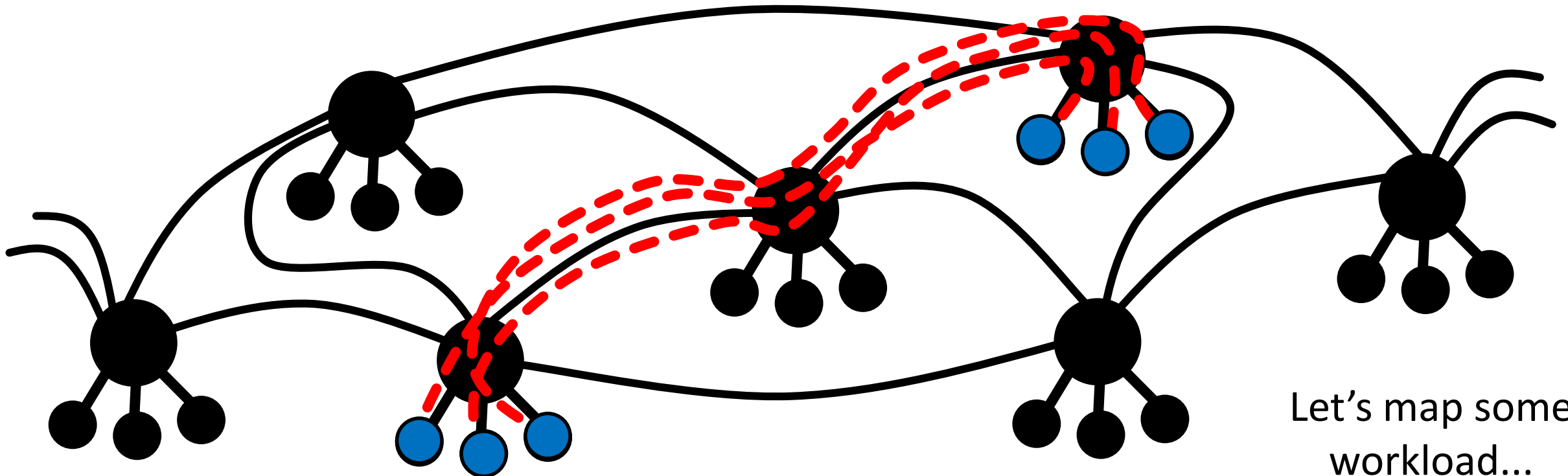


# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

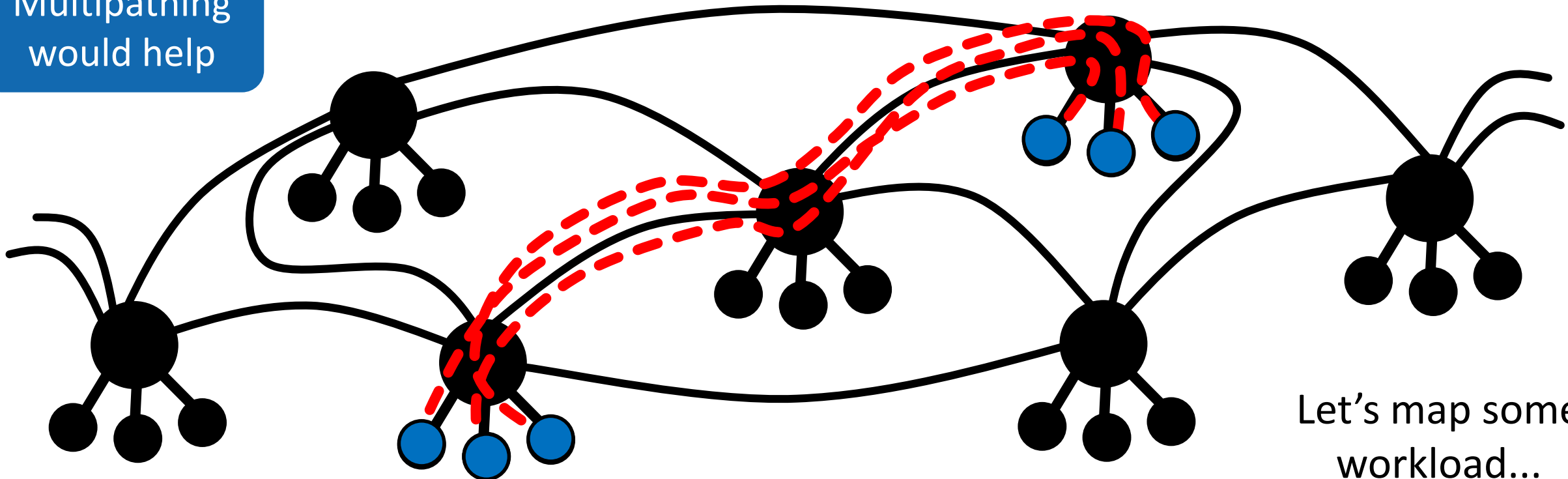


**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

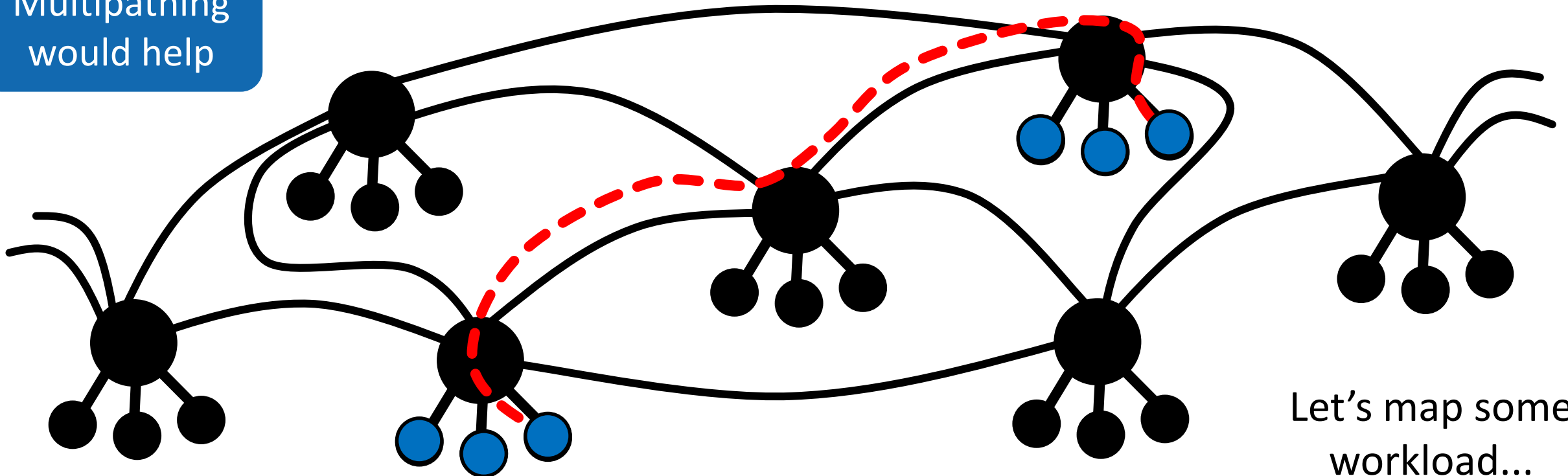


**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

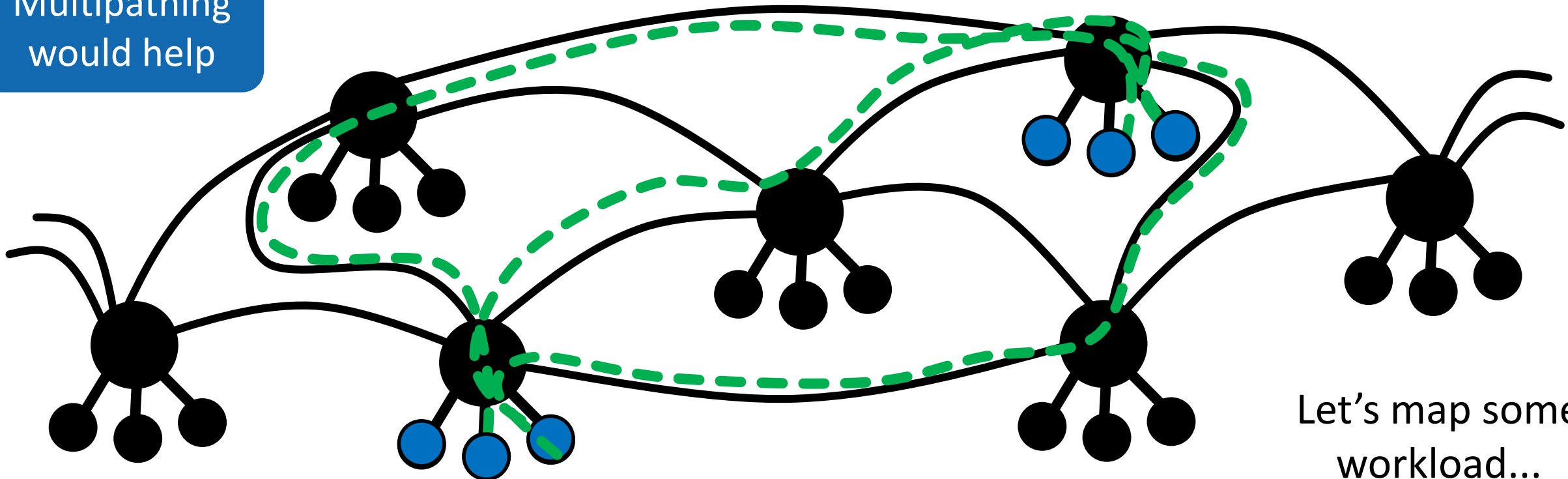


**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

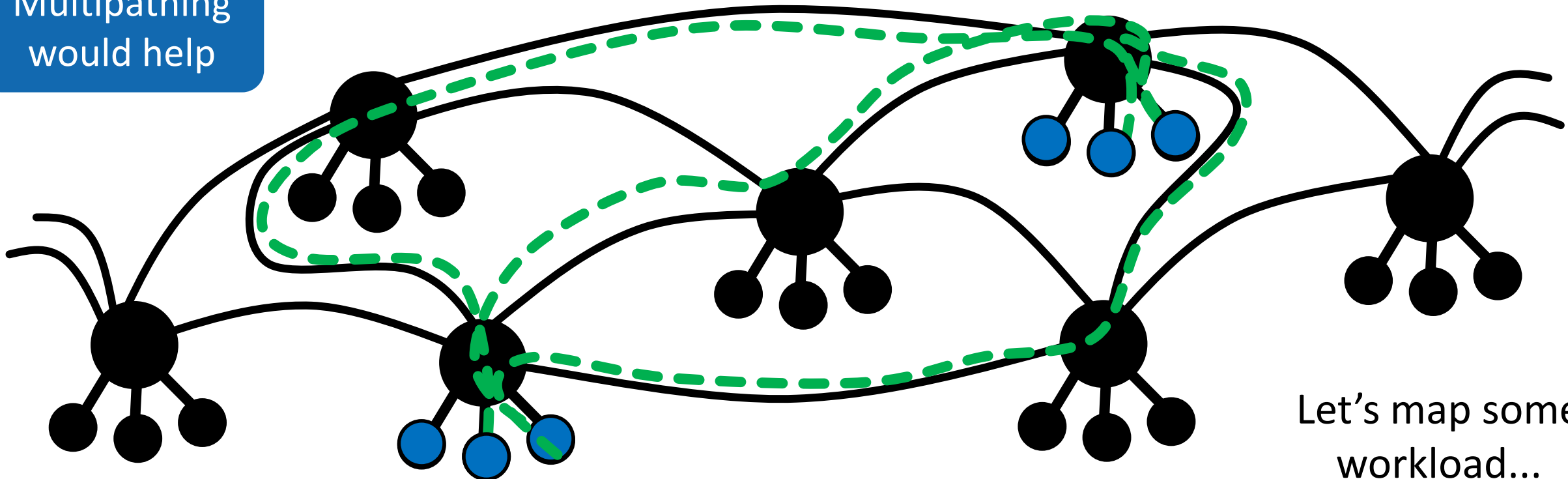


**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate“ collisions?

Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

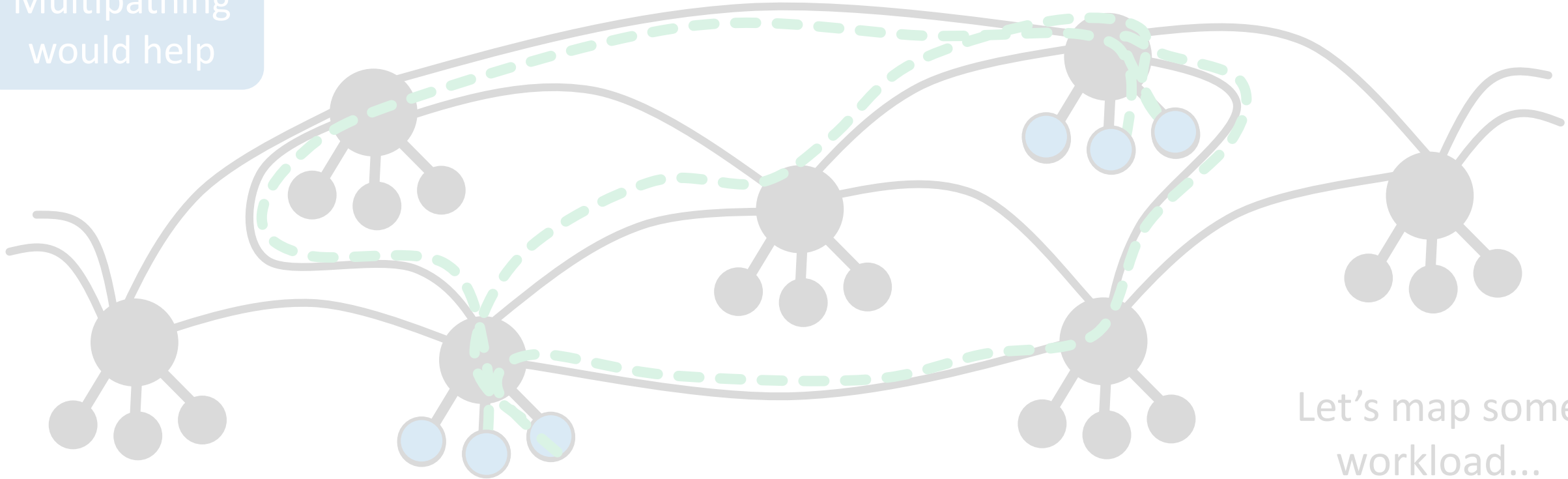


**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate” collisions?

Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

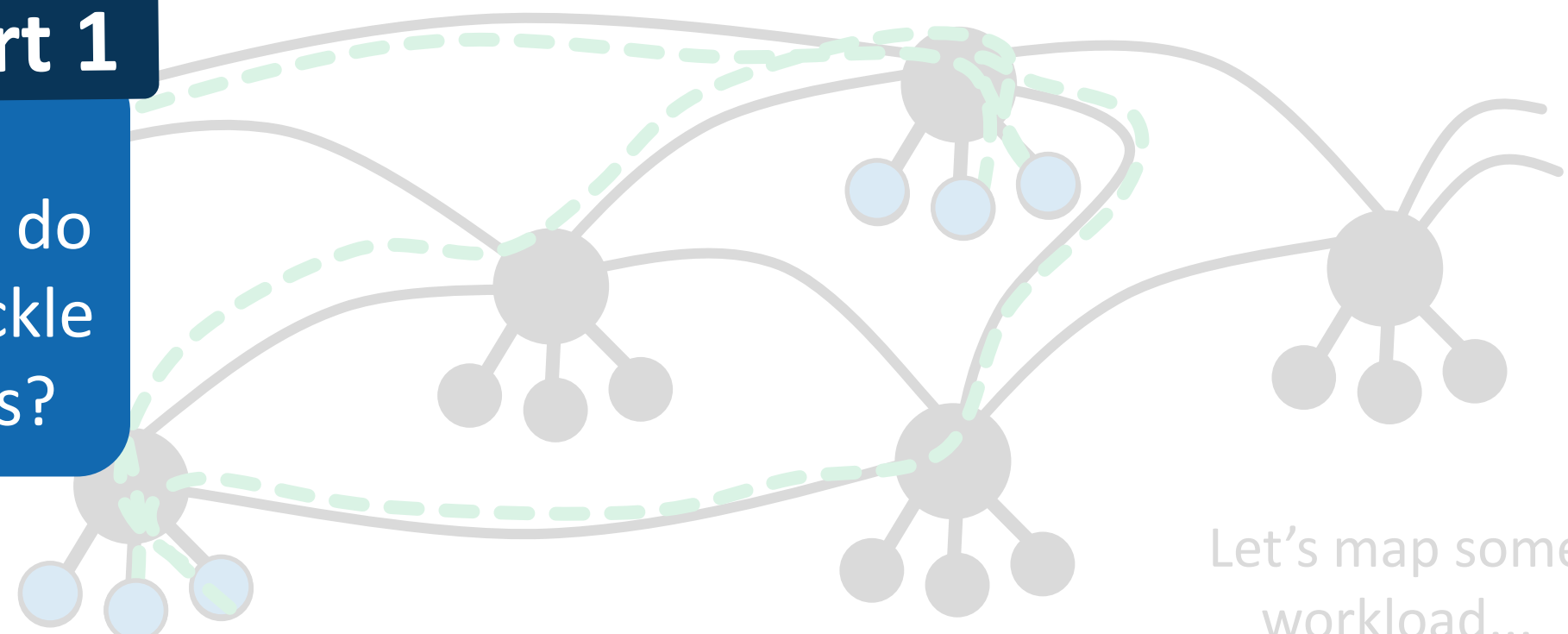
What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate“ collisions?

Multipathing

## Part 1

How many multiple paths do we need to tackle flow collisions?



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate” collisions?

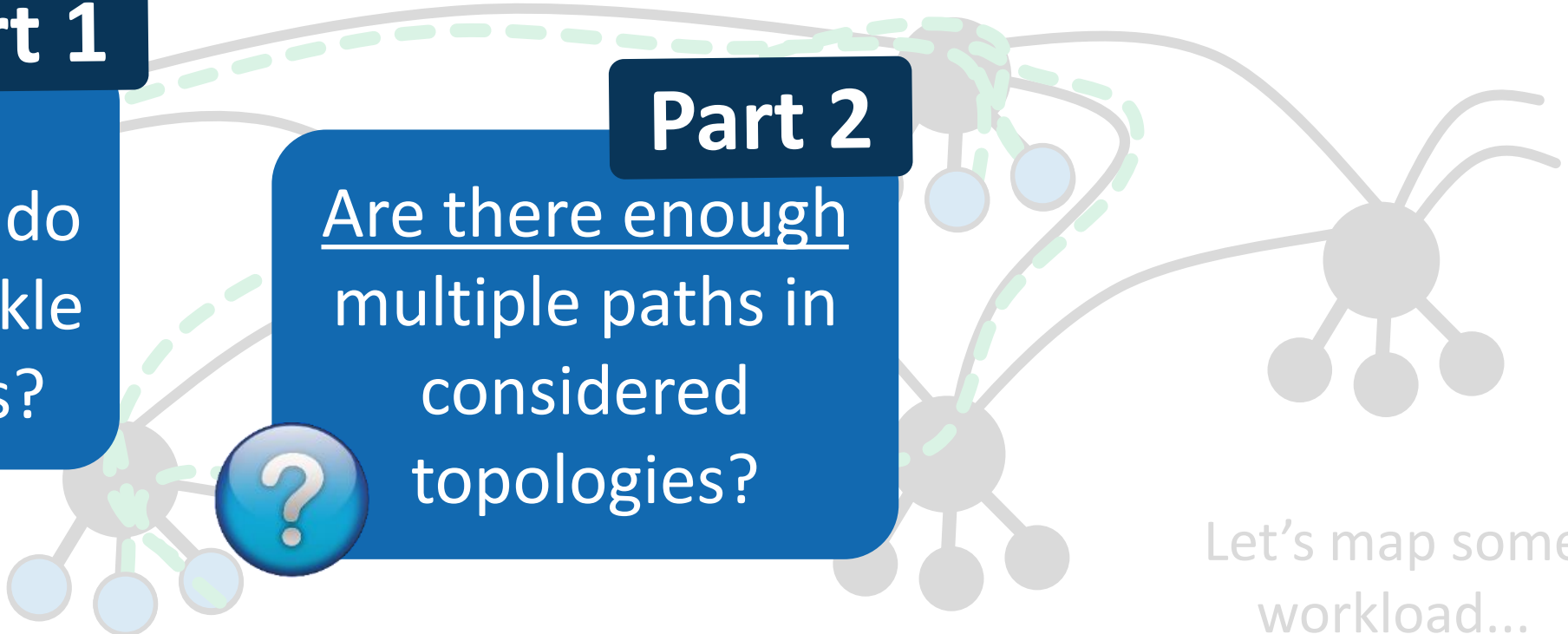
Multipathing

## Part 1

How many multiple paths do we need to tackle flow collisions?

## Part 2

Are there enough multiple paths in considered topologies?



Let's map some workload...



# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate“ collisions?

Multipathing

## Part 1

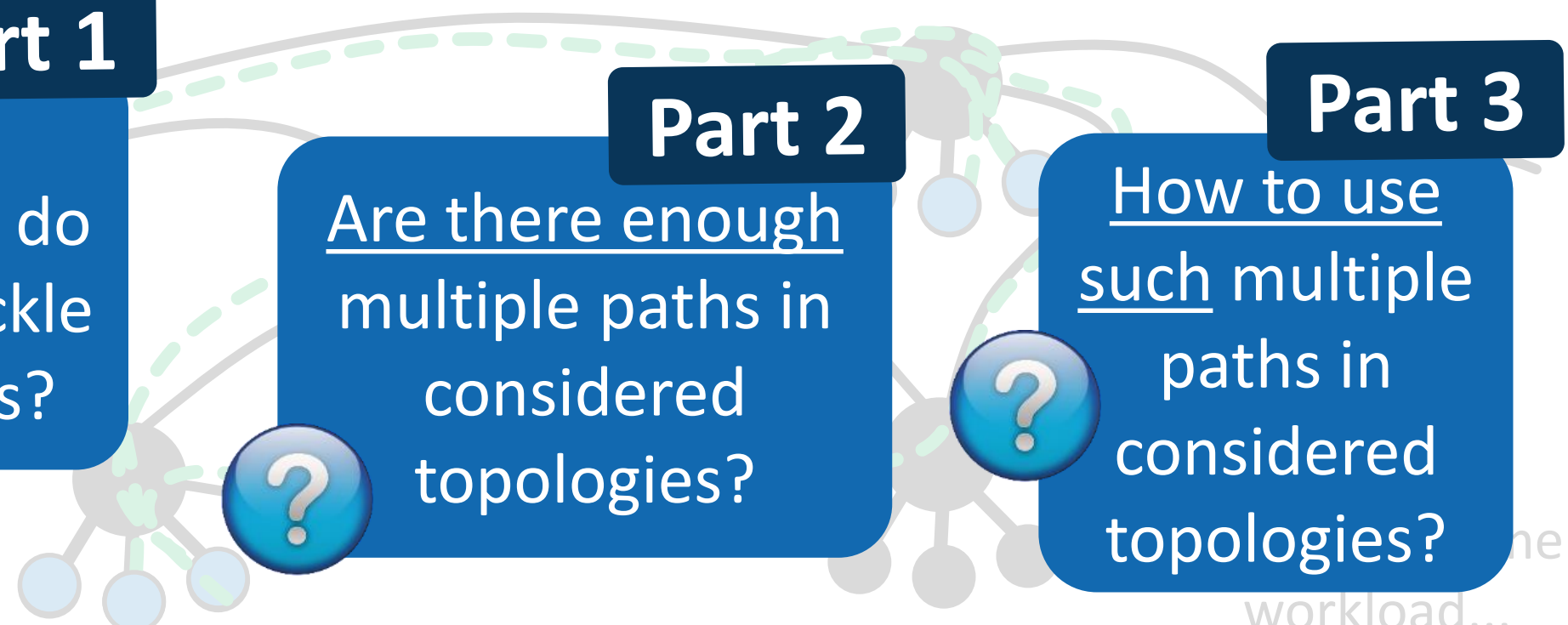
How many multiple paths do we need to tackle flow collisions?

## Part 2

Are there enough multiple paths in considered topologies?

## Part 3

How to use such multiple paths in considered topologies?

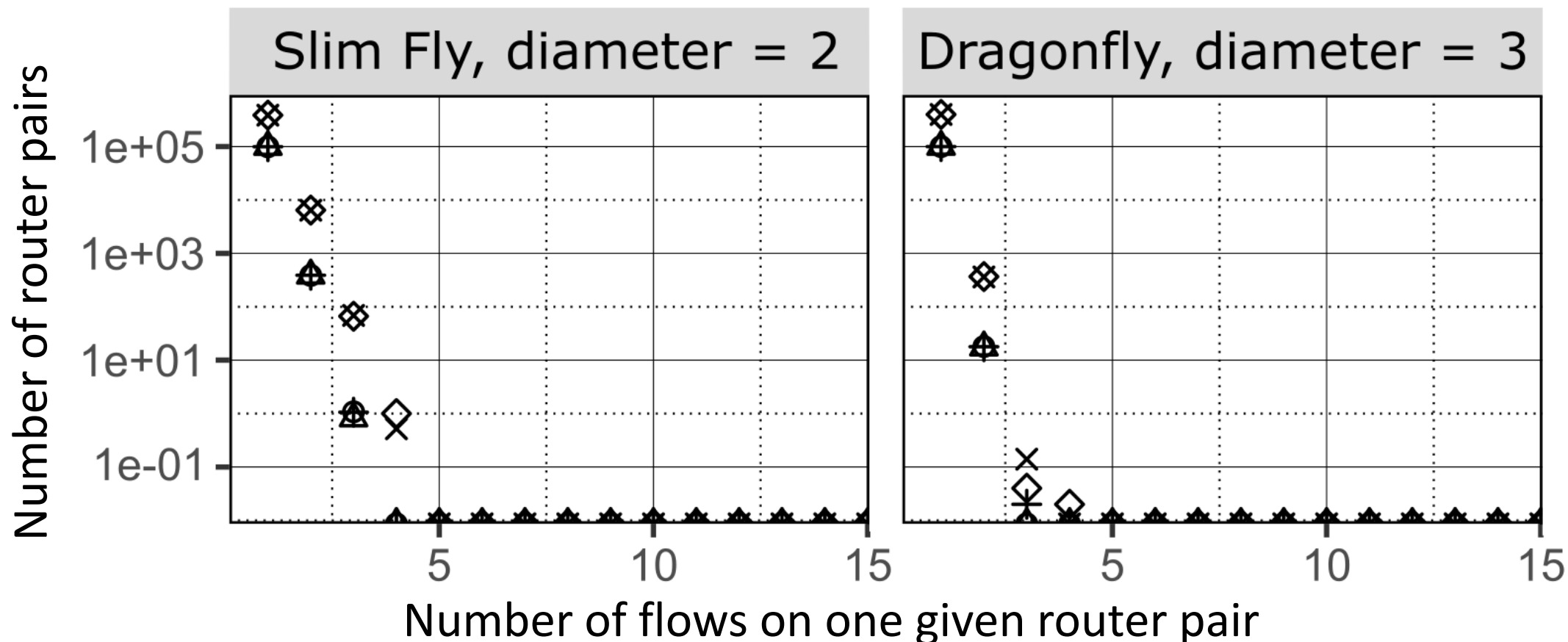


# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

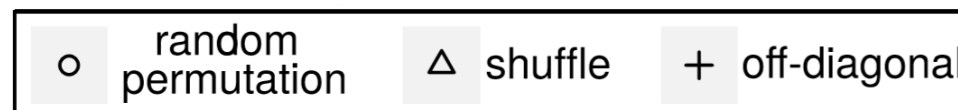
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

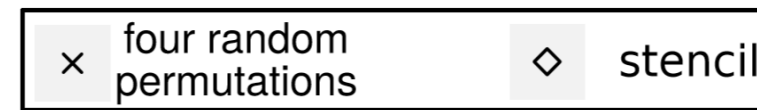


≈10,000 servers;  
comparable cost

### Workloads

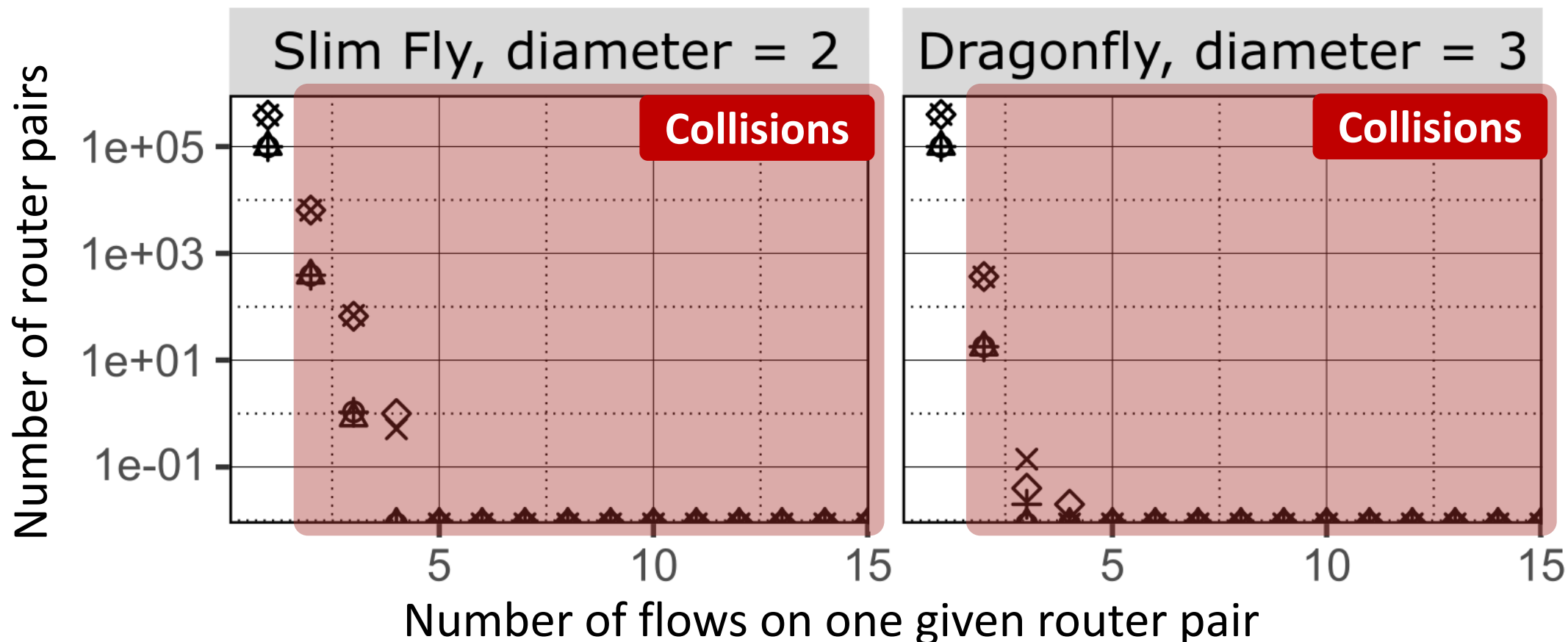


### 4x-oversubscribed workloads



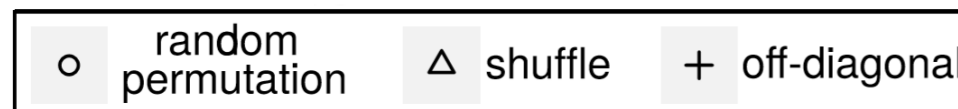
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

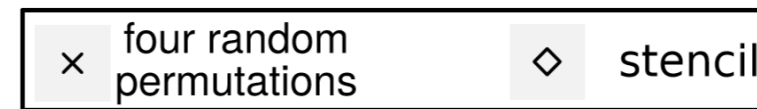


≈10,000 servers;  
comparable cost

### Workloads

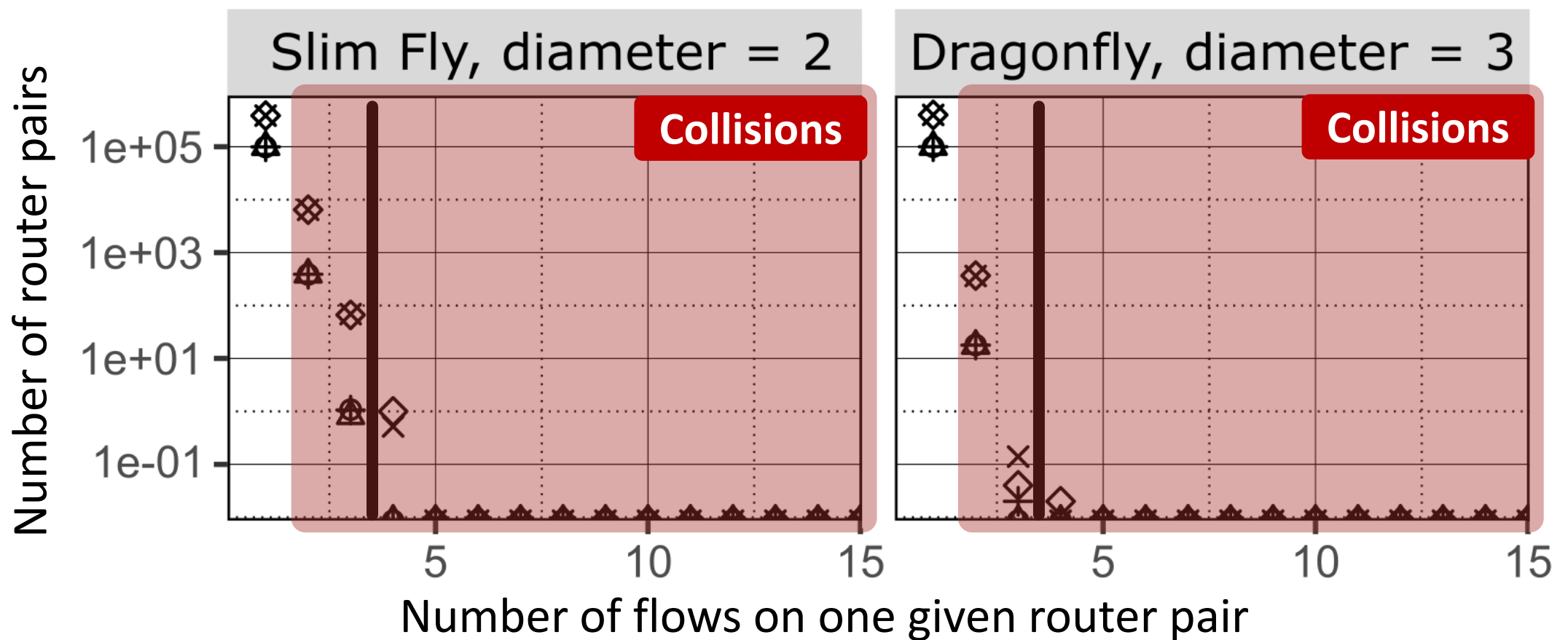


### 4x-oversubscribed workloads



# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1



≈10,000 servers;  
comparable cost

Workloads

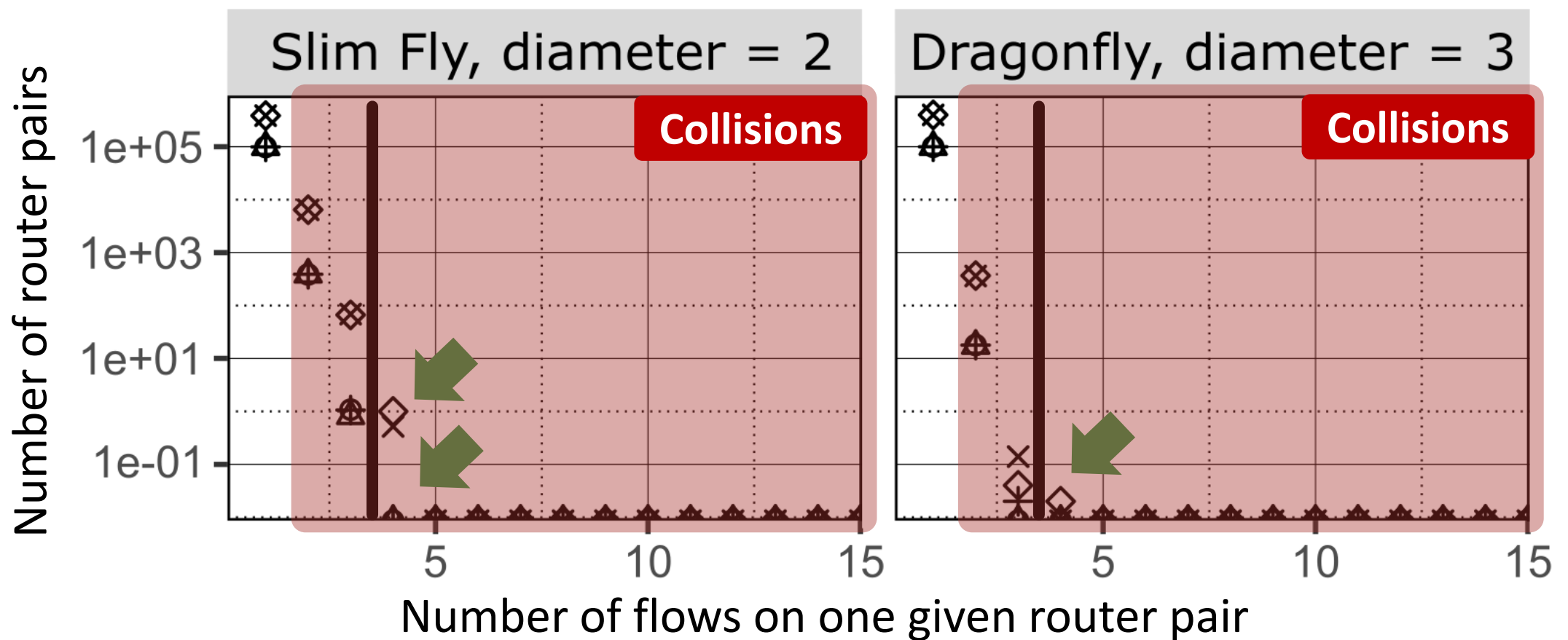
- random permutation
- △ shuffle
- + off-diagonal

4x-oversubscribed workloads

- × four random permutations
- ◇ stencil

# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1



≈10,000 servers;  
comparable cost

Workloads

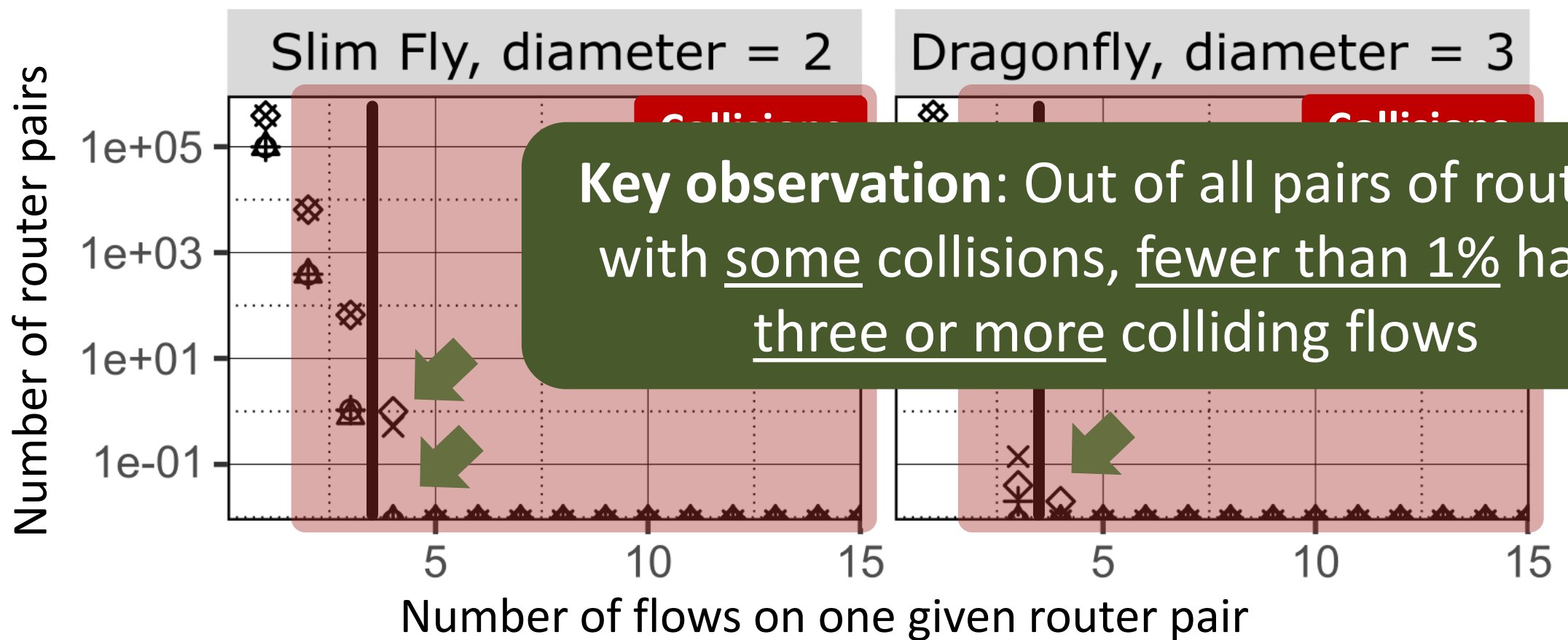
- random permutation
- △ shuffle
- + off-diagonal

4x-oversubscribed workloads

- × four random permutations
- ◇ stencil

# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1



≈10,000 servers;  
comparable cost

Workloads

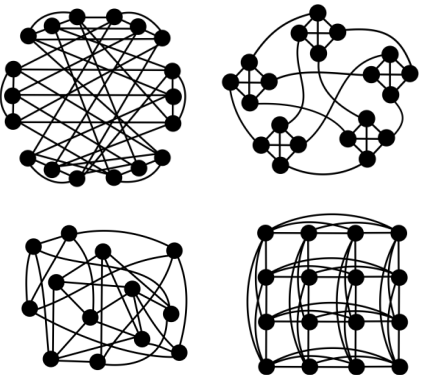
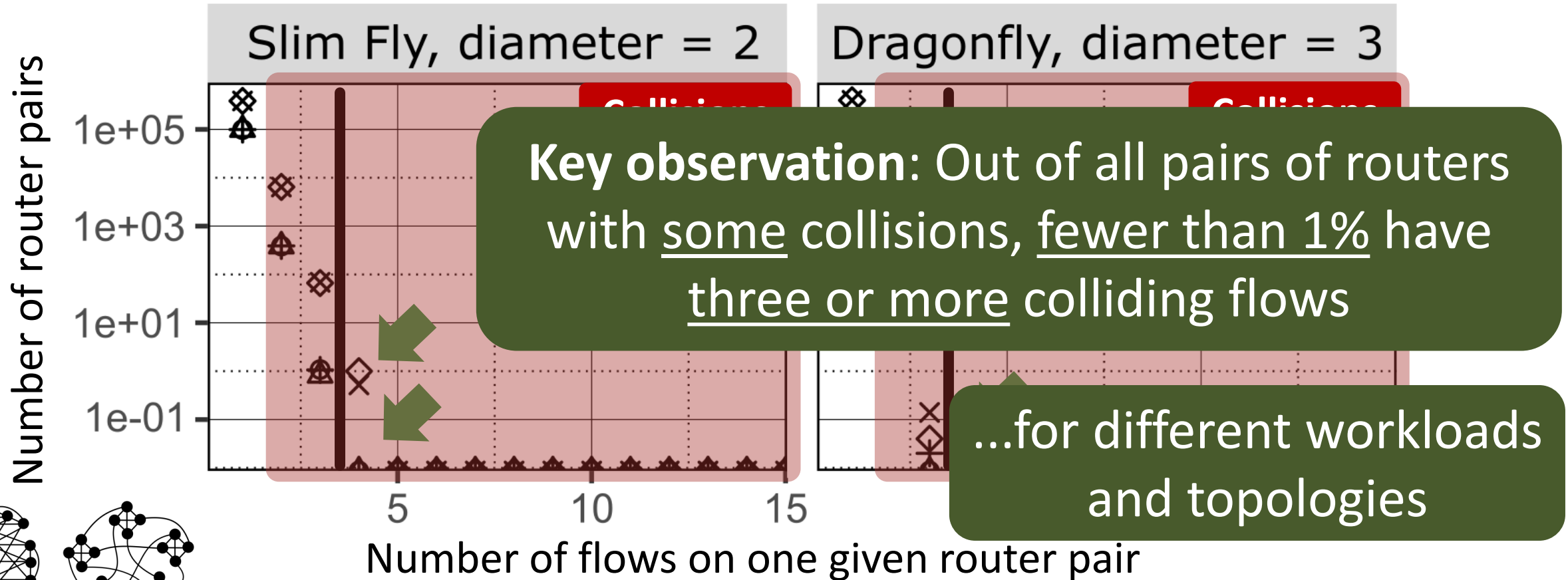
- random permutation
- △ shuffle
- + off-diagonal

4x-oversubscribed workloads

- × four random permutations
- ◇ stencil

# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1



### Workloads

○ random permutation	△ shuffle	+ off-diagonal
----------------------	-----------	----------------

### 4x-oversubscribed workloads

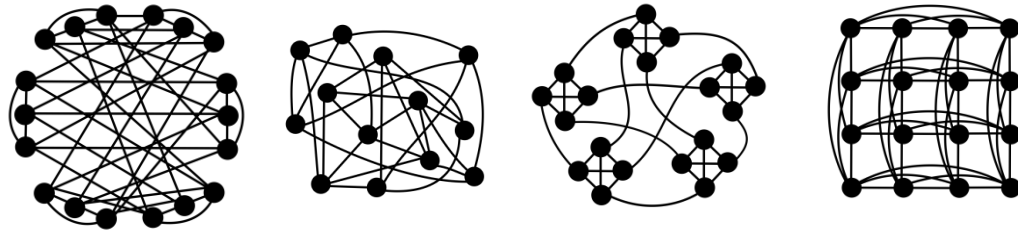
× four random permutations	◇ stencil
----------------------------	-----------



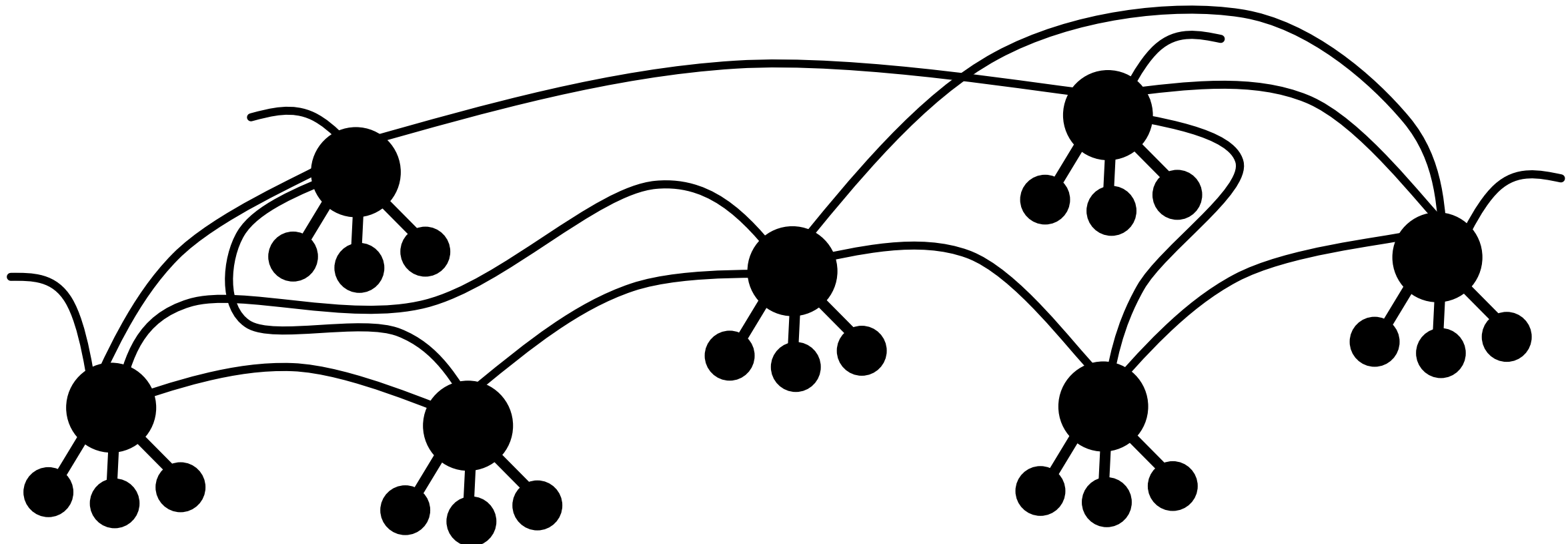
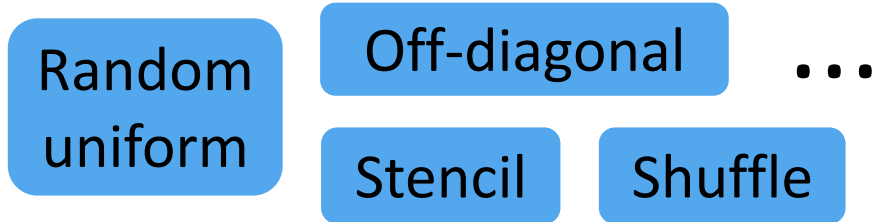
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



...For different workloads



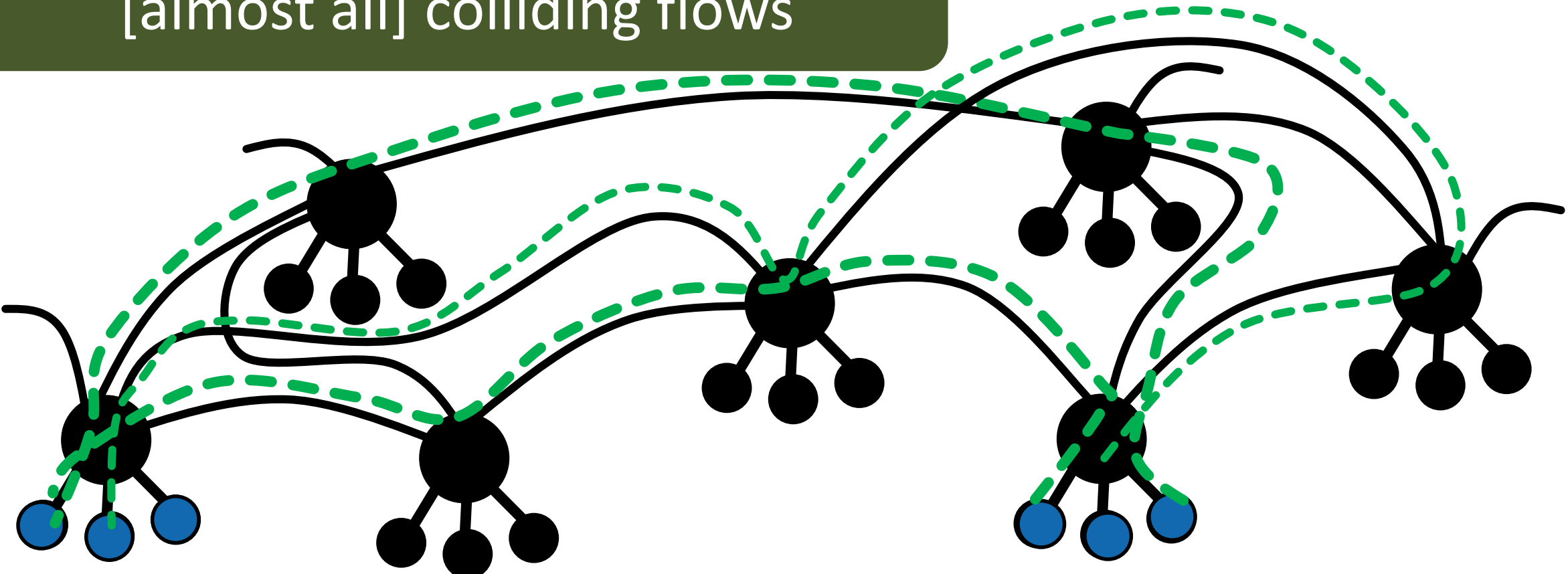
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

**Key Insight:** We need three disjoint paths per router pair to handle [almost all] colliding flows

For different workloads

- Random uniform
- Off-diagonal ...
- Stencil
- Shuffle

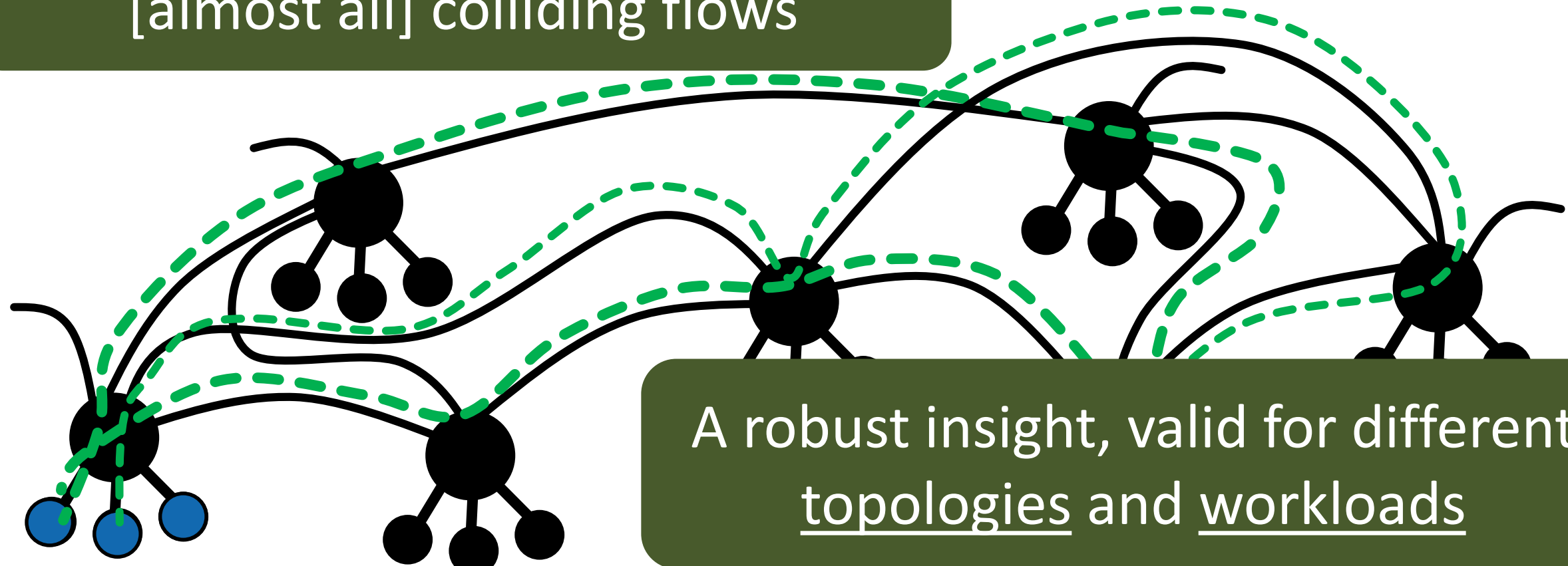


# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

**Key Insight:** We need three disjoint paths per router pair to handle [almost all] colliding flows

- For different workloads
- Random uniform
  - Off-diagonal
  - Stencil
  - Shuffle
  - ...



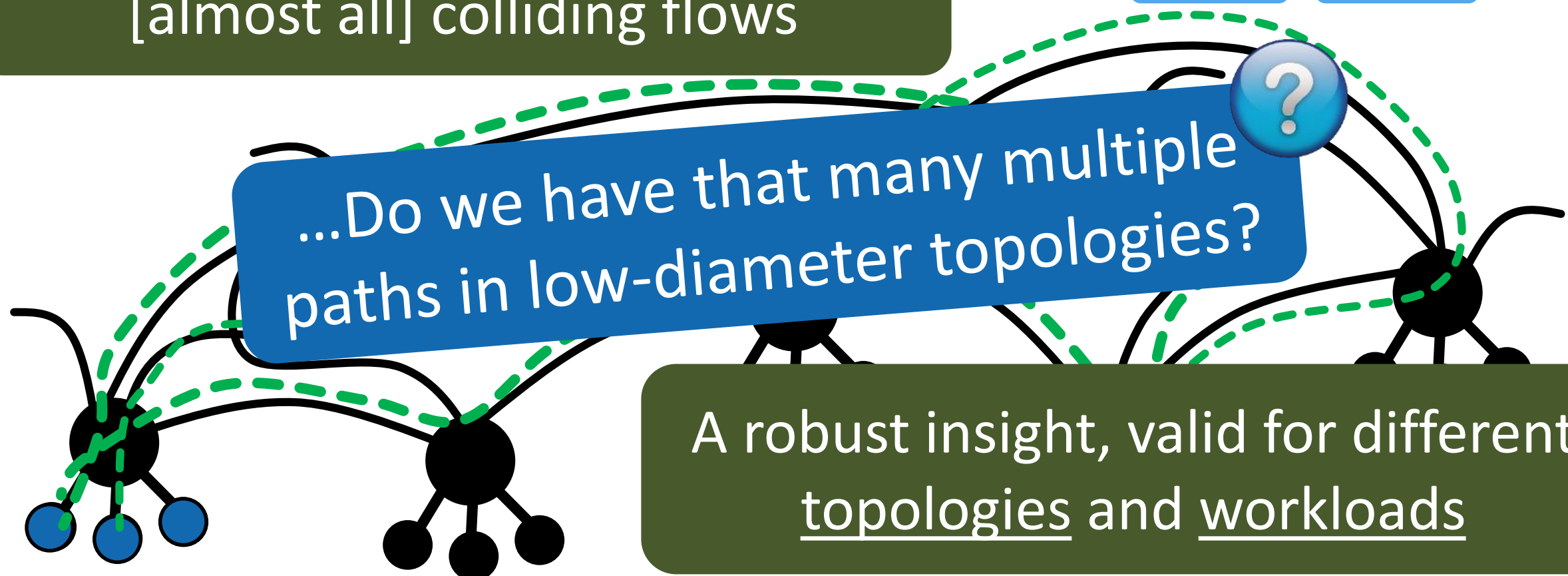
A robust insight, valid for different topologies and workloads

# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

**Key Insight:** We need three disjoint paths per router pair to handle [almost all] colliding flows

- For different workloads
- Random uniform
  - Off-diagonal ...
  - Stencil
  - Shuffle



A robust insight, valid for different topologies and workloads

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

## Part 2

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

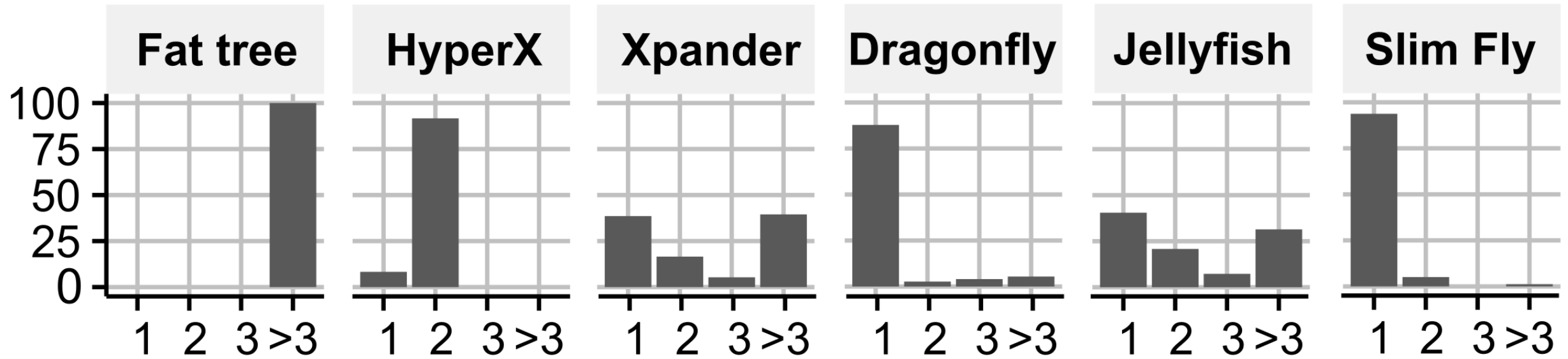
Part 2

Shortest  
paths

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**  
 Shortest paths

Percentage of router pairs



Number of shortest paths (disjoint) between a router pair

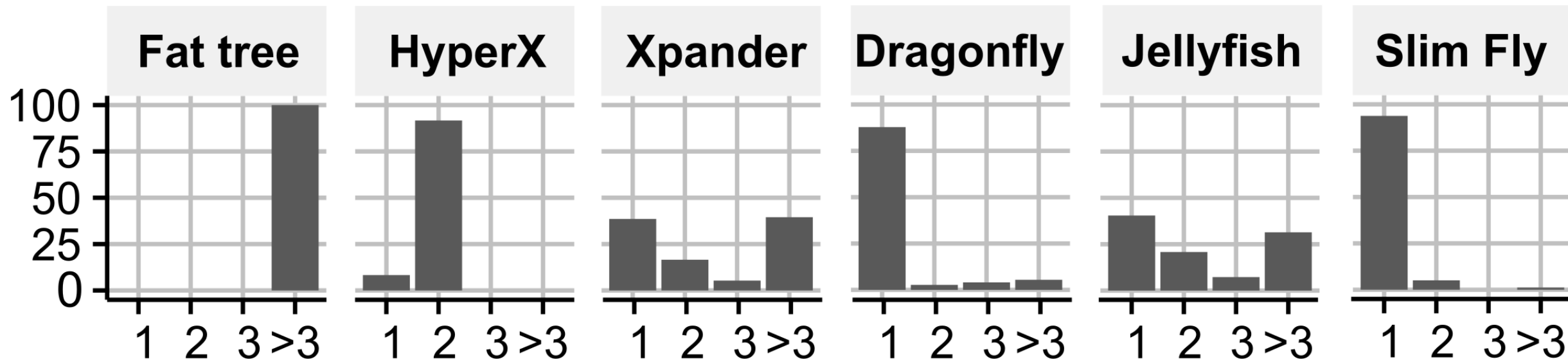
# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

Part 2

Shortest paths

≈10,000 servers;  
comparable cost

Percentage of router pairs



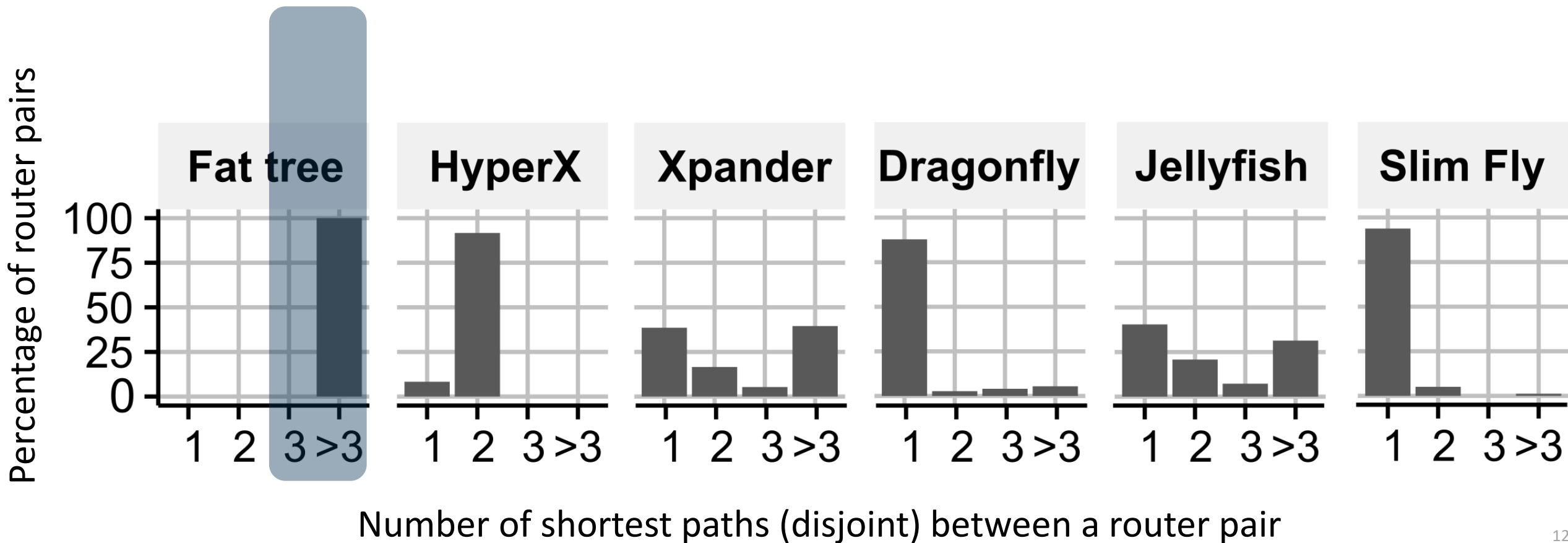
Number of shortest paths (disjoint) between a router pair



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**  
**Shortest paths**

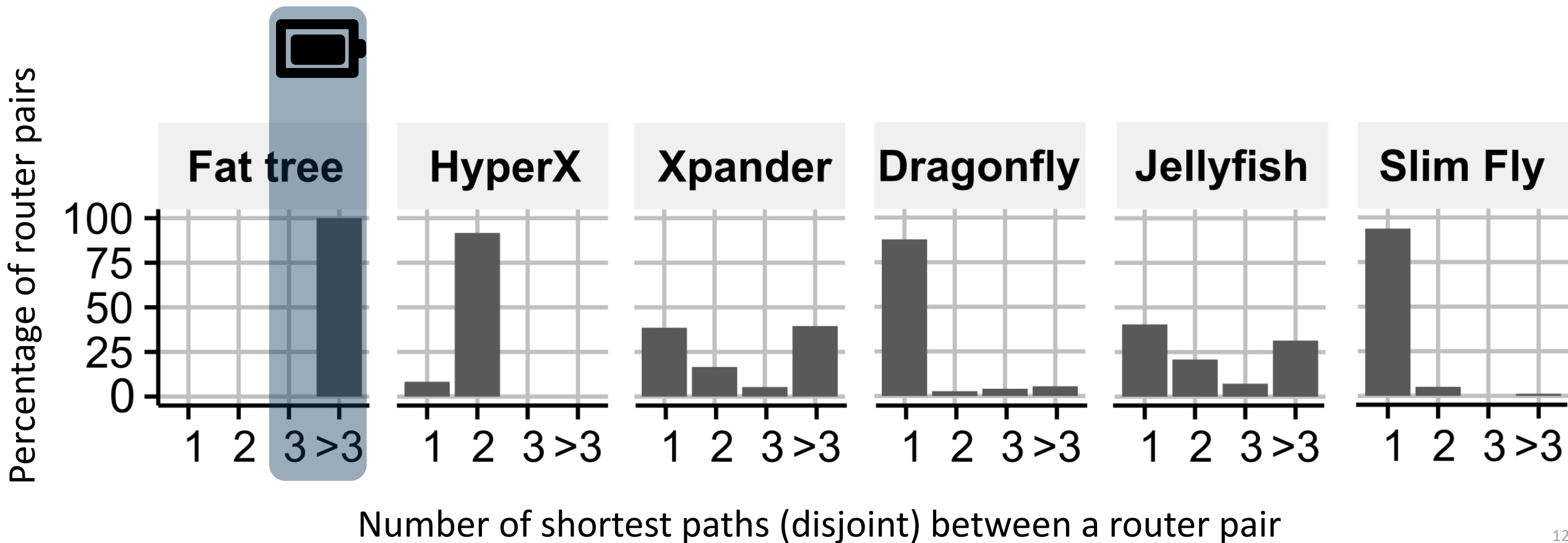
≈10,000 servers;  
 comparable cost



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**  
**Shortest paths**

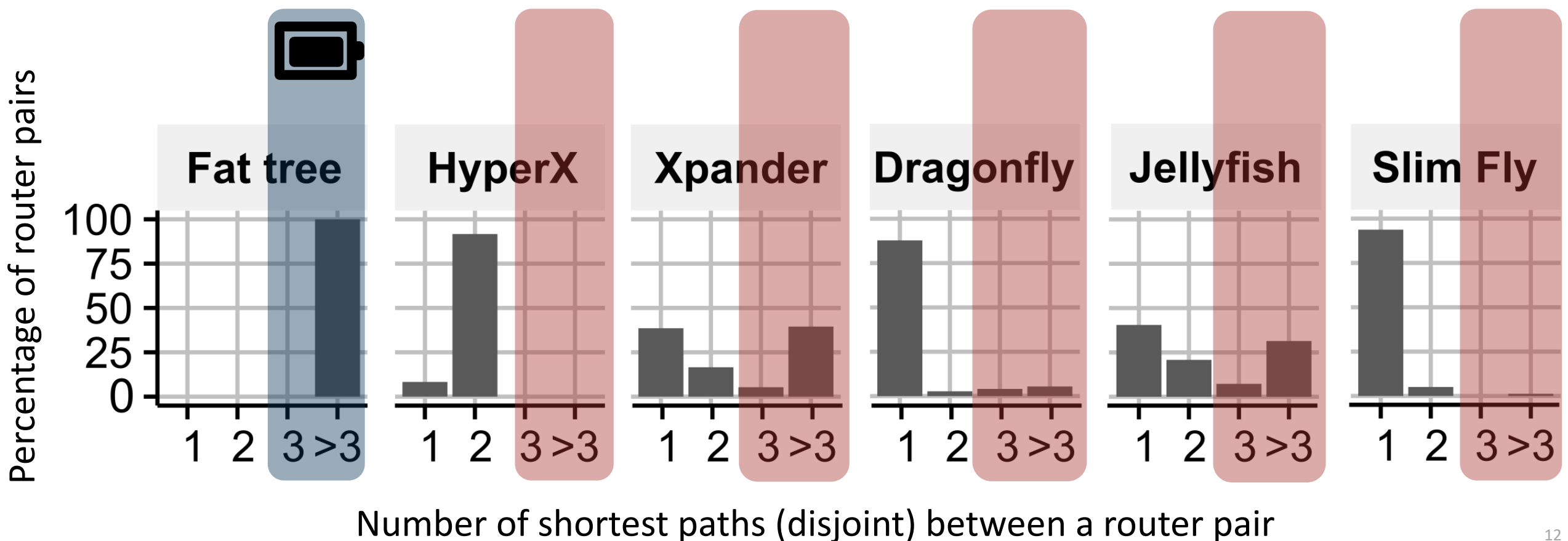
≈10,000 servers;  
comparable cost



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**  
**Shortest paths**

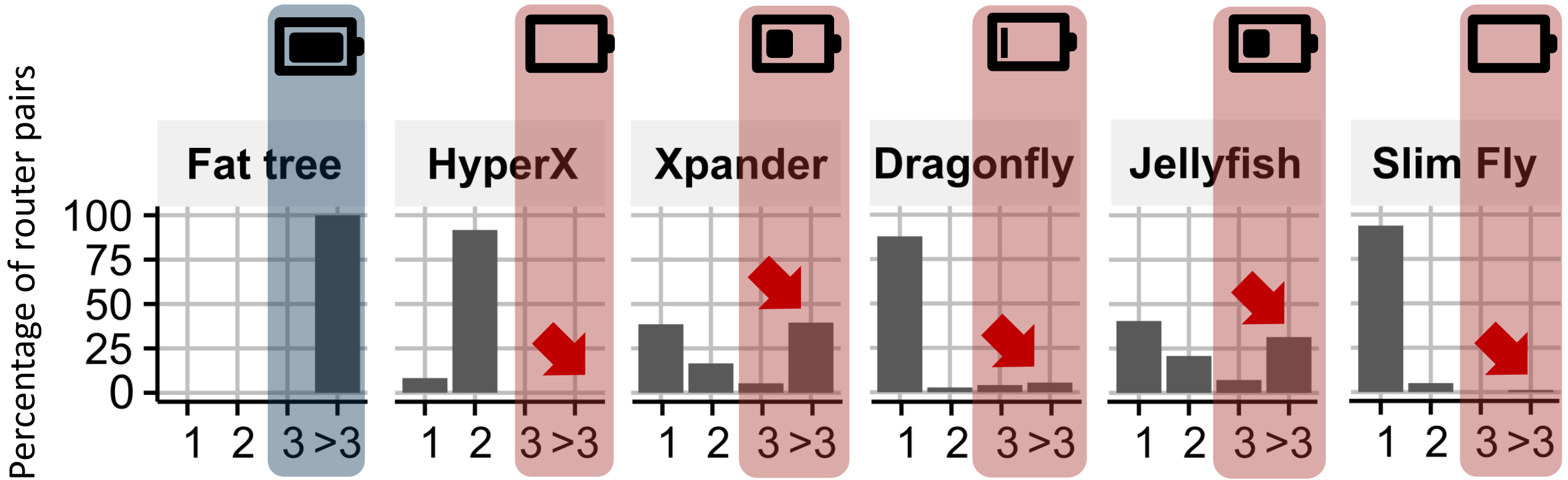
≈10,000 servers;  
 comparable cost



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

≈10,000 servers;  
 comparable cost

**Part 2**  
**Shortest paths**



Number of shortest paths (disjoint) between a router pair

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

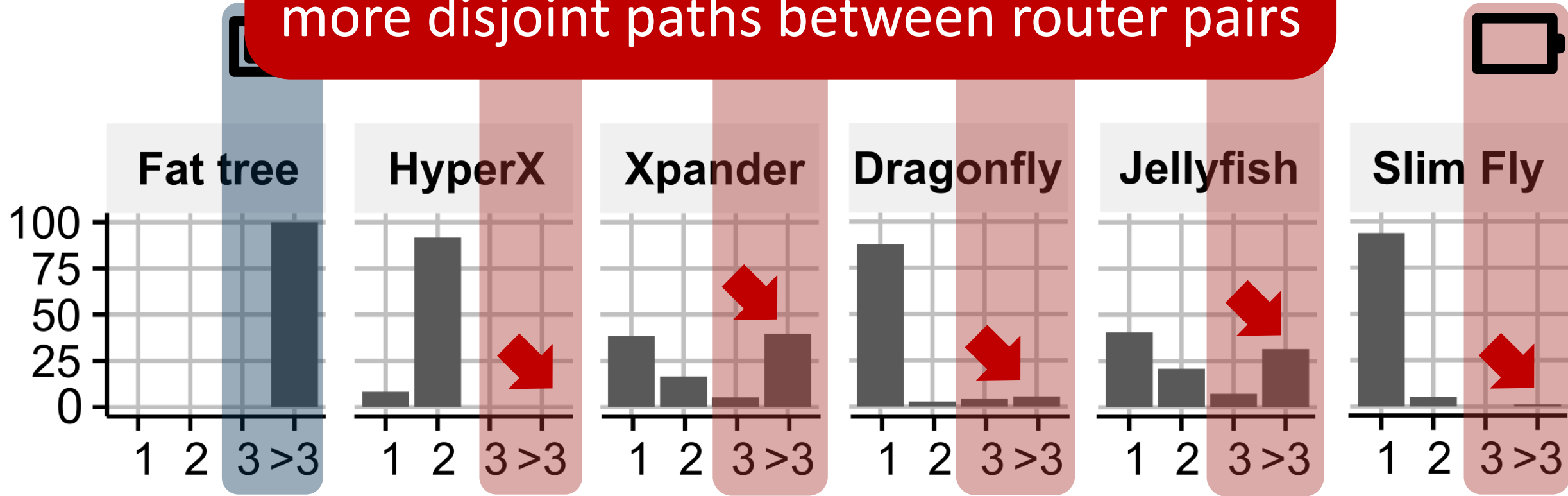
**Part 2**

**Shortest paths**

≈10,000 servers;  
comparable cost

**Key observation:** In most cases, there is not enough path diversity for three or more disjoint paths between router pairs

Percentage of router pairs



Number of shortest paths (disjoint) between a router pair

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non- shortest  
paths**

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

Part 2

Non-shortest  
paths

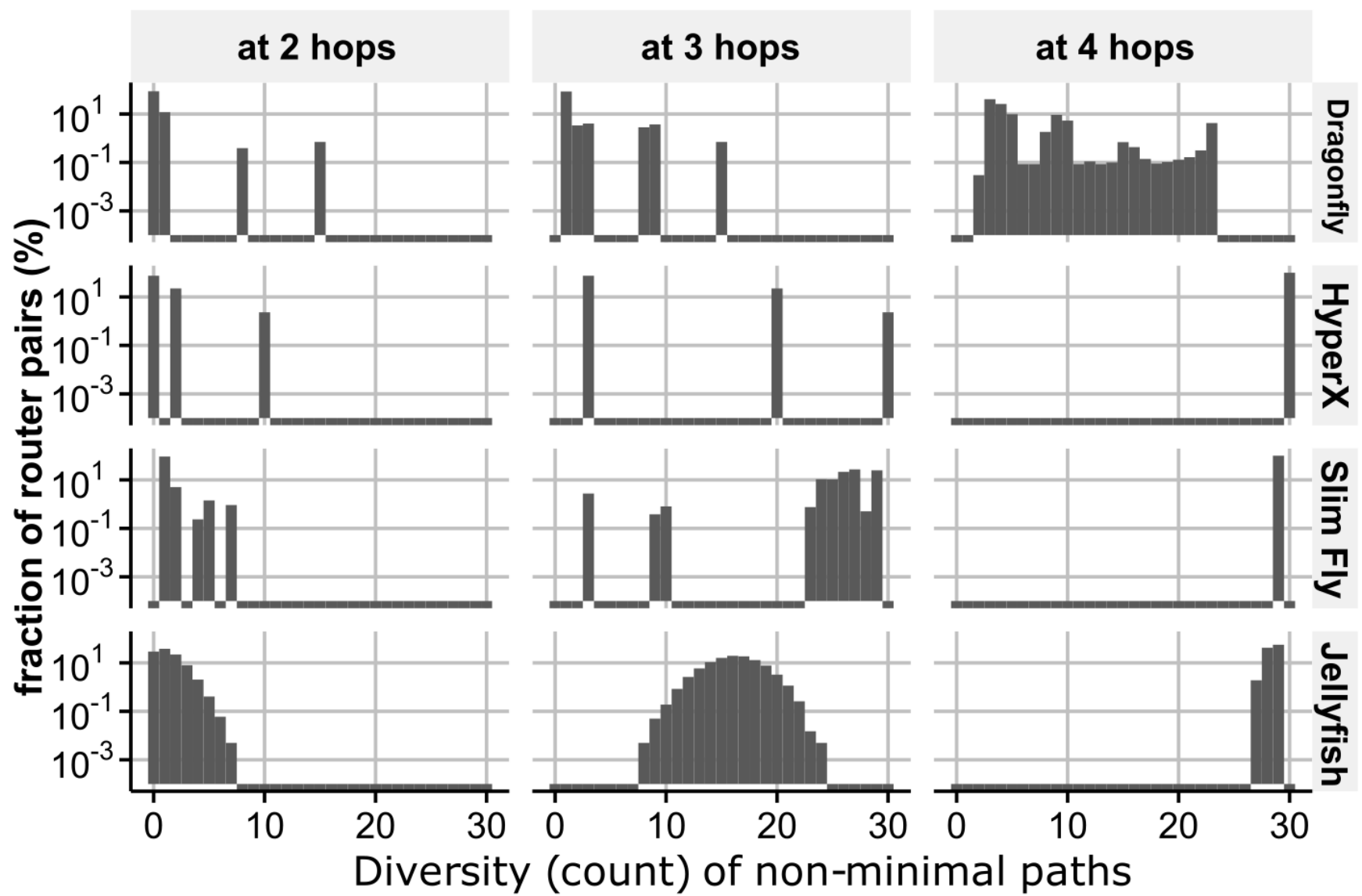
How about  
non-shortest  
paths?



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**



How about non-shortest paths?



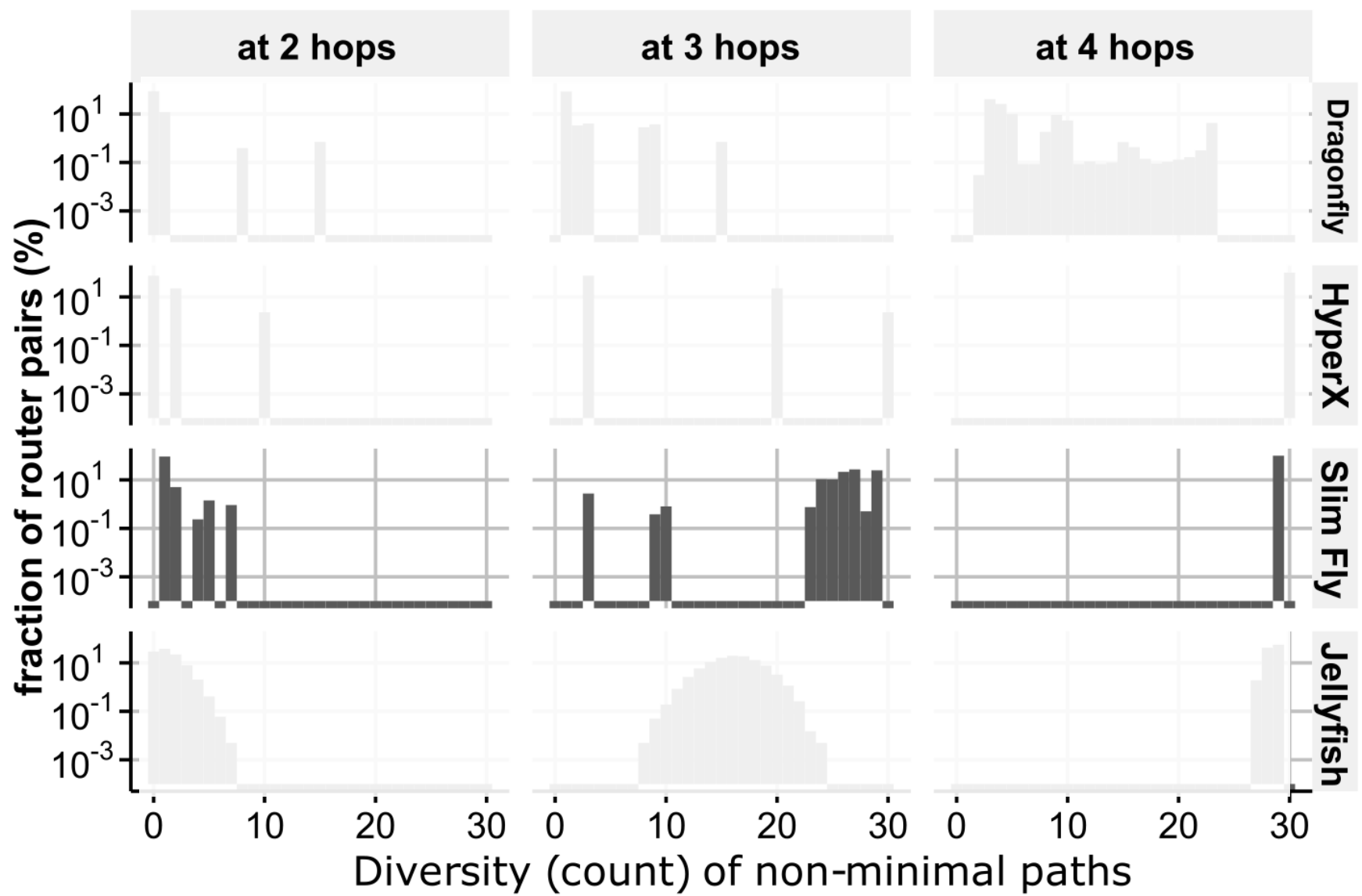


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non- shortest paths**

How about non-shortest paths?



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**

How about non-shortest paths?

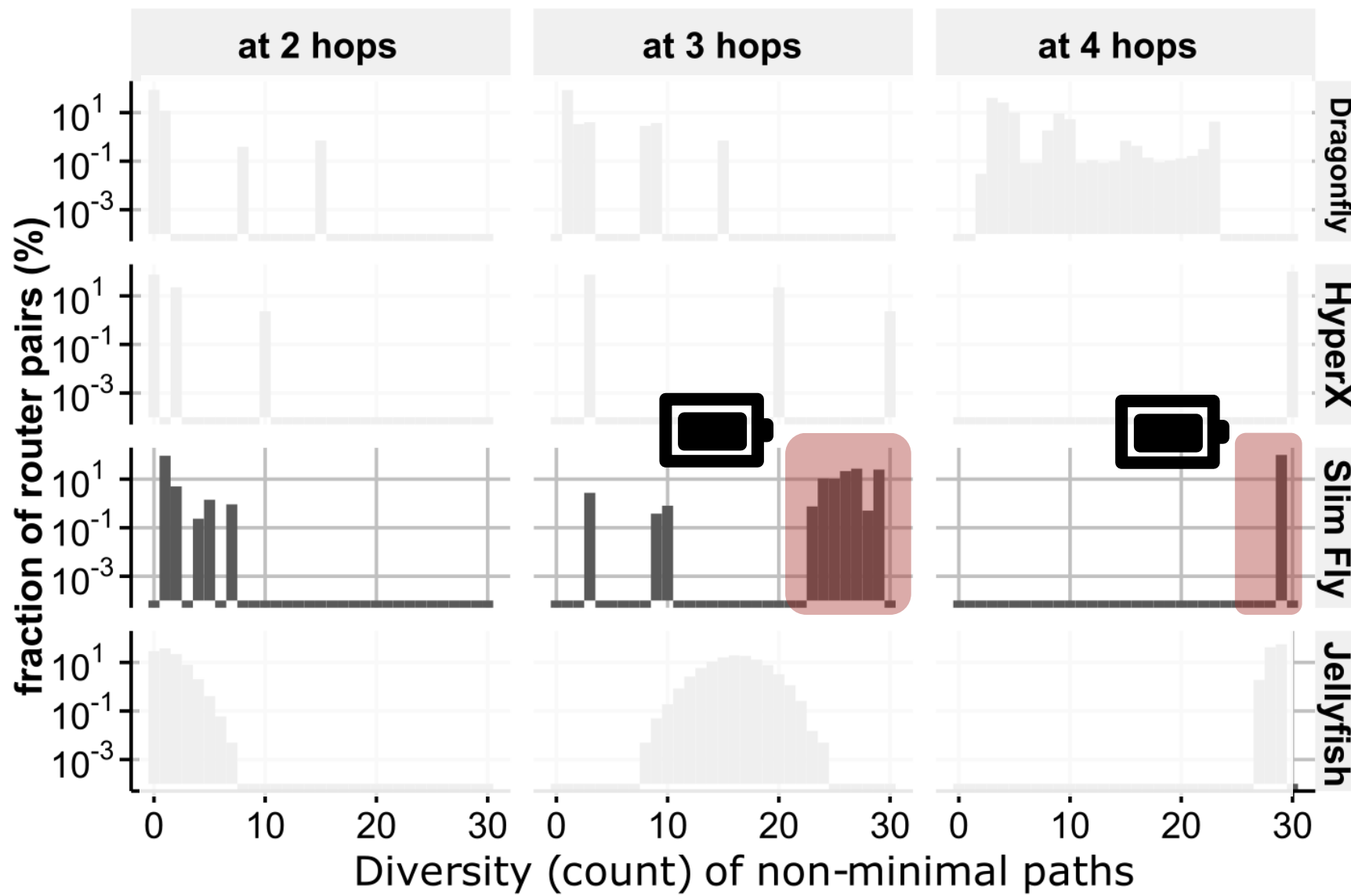


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**

How about non-shortest paths?

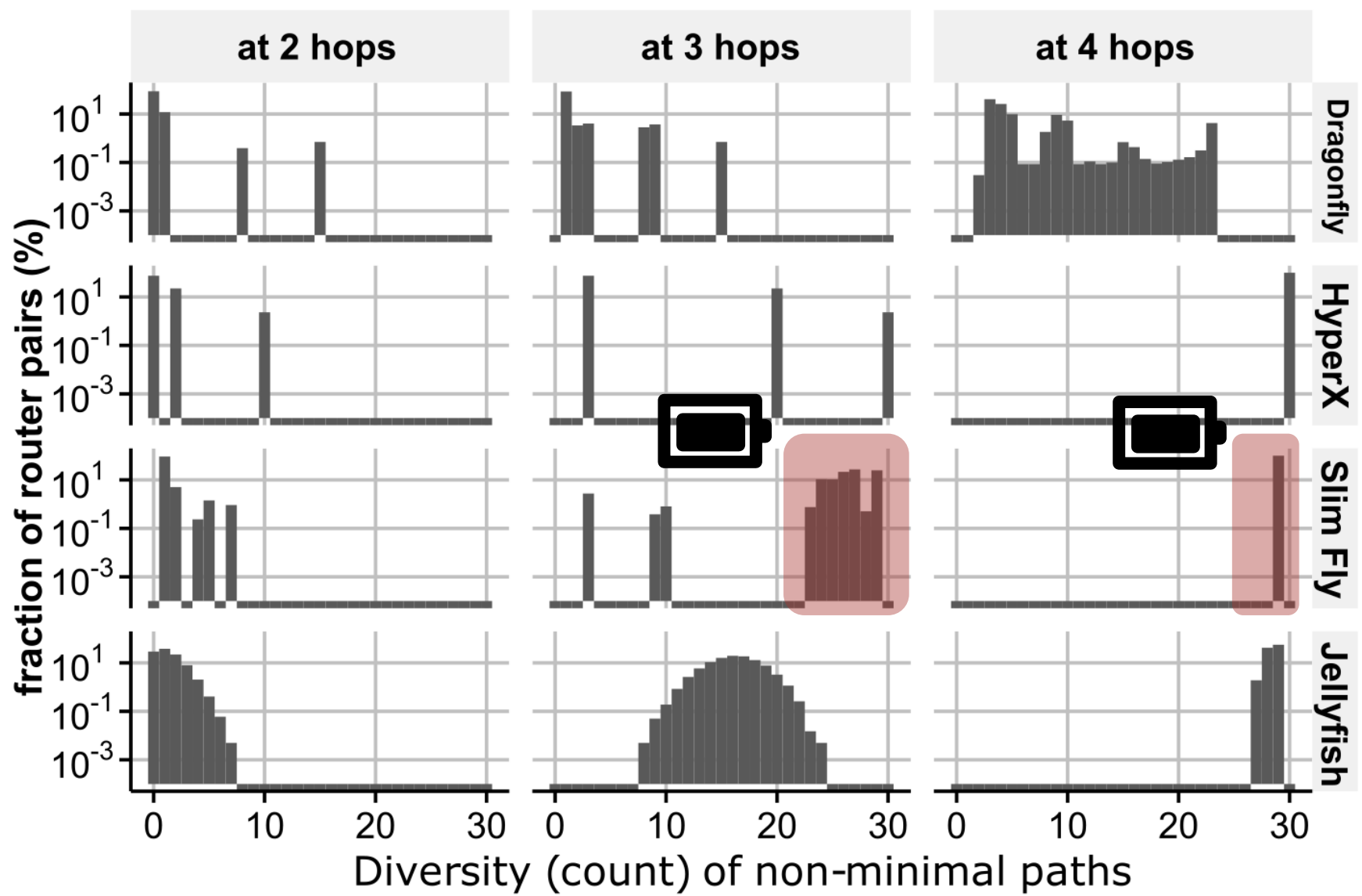


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**

How about non-shortest paths?



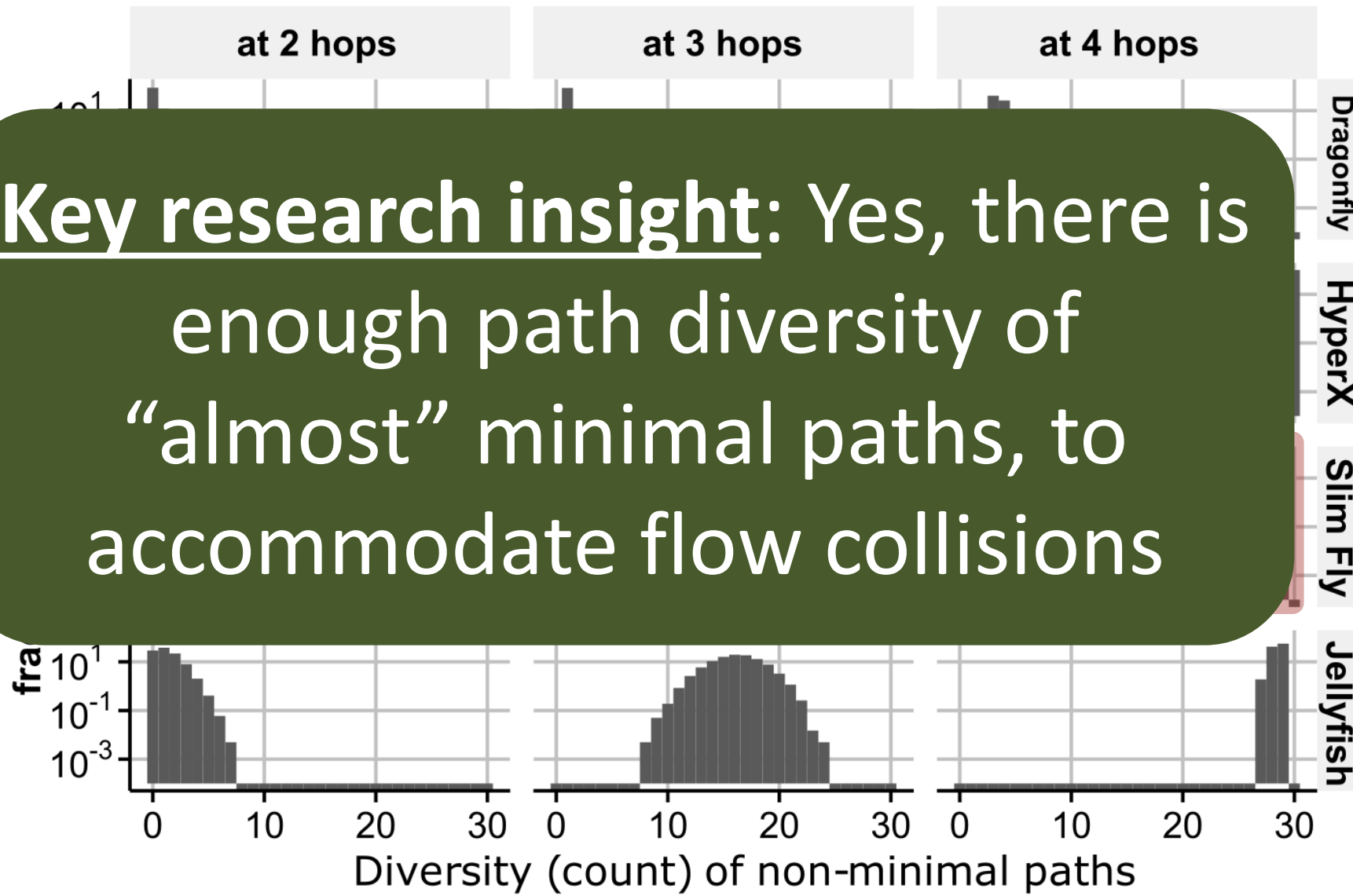
# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

Part 2

Non-shortest paths

Key research insight: Yes, there is enough path diversity of “almost” minimal paths, to accommodate flow collisions

How about non-shortest paths?



# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

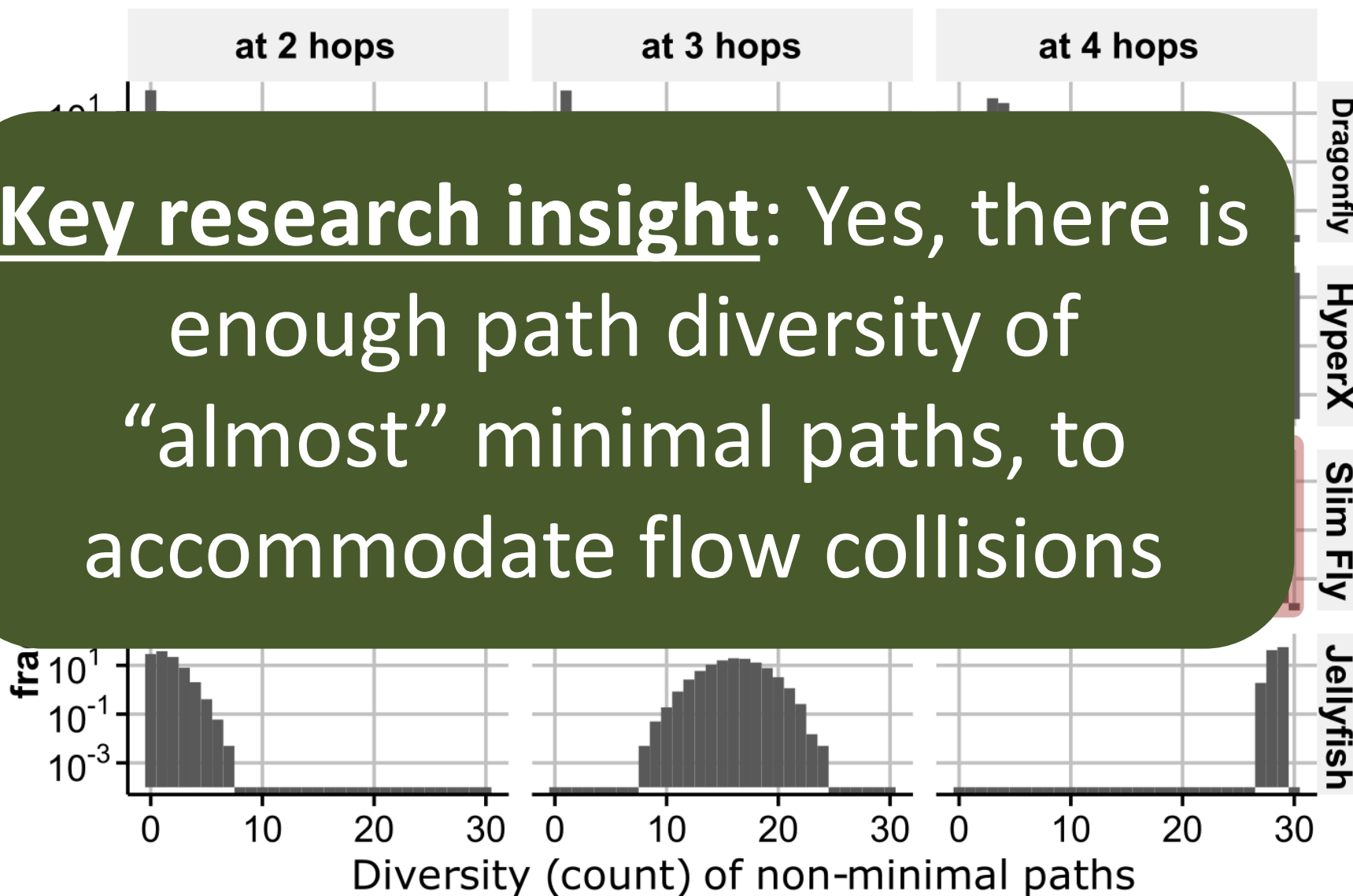
Part 2

Non-shortest paths

How about non-shortest paths?

Let's sum up...

Key research insight: Yes, there is enough path diversity of "almost" minimal paths, to accommodate flow collisions



# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

How many paths (in the network) do we need to „accommodate” collisions?

Multipathing

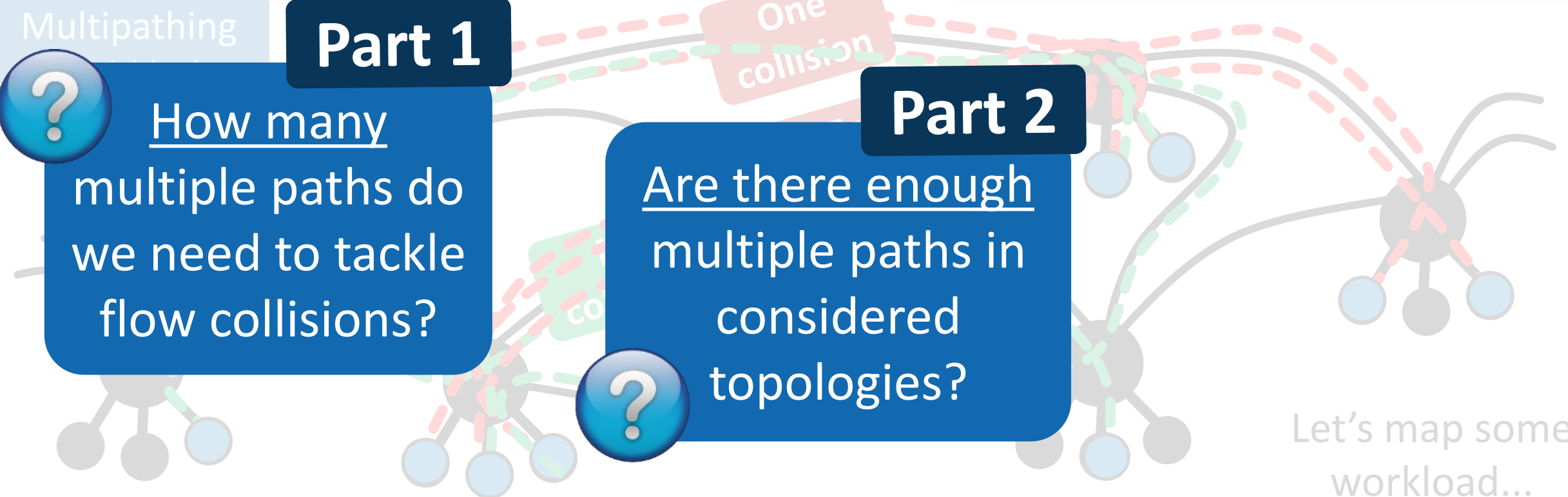
## Part 1

How many multiple paths do we need to tackle flow collisions?

One collision

## Part 2

Are there enough multiple paths in considered topologies?



Let's map some workload...

Multipath MOTIVATION



✗ Flows collide!

Three per router pair

How many paths do we need to tackle with multipathing?

? How many paths (in the network) do we need to „accommodate” collisions?

Multipathing

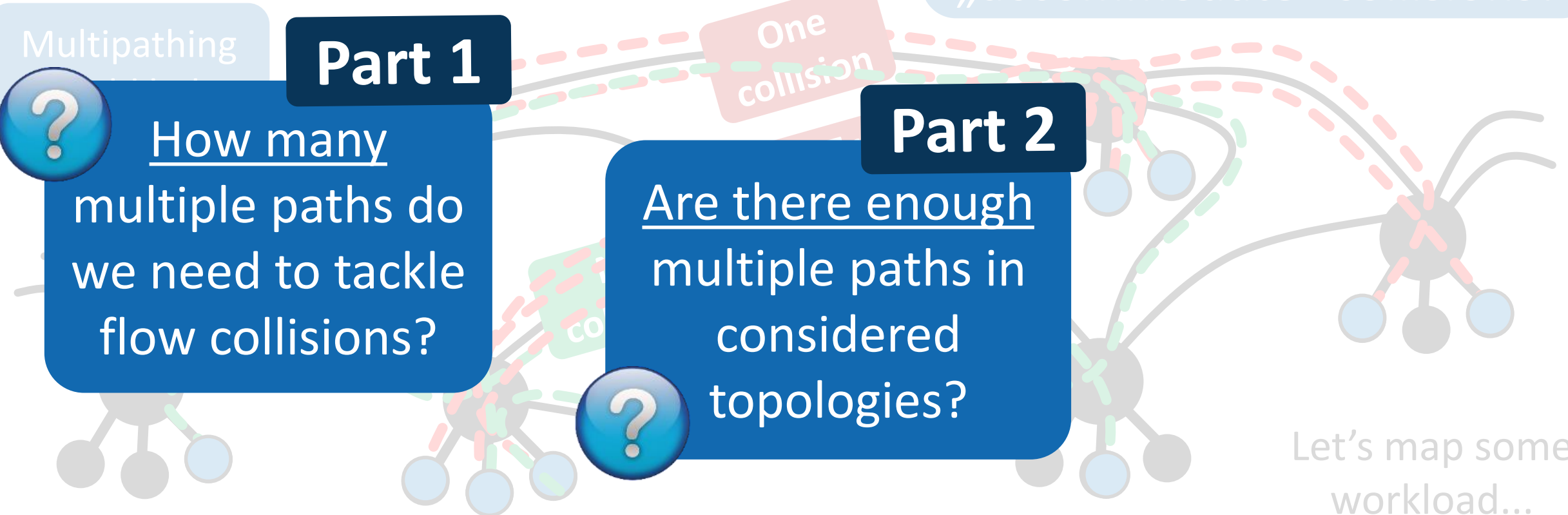
Part 1

? How many multiple paths do we need to tackle flow collisions?

One collision

Part 2

? Are there enough multiple paths in considered topologies?



Let's map some workload...



MULTIPATHING MOTIVATION

✗ Flows collide!

Three per router pair

Yes

When there are multiple paths between two routers, how many paths do we need to tackle with multipathing?

How many paths (in the network) do we need to „accommodate” collisions?

Multipathing

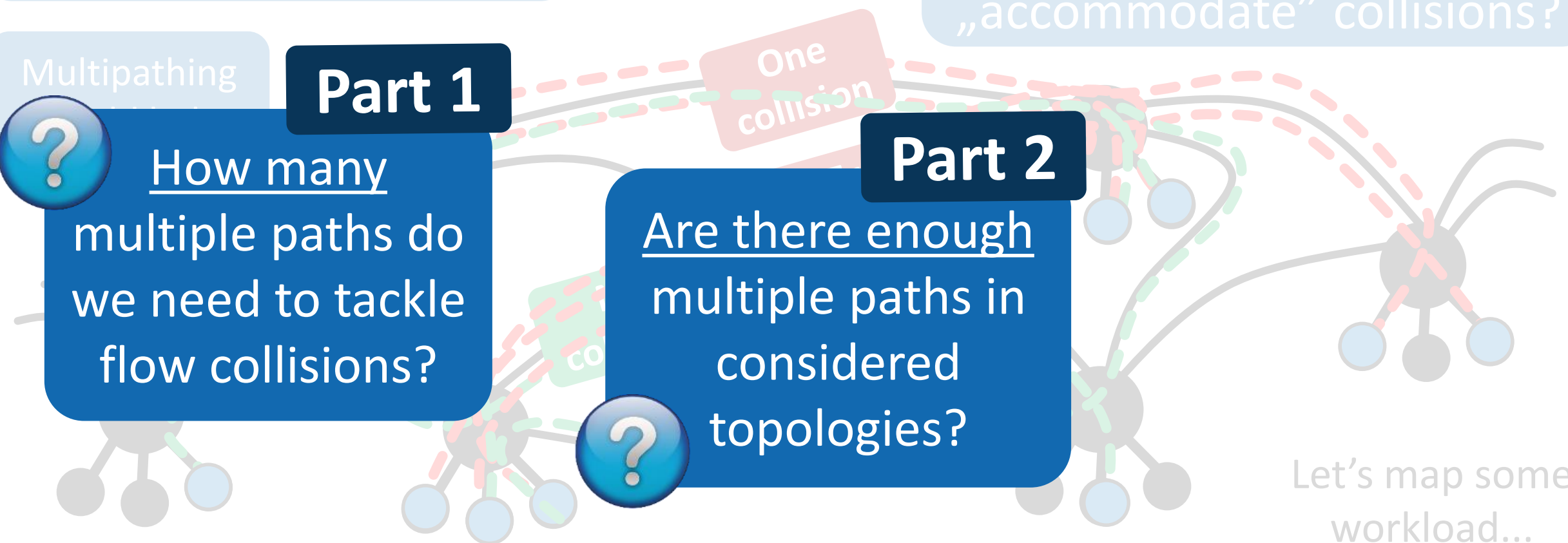
Part 1

? How many multiple paths do we need to tackle flow collisions?

One collision

Part 2

? Are there enough multiple paths in considered topologies?



Let's map some workload... 14

MULTIPATHING MOTIVATION

✗ Flows collide!

Three per router pair

Yes

Why do we need more than three paths to tackle flows with multipathing?

How many paths (in the network) do we need to „accommodate“ collisions?

Multipathing

Part 1

? How many multiple paths do we need to tackle flow collisions?

One collision

Part 2

? Are there enough multiple paths in considered topologies?

Part 3

? How to use such multiple paths in considered topologies?

Three per  
router pair

Yes

Let's develop a  
multipathing-  
focused routing  
architecture using  
these insights

Part 1

? How many  
multiple paths do  
we need to tackle  
flow collisions?

Part 2

? Are there enough  
multiple paths in  
considered  
topologies?

Part 3

? How to use  
such multiple  
paths in  
considered  
topologies?

One  
collision

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

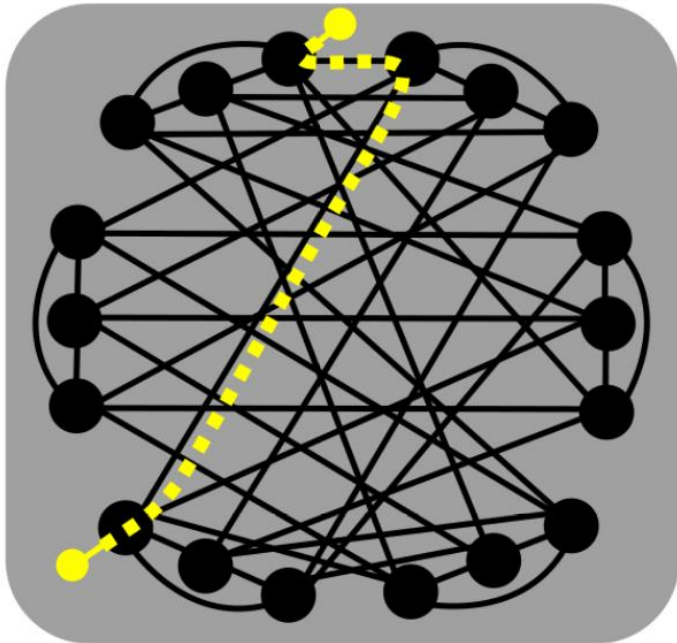
## Part 3.1

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

**Default topology**

Shortest path: 2 hops



**Default topology**

Shortest path: 2 hops

**Part 3.1**

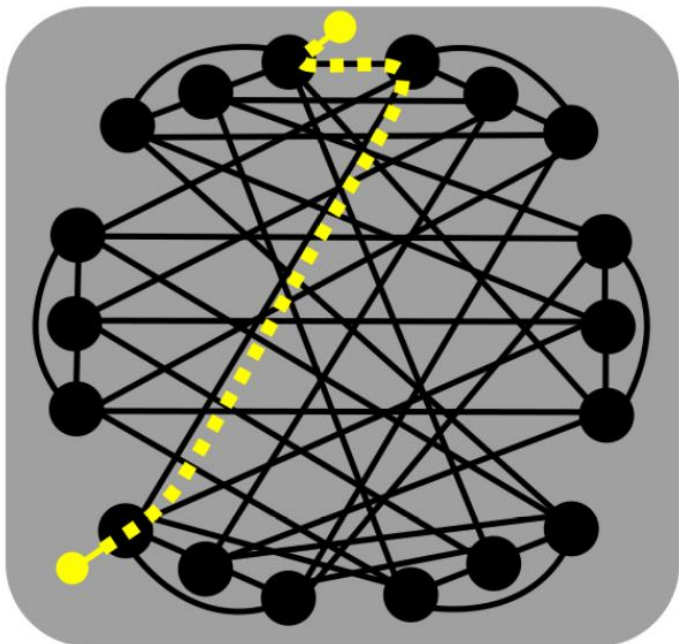
# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

(a.1) Divide links into subsets (layers).

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

## Default topology

Shortest path: 2 hops



## Default topology

Shortest path: 2 hops

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

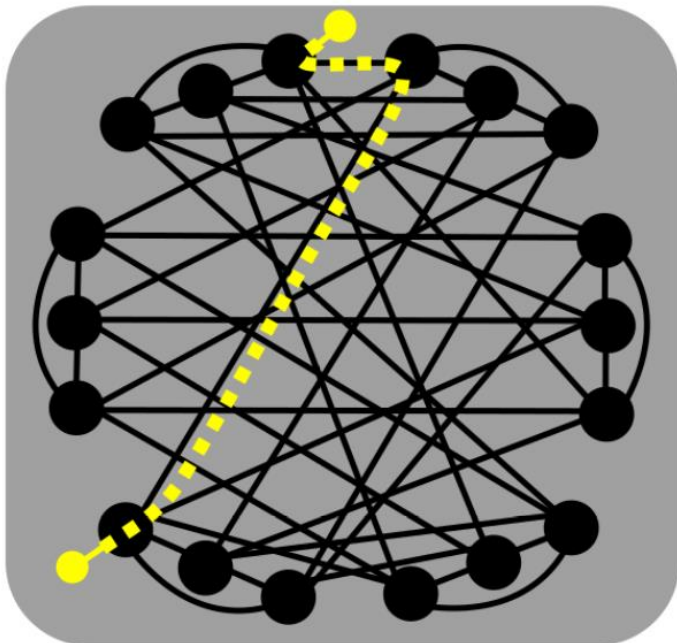
## Part 3.1

(a.1) Divide links into subsets (layers).

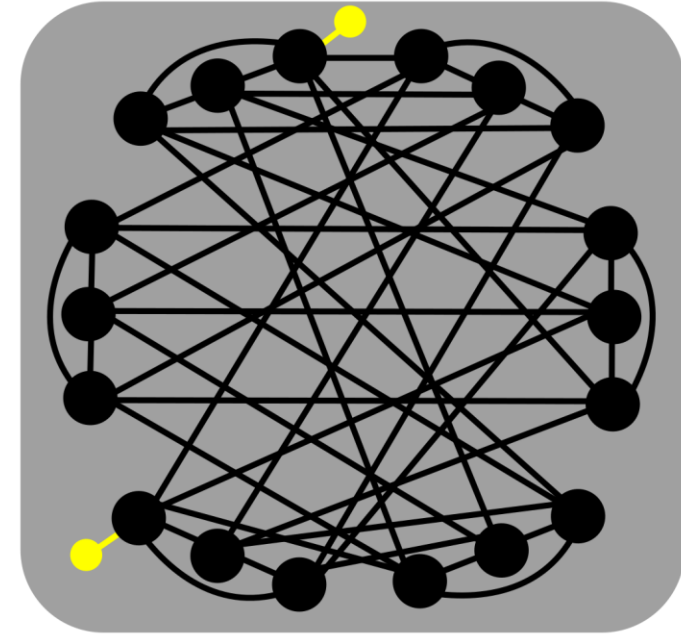
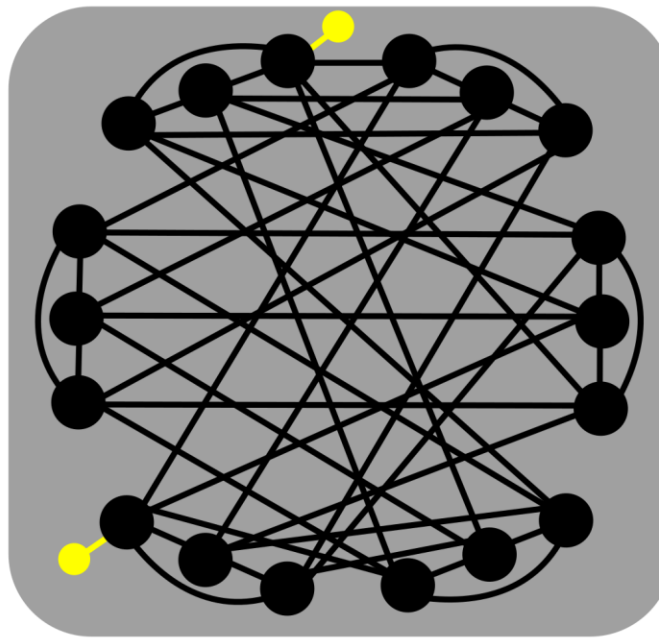
(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

**Default topology**

Shortest path: 2 hops



**Default topology**



# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

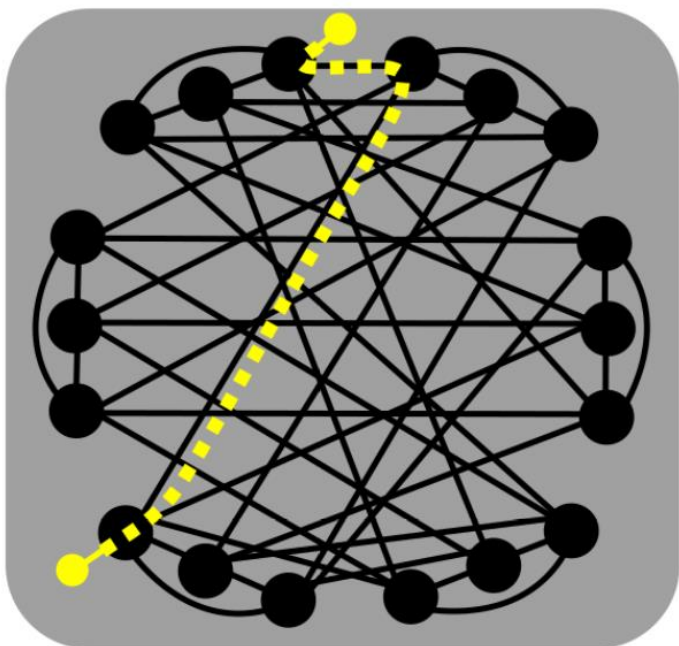
## Part 3.1

(a.1) Divide links into subsets (layers).

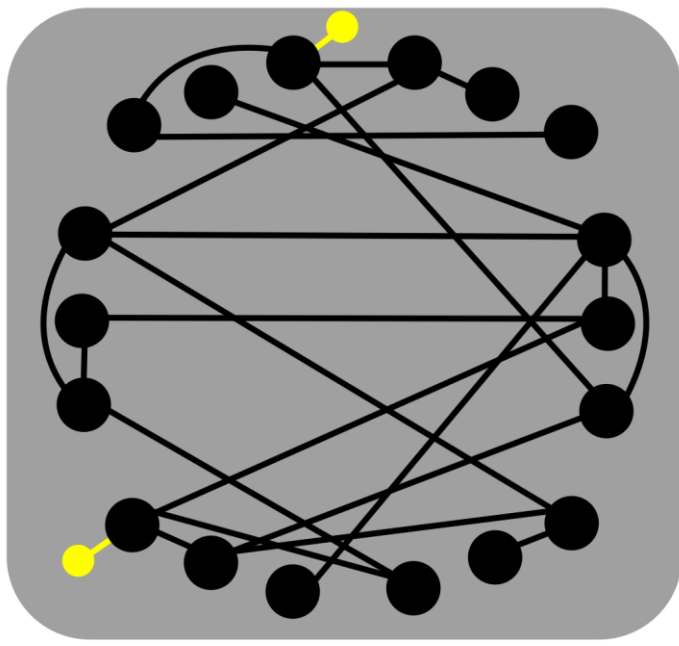
(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

**Default topology**

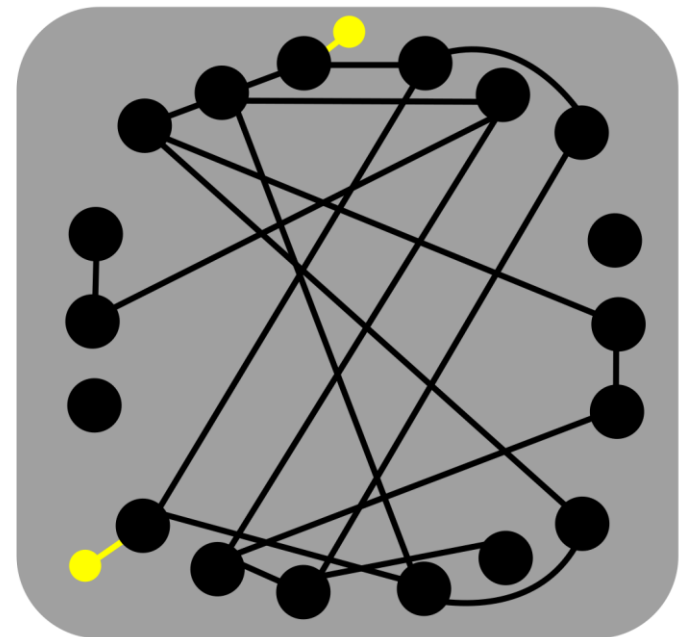
Shortest path: 2 hops



**Layer 1: "Almost"-Default topology**  
-shortest path: 5 hops



**Layer 2: "Almost"-shortest path: 3 hops**



...



# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

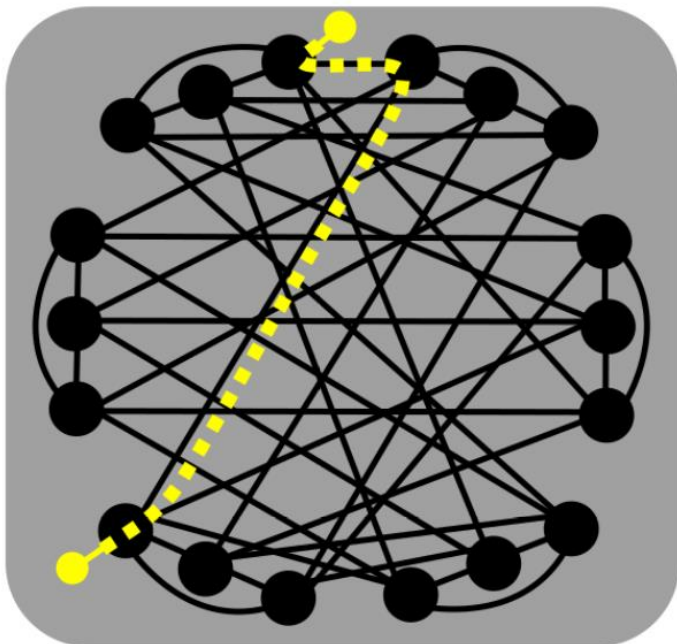
(a.1) Divide links into subsets (layers).

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

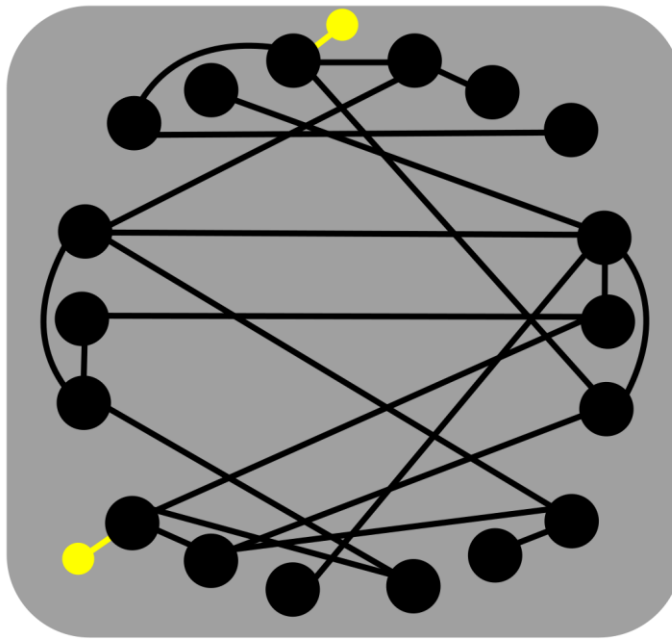
(b) Divide one flow into subflows and send subflows across different layers

**Default topology**

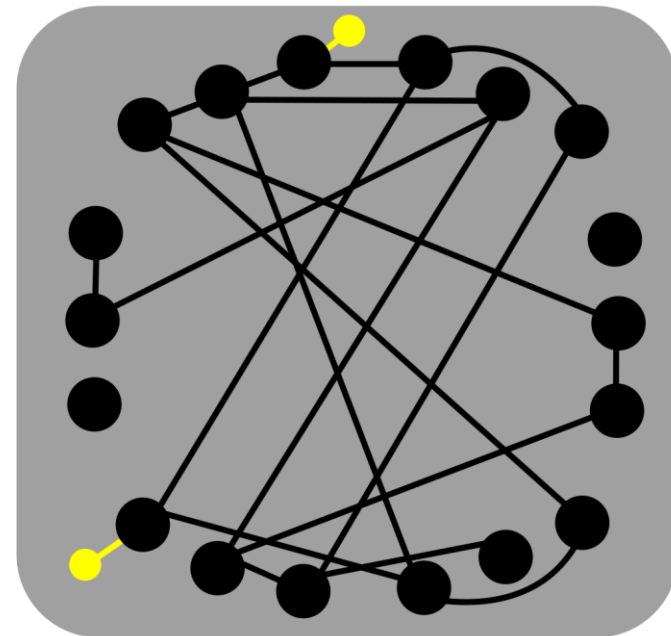
Shortest path: 2 hops



**Layer 1: "Almost"-  
Default topology**  
-shortest path: 5 hops



**Layer 2: "Almost"-  
-shortest path: 3 hops**



...

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

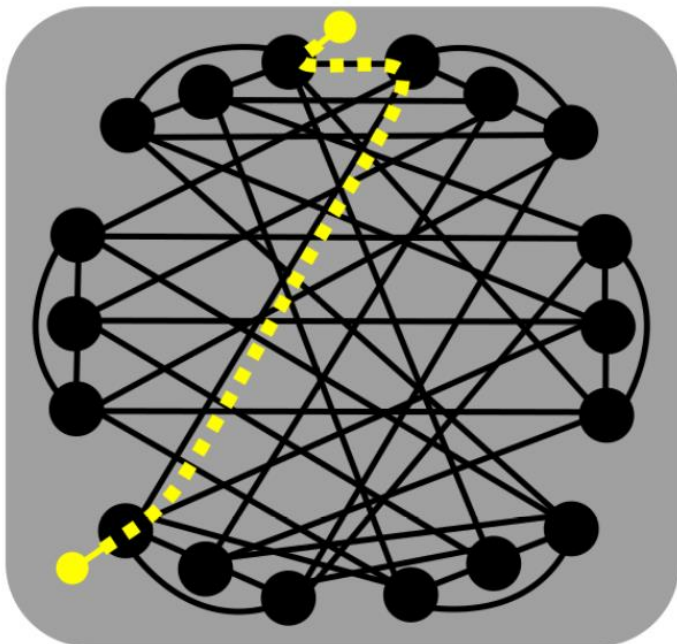
(a.1) Divide links into subsets (layers).

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

(b) Divide one flow into subflows and send subflows across different layers

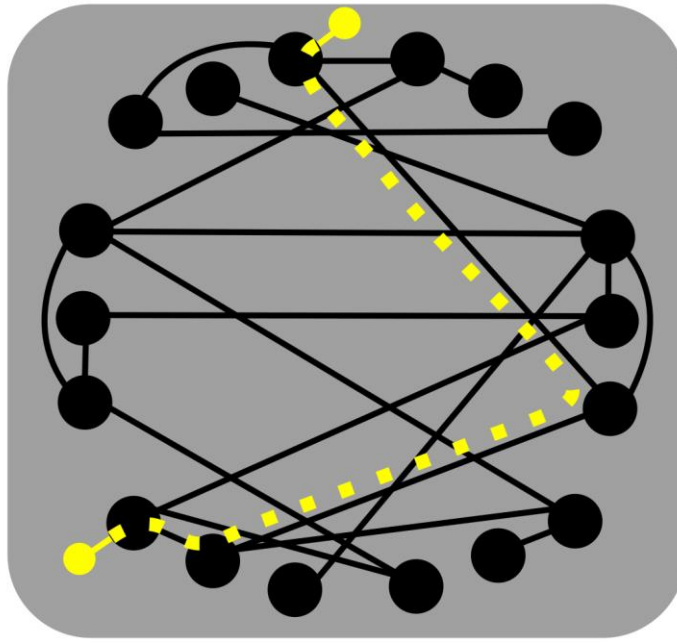
**Default topology**

Shortest path: 2 hops

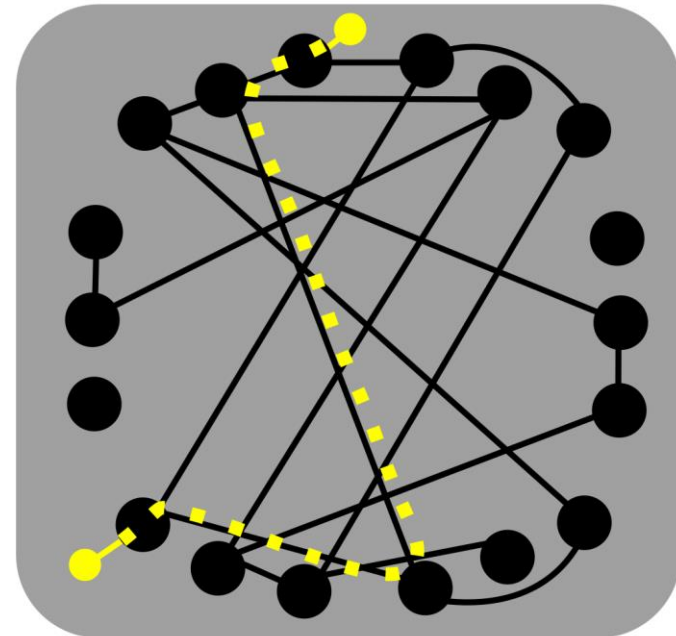


**Layer 1: "Almost"-  
Default topology**

-shortest path: 3 hops



**Layer 2: "Almost"-  
-shortest path: 3 hops**



...

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

(a.1) Divide links into subsets (layers).

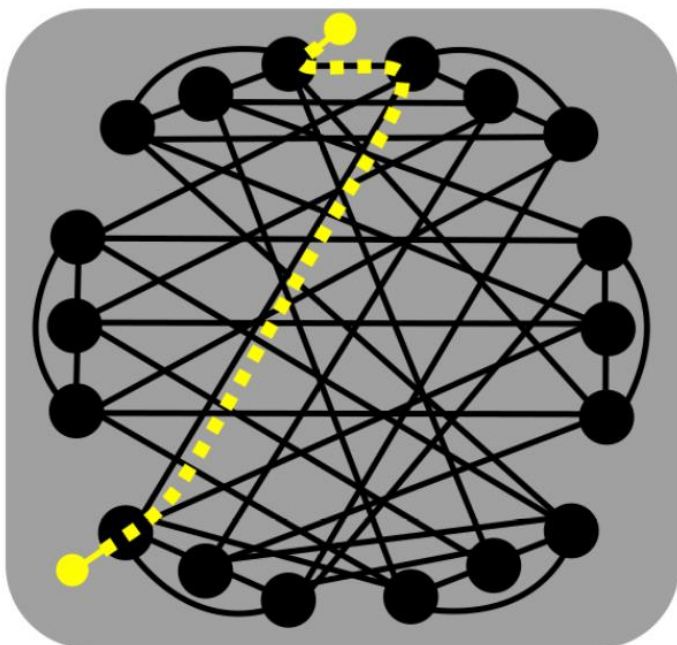
(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

(b) Divide one flow into subflows and send subflows across different layers

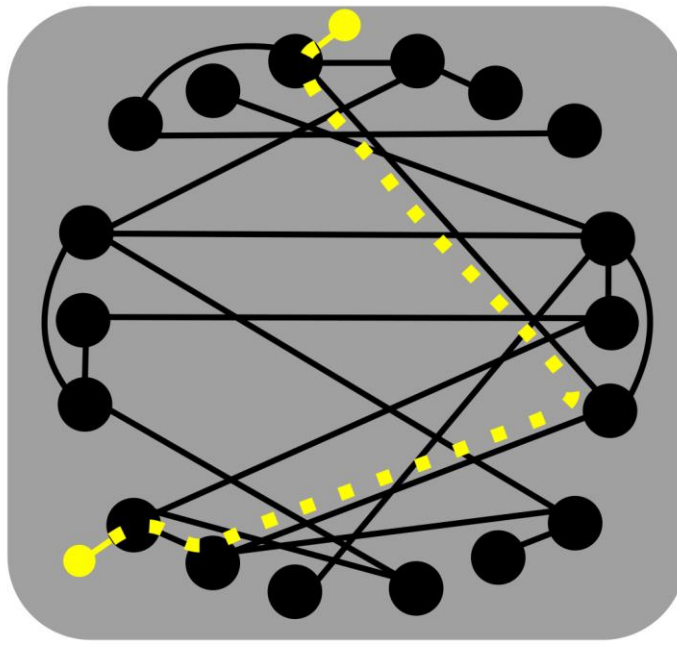
(c) Route minimally in each layer. A minimal route in one layer is often non-minimal when considering all links

**Default topology**

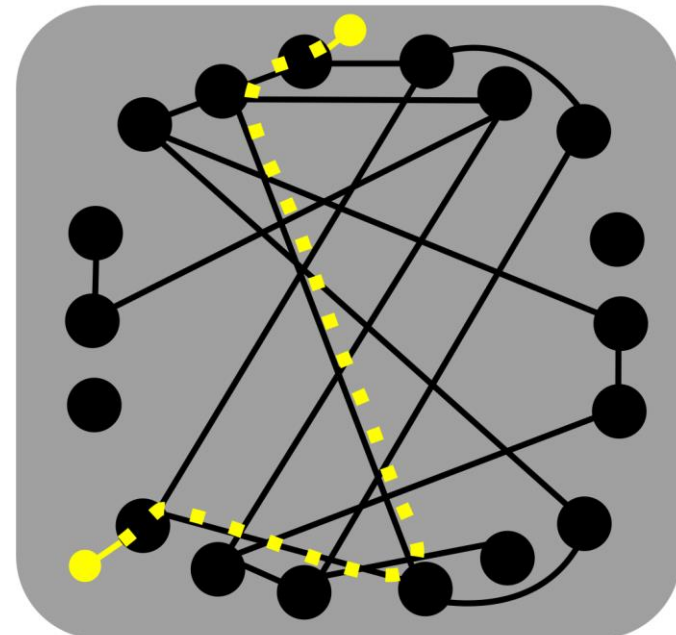
Shortest path: 2 hops



**Layer 1: "Almost"-  
Default topology**  
-shortest path: 5 hops



**Layer 2: "Almost"-  
-shortest path: 3 hops**



...

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

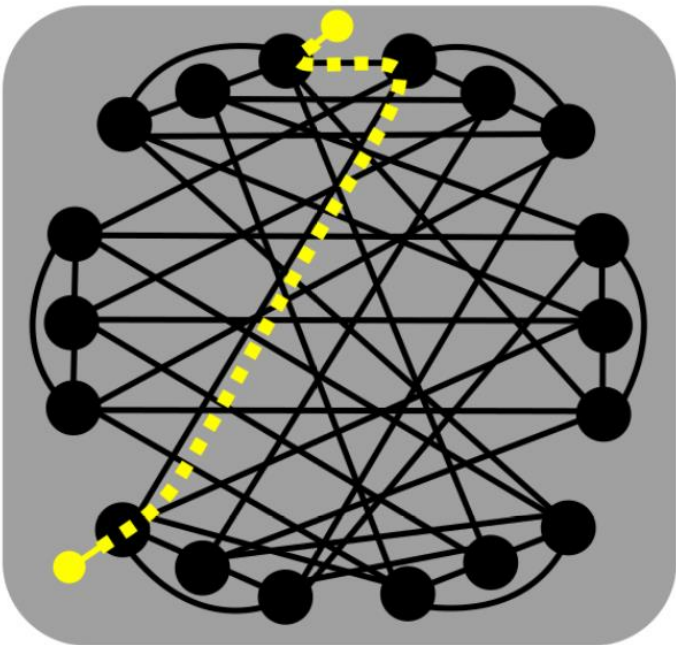
(a.1) Divide links into subsets (layers).

(a.2) ...

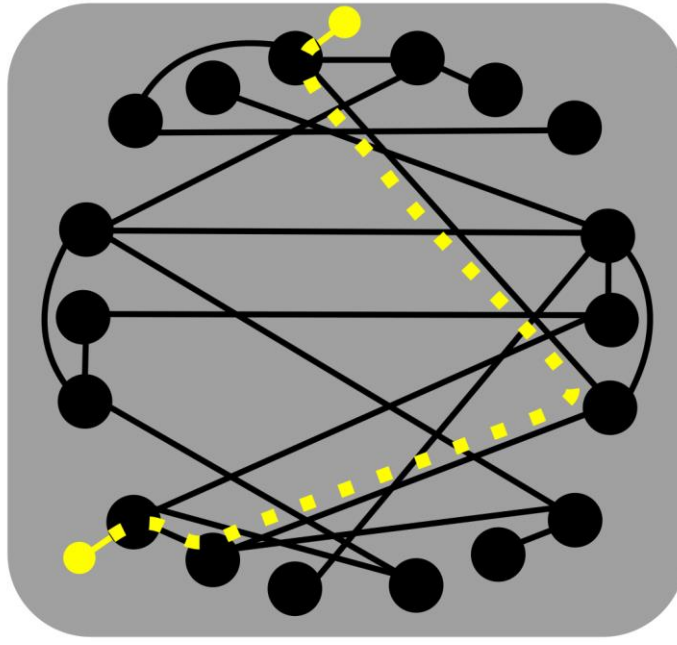
Design justification: Different layers enable multipathing

minimally in each ...  
 ... route in one ...  
 ... non-minimal ...  
 ... considering all links

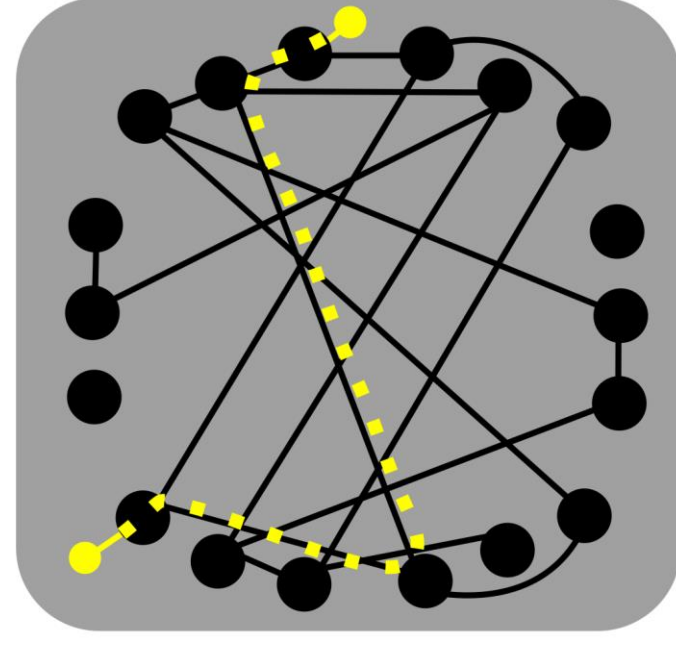
**Default topology**  
 Shortest path: 2 hops



**Layer 1: "Almost"-  
 Default topology**  
 -shortest path: 5 hops



**Layer 2: "Almost"-  
 -shortest path: 3 hops**



...

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

(a.1) Divide links into subsets (layers).

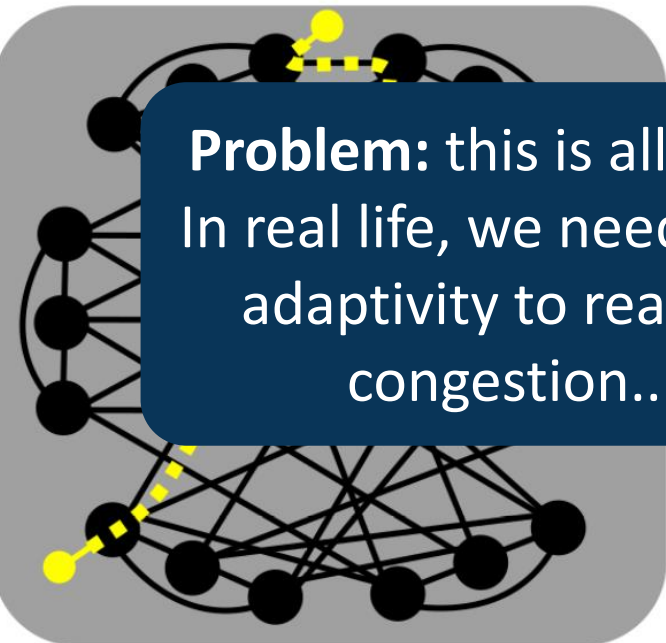
(a.2) ...

Design justification: Different layers enable multipathing

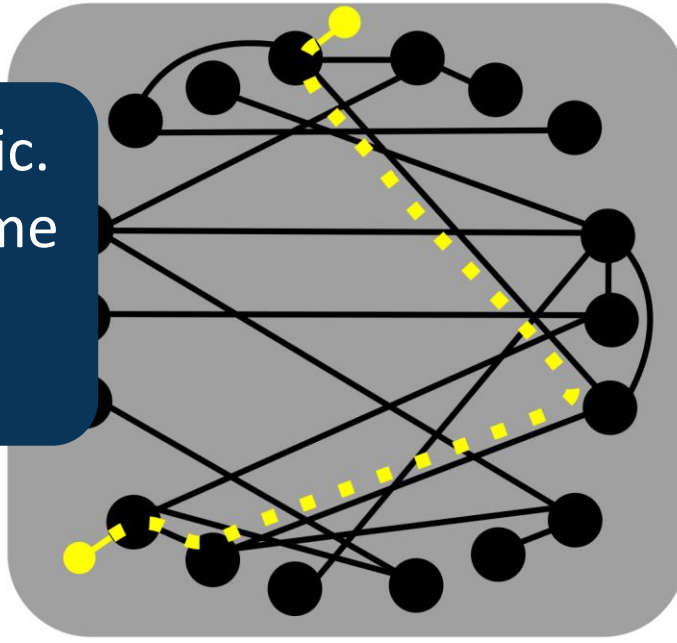
minimally in each ...  
 ... route in one ...  
 ... non-minimal ...  
 ... considering all links

**Default topology**

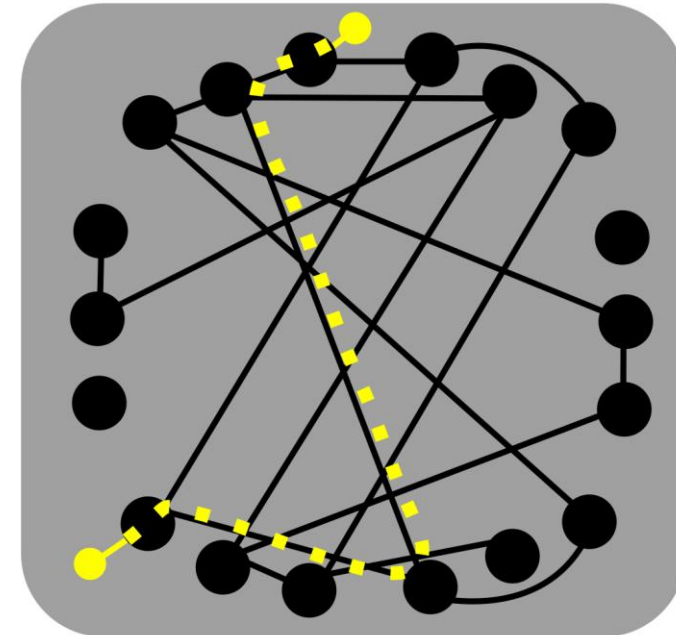
Shortest path: 2 hops



**Layer 1: "Almost"-shortest path: 3 hops**



**Layer 2: "Almost"-shortest path: 3 hops**



...

**Problem:** this is all static. In real life, we need some adaptivity to react to congestion...

# FATPATHS ARCHITECTURE: LAYERED ROUTING PROTOCOL

## Part 3.1

(a.1) Divide links into subsets (layers).

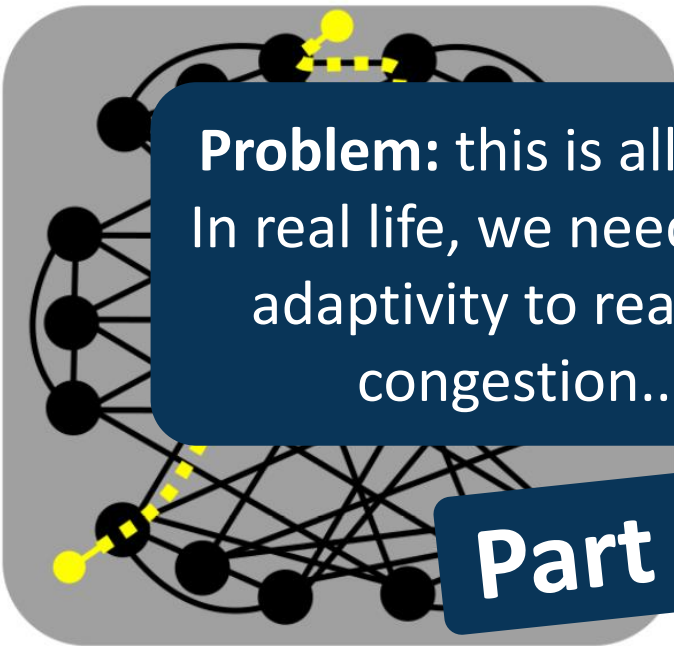
Design justification: Different layers enable multipathing

minimally in each  
normal route in one  
then non-minimal  
considering all links

**Default topology**  
Shortest path: 2 hops

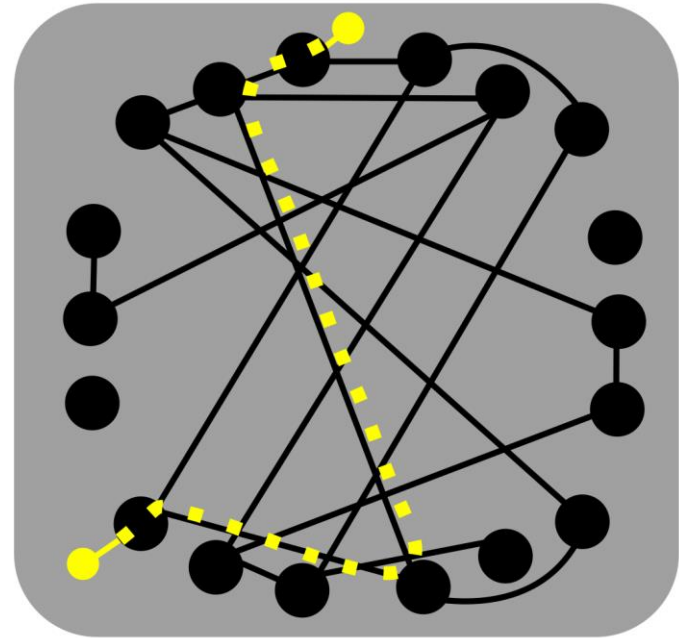
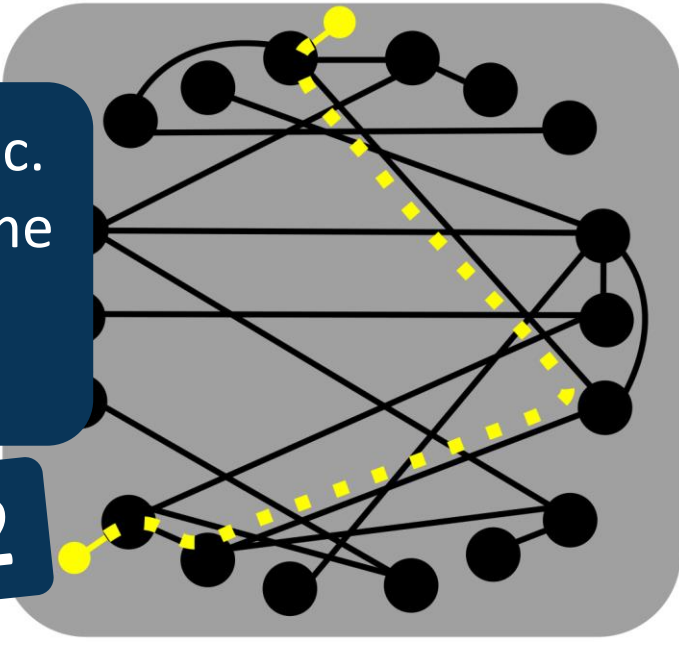
**Layer 1: "Almost"-  
Default topology**  
-shortest path: 5 hops

**Layer 2: "Almost"-  
-shortest path: 3 hops**



**Problem:** this is all static. In real life, we need some adaptivity to react to congestion...

## Part 3.2



...

# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

## FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

How to load balance traffic  
over different paths of  
different lengths?





## FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

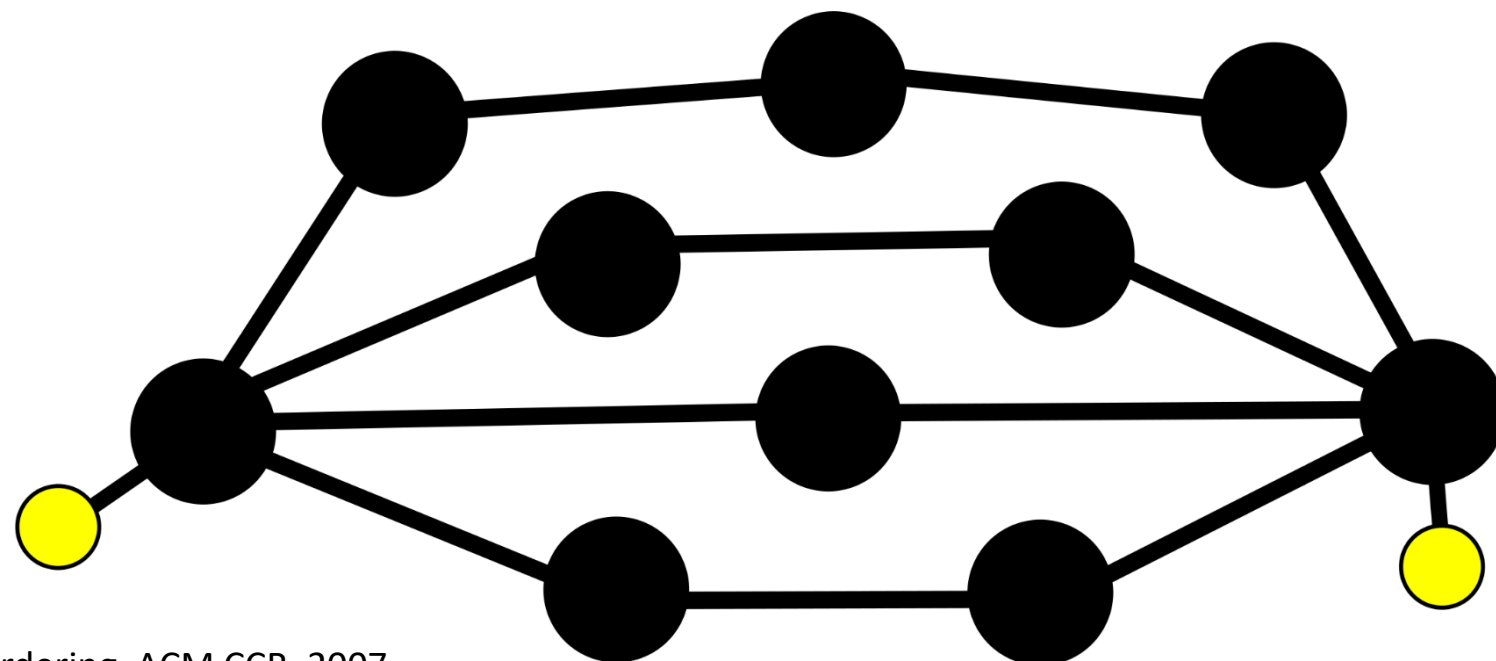
[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

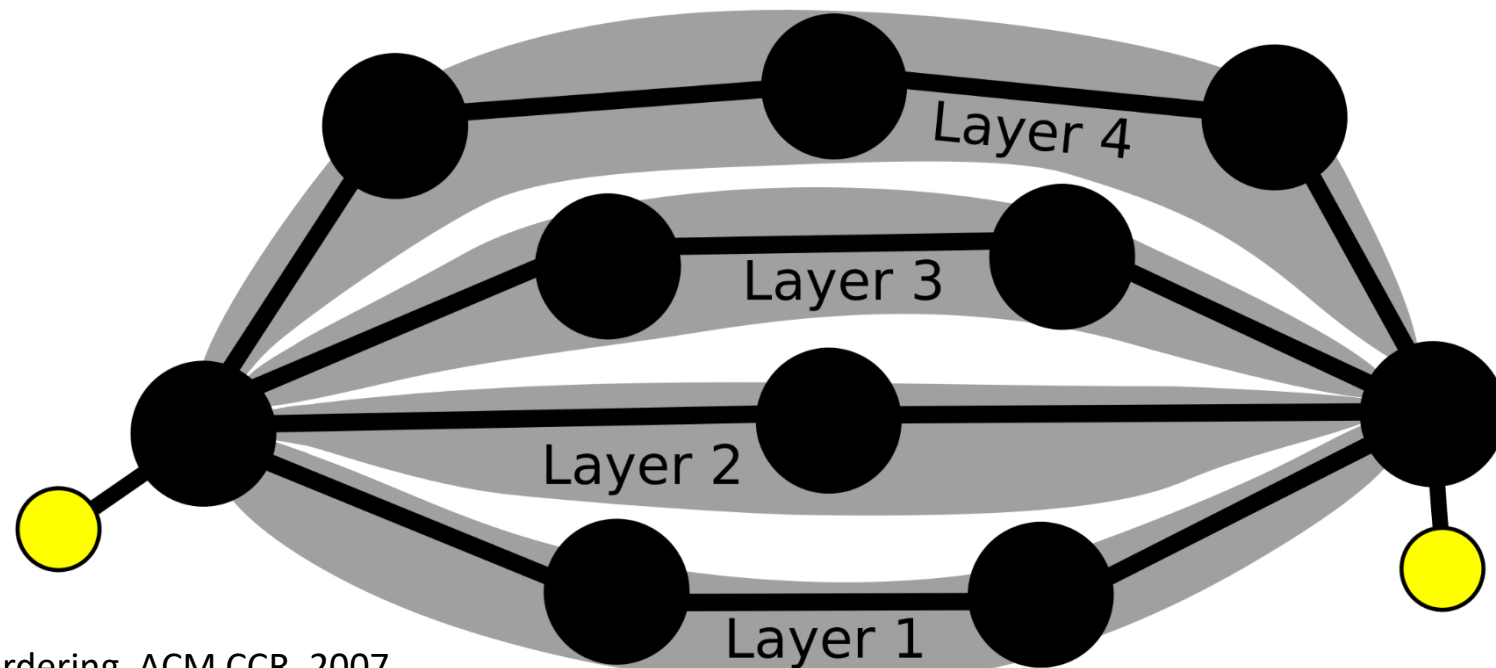
[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

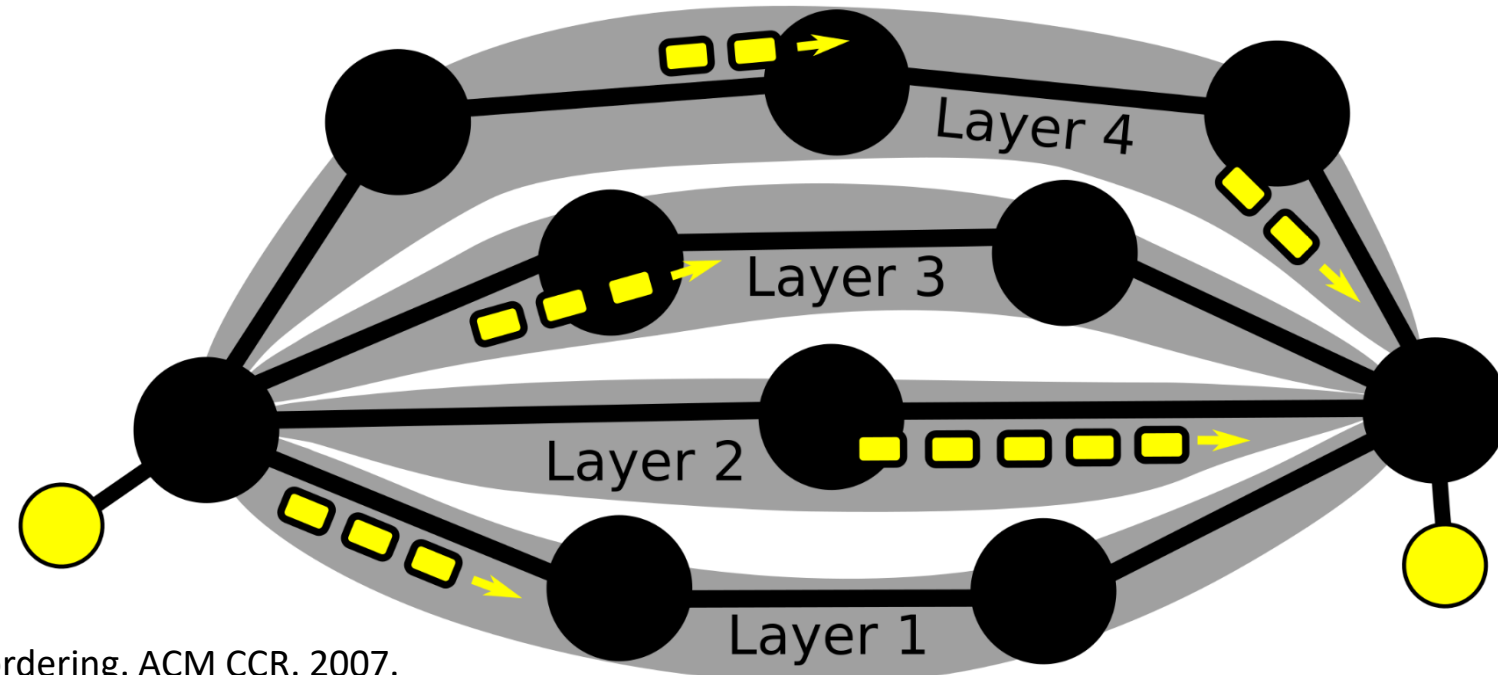
[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

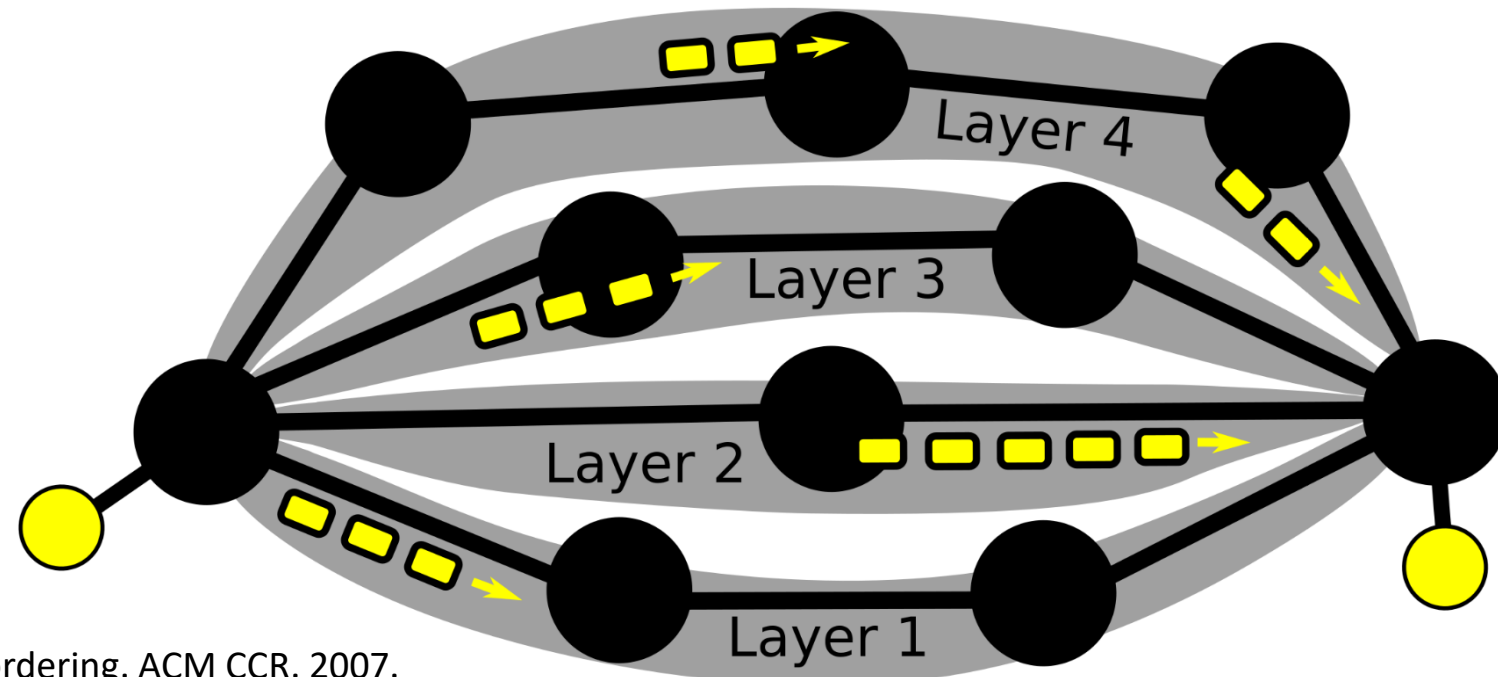
# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?

**Very simple load balancing:** when one detects congestion, a router just picks a random path (layer) for the packets to be sent (they become a new flowlet)



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

# FATPATHS ARCHITECTURE: LOAD BALANCING WITH FLOWLETS

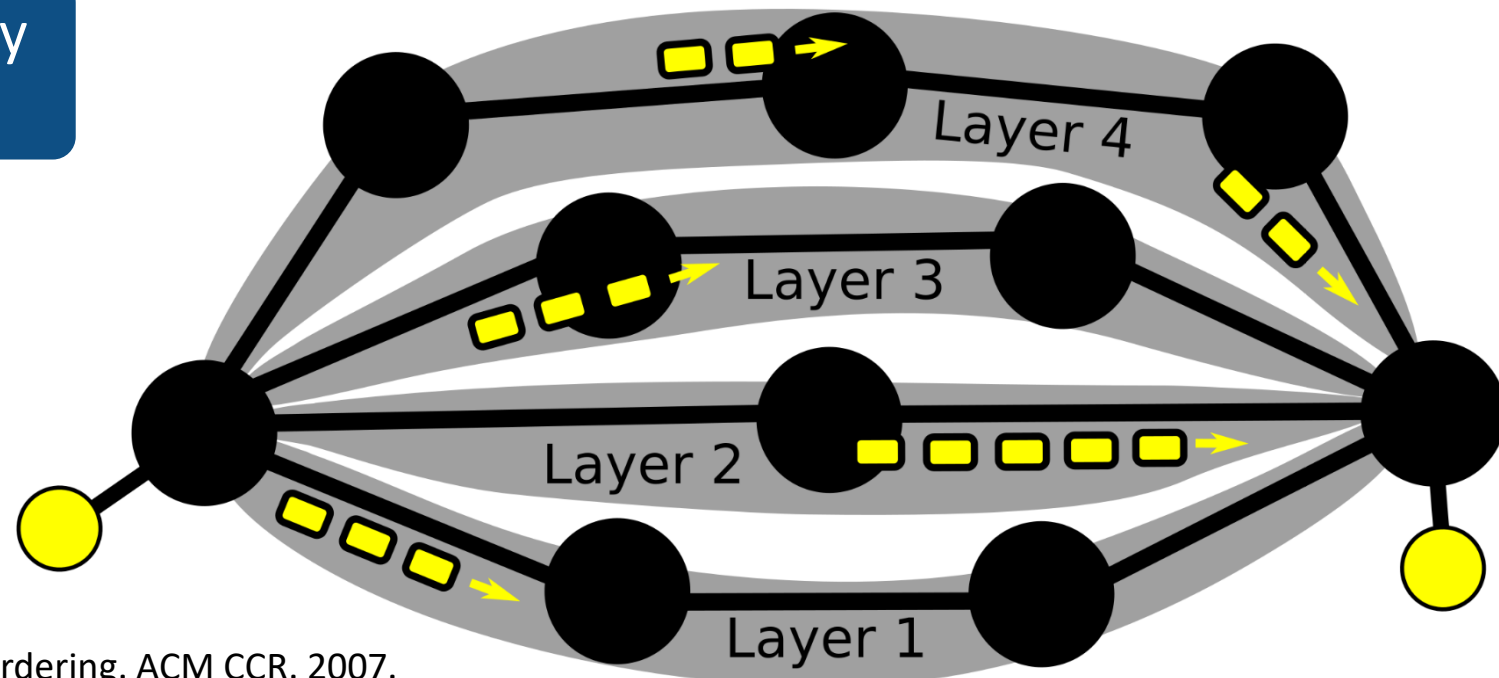
A flowlet [1,2] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.



How to load balance traffic over different paths of different lengths?

Size of flowlets changes automatically based on conditions in the network

**Very simple load balancing:** when one detects congestion, a router just picks a random path (layer) for the packets to be sent (they become a new flowlet)



[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR. 2007.

[2] E. Vanini et al. Let It Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

# FATPATHS ARCHITECTURE

Design details

# FATPATHS ARCHITECTURE

Design details

```

1 L = {E} //Init a set of layers L; we start with E that corresponds to  $\sigma_1$ .
2  $\Pi = \{\pi_1(V), \dots, \pi_{n-1}(V)\}$  //Generate n-1 random permutations of vertices.
3
4 //Init a matrix W containing weights of edges  $(u, v) \in E$ .
5  $W = \{[w_{uv}] \mid \forall u, v \in V : w_{uv} = 0\}$ 
6
7 foreach  $\pi \in \Pi$  do: //One iteration of the main loop derives one layer.
8    $E' = \text{create\_layer}(\pi, E, W, L_{min}, L_{max})$  //Generate a layer.
9    $L = L \cup \{E'\}$  //Record the layer
10
11 //Derive a layer that corresponds to some  $\sigma_i$ .
12 create_layer( $\pi, E, W, L_{min}, L_{max}$ ):
13   //A condition " $\pi(u) < \pi(v)$ " ensures layer's acyclicity.
14    $\mathcal{V} = \{(u, v) \in V \times V : \pi(u) < \pi(v)\}$ 
15   //Init a priority queue Q. One queue element is a pair of vertices from  $\mathcal{V}$ .
16    $Q.\text{init}(\mathcal{V})$ 
17    $\text{incidence}_G = \text{incidence\_matrix}(G)$  //Generate an incidence matrix of G.
18    $p_{cnt} = 0$  //Init path count variable pcnt.
19   while  $(\mathcal{V} \neq \emptyset \text{ and } p_{cnt} < M)$  do:
20      $(u, v) = Q.\text{pop}()$  //Get a pair of vertices with lowest priority (lowest number of added paths).
21      $\text{path} = (v_1, v_2, \dots, v_d) \leftarrow \text{find\_path}(u, v, W, \text{incidence}_G, L_{min}, L_{max})$  //Find path from u to v, with path of
        length  $\in [L_{min}, L_{max}]$ , minimizing the sum of edge weights, such that for  $i < j$ :  $\pi(v_i) < \pi(v_j)$  and the
        available edges are given by the  $\text{incidence}_G$  matrix.
        path exists: //Check if u and v are connected, in the graph given by the  $\text{incidence}_G$  matrix.
         $p_{cnt} = p_{cnt} + 1$ 
        foreach  $\text{link} \in \text{path}$  do:
             $L = E' \cup \{\text{link}\}$  //Add each edge from path to the current layer.
        foreach  $v_i, v_j \in \text{path}$  where  $|i - j| > 1$  do:
             $\text{incidence}_G[v_i][v_j] = 0$  //Exclude all edges, which will force the traffic from  $v_1$  to  $v_d$  to use a
            different path than one that was found.
        foreach  $v_i, v_j \in \text{path}$  where  $j - i < L_{min}$  do:
             $\mathcal{V} = \mathcal{V} \setminus (v_i, v_j)$  // Even if an additional path from  $v_i$  to  $v_j$  of length at least  $L_{min}$  will be added,
            there will still exist this shorter path between these vertices (the one contained by the
            path  $(v_1, v_2, \dots, v_d)$ ), therefore the pair  $(v_i, v_j)$  should be further excluded from the  $\mathcal{V}$  set.
22 find_path( $\text{src}, \text{dst}, W, \text{incidence}_G, L_{min}, L_{max}$ ):
23    $\text{path} = \text{null}$  //Init variable containing the path with lowest cost.
24    $Q = \emptyset$  //Init an empty queue.
25    $Q.\text{push}(\text{src}, \text{src})$ 
26   while  $Q \neq \emptyset$  do:
27      $(u, \text{cost}) = Q.\text{pop}()$ 
28     if  $\text{path.last()} == \text{dst}$  then:
29       return  $\text{path}$ 
30     foreach  $(v, w) \in \text{incidence}_G[u]$  do:
31        $\text{edge} = \text{path} \cup \{(u, v)\}$ 
32        $\text{cost}' = \text{cost} + w$ 
33        $Q.\text{push}(v, \text{cost}')$ 
34   return  $\text{path}$ 

```

Creating layers

Deploying layers (VLAN, ...)

```

1 L = {E} //Init a set of layers L; we start with E that corresponds to  $\sigma_1$ .
2  $P = \{\pi_1(V), \dots, \pi_{n-1}(V)\}$  //Generate n-1 random permutations of vertices.
3 foreach  $\pi \in P$  do: //One iteration of the main loop derives one layer.
4    $E' = \{\}$ ;
5   //Generate a layer.
6   //Record the layer
7
8
9

```



# FATPATHS ARCHITECTURE

Design details

```

1 L = {E} //Init a set of layers L; we start with E that corresponds to  $\sigma_1$ .
2  $\Pi = \{\pi_1(V), \dots, \pi_{n-1}(V)\}$  //Generate  $n-1$  random permutations of vertices.
3
4 //Init a matrix W containing weights of edges  $(u,v) \in E$ .
5  $W = \{[w_{uv}] \mid \forall u,v \in V : w_{uv} = 0\}$ 
6
7 foreach  $\pi \in \Pi$ 
8    $E' = \text{create\_layer}(E, \pi)$ 
9    $L = L \cup \{E'\}$ 
10
11 //Derive a layer  $i$ 
12 create_layer( $E, \pi$ )
13   //A condition for creating a layer
14    $V = \{(u,v) \in E \mid \pi(u) < \pi(v)\}$ 
15   //Init a priority queue Q
16    $Q.\text{init}(V)$ 
17    $\text{incidence}_G = \{ \}$ 
18    $\text{pcnt} = 0$  //Initial counter
19   while  $(V \neq \emptyset)$ 
20      $(u,v) = Q.\text{pop}()$ 
21      $\text{path} = (v_1, \dots, v_k)$ 

```

Creating layers

Deploying layers (VLAN, ...)

```

1 //Input:  $E'_i, i \in \{1, \dots, n\}$ : the specification of each layer.  $E'_i$  are router-router links in each layer  $i$ ,
2 //produced by FatPaths layer construction algorithms specified in Listing 1 or Listing 2.
3 //Output:  $F_{i,s}$ : a forwarding table in router  $s \in V$  within layer  $i$  ( $1 \leq i \leq n$ ).
4 //  $F_{i,s}[t] \in \{1, \dots, k\}$  is a port number (in router  $s \in V$ ) that leads to router  $t \in V$  (within layer  $i$  ( $1 \leq i \leq n$ )).
5
6 //Compute shortest paths between routers in each layer  $i$ . First, initialize the function  $\sigma_i$ .
7 //  $\sigma_i(s,t)$  is a port number (in router  $s$ ) that leads to router  $t$  (in layer  $i$ );  $d$  is an auxiliary structure.
8 foreach  $i \in \{1, \dots, n\}$  do: //For each layer...
9   foreach  $(s,t) \in V \times V$  do: //For each router pair...
10    if  $s == t$  then:
11       $d_{st} = 0$  //The distance between  $u$  and  $v$  is zero, if  $u == v$ .
12       $\sigma_i(s,t) = s$ 
13    else if  $(s,t) \in E'_i$  then:
14       $d_{st} = 1$  //The distance between directly connected routers is 1.
15       $\sigma_i(s,t) = t$ 
16    else:
17       $d_{st} = +\infty$  //The initial distance between non-adjacent routers is infinity.
18       $\sigma_i(s,t) = \text{null}$ 
19
20 //For each layer, derive the routing functions:
21 foreach  $z \in |V|$  do:
22   foreach  $s \in |V|$  do:
23     foreach  $t \in |V|$  do:
24       if  $d_{st} > d_{sz} + d_{zt}$  then:
25          $d_{st} = d_{sz} + d_{zt}$ 
26          $\sigma_i(s,t) = \sigma_i(s,z)$ 
27
28 //Once  $\sigma_i$  is computed, populate a forwarding table  $F_{i,s}$  for each router  $s \in V$ .
29 foreach  $s \in V$  do: //For every router  $s \in V$ 
30    $F_s = \{ \}$  //Initialize the forwarding table
31
32 //Build a forwarding table within layer  $i$  using the previously devised structures:
33 foreach  $s \in V$  do:
34   foreach  $t \in V, t \neq s$  do:
35      $F_s[t] = \sigma_i(s,t)$ 

```

Populating forwarding entries

# FATPATHS ARCHITECTURE

Design details

```

1 L = {E} //Init a set of
2 Π = {π1(V), ..., πn-1(V)} //G
3
4 //Init a matrix W containi
5 W = {{wuv} | ∀u,v ∈ V: wuv
6
7 foreach π ∈ Π
8   E' = create
9   L = L ∪ {E'}
10
11 //Derive a la
12 create_layer(
13 //A conditi
14 V = {(u,v) ∈
15 //Init a pr
16 Q.init(V)
17 incidenceG =
18 pcnt = 0 //I
19 while (V ≠
20 (u,v) = Q.p
21 path = (v1

```

Creating layers

Deploying layers (VLAN, ...)

## 3 Purified Transport based on NDP [47] (§3)

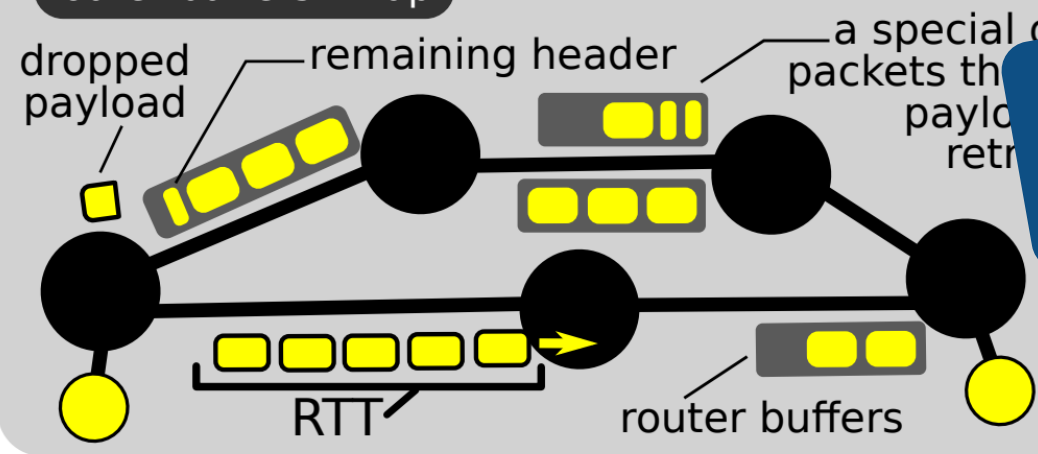
What is a design of high-performance transport layer?

Use shallow buffers in routers

Start sending at line rate

Drop only payload if router buffers fill up

Prioritize packets with dropped payload and retransmitted packets



Flow control

Populating forwarding entries

```

1 //Input
2 //prod
3 //Output
4 //Fi,s[t]
5
6 //Comp
7 //σi(s,t)
8 foreach
9   fore
10  if
11
12
13  el
14
15  el
16
17
18
19
20 //For each layer, derive the routing functions:
21 foreach z ∈ |V| do:
22   foreach s ∈ |V| do:
23     foreach t ∈ |V| do:
24       if dst > dsz + dzt then:
25         dst = dsz + dzt
26         σi(s,t) = σi(s,z)
27
28 //Once σi is computed, populate a f
29 foreach s ∈ V do: //For every router
30   Fs = {} //Initialize the forwarding
31
32 //Build a forwarding table within la
33 foreach s ∈ V do:
34   foreach t ∈ V, t ≠ s do:
35     Fs[t] = σi(s,t)

```

# FATPATHS ARCHITECTURE

Design details

# Check the paper ☺

**3 Purified Transport based on NDP [47] (§3)**

What is the design of the transport layer?

- Use shallow congestion control
- Control sending rate
- Prioritize packets with dropped payload and retransmitted packets

Drop only payload if router buffers fill up

Flow control

```

1 L = {E} //Init a set of layers
2 Π = {π1(V), ..., πn-1(V)} //Init a set of paths
3
4 //Init a matrix W containing weights
5 W = {{wuv} | ∀u,v ∈ V: wuv > 0}
6
7 foreach π ∈ Π
8   E' = create_layer(π)
9   L = L ∪ {E'}
10
11 //Derive a layer L'
12 create_layer(L')
13 //A condition for a layer to be added
14 V = {(u,v) ∈ E' | wuv > 0}
15 //Init a priority queue Q
16 Q.init(V)
17 incidenceG = G
18 pcut = 0 //Init a counter
19 while (V ≠ ∅)
20   (u,v) = Q.pop()
21   path = (v1, ..., vk)

```

```

1 //Input
2 //produce a set of paths
3 //Output
4 //F_i,s[t]
5
6 //Compute the routing functions
7 //σ_i(s,t)
8 foreach i ∈ {1, ..., n}
9   foreach s ∈ V
10    if s ≠ t
11     d_st = d_sz + d_zt
12     σ_i(s,t) = σ_i(s,z)
13
14 //Once σ_i is computed, populate a forwarding table
15 foreach s ∈ V do: //For every router
16   F_s = {} //Initialize the forwarding table
17
18 //Build a forwarding table within a layer
19 //using the previously devised structures:
20 //For each layer, derive the routing functions:
21 foreach z ∈ |V| do:
22   foreach s ∈ |V| do:
23     foreach t ∈ |V| do:
24       if d_st > d_sz + d_zt then:
25         d_st = d_sz + d_zt
26         σ_i(s,t) = σ_i(s,z)
27
28 //Once σ_i is computed, populate a forwarding table
29 foreach s ∈ V do: //For every router
30   F_s = {} //Initialize the forwarding table
31
32 //Build a forwarding table within a layer
33 //using the previously devised structures:
34 foreach s ∈ V do:
35   foreach t ∈ V, t ≠ s do:
36     F_s[t] = σ_i(s,t)

```

Creating layers

Deploying layers (VLAN, ...)

Populating forwarding entries

```

1 L = {E} //Init a set of layers
2 P = {π1(V), ..., πn-1(V)} //Init a set of paths
3 foreach π ∈ P do: //One iteration
4   E' = {}; foreach v ∈ V
5     //Init a path
6     path(src, v) = π
7     Q = Q ∪ {v}
8   //Build a forwarding table within a layer
9   //using the previously devised structures:

```

# FATPATHS ARCHITECTURE

Design details

Check the paper 😊

3 Purified Transport based on NDP [47] (§3)

Use shallow... sending... line rate

Prioritize packets with dropped payload and retransmitted packets

Drop only payload if router buffers fill up

dropped payload remaining header a special queue packets th...

Flow

Key design insight: Layered routing enables (1) easy encoding of the diversity of multiple minimal and non-minimal paths, and (2) enables simple and robust load balancing

```

7
8
9
(VLAN, ...) edge
33 foreach s ∈ V do:
34   foreach t ∈ V, t ≠ s do:
35     F_s[t] = σ_i(s, t)

```

# EVALUATION



## EVALUATION


**Question 1:** Can we get more performance on low-diameter networks than on Fat trees?



## EVALUATION



**Question 1:** Can we get more performance on low-diameter networks than on Fat trees?

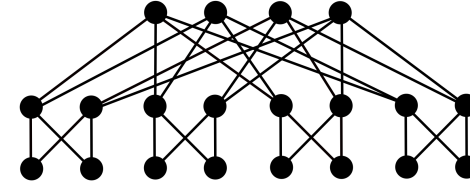
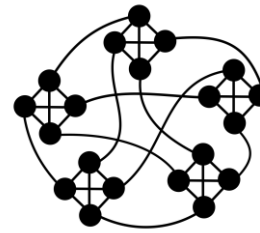
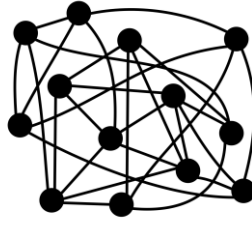
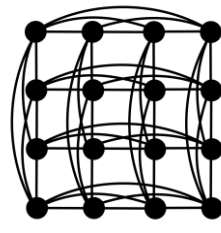
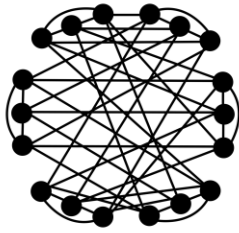


**Question 2:** Does Fatpaths give more performance than existing routing schemes?

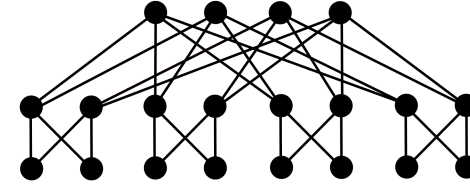
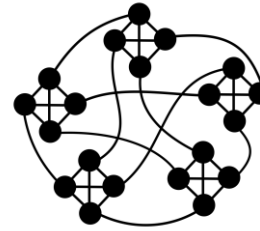
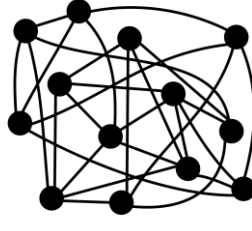
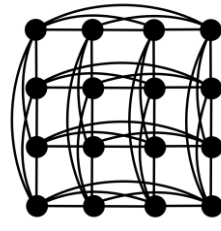
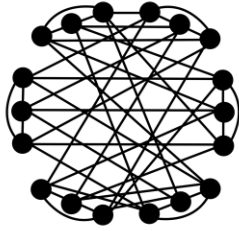
# EVALUATION



# EVALUATION

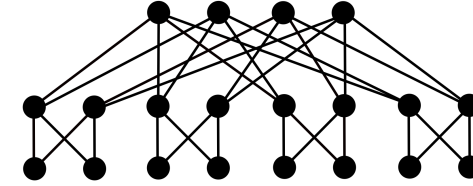
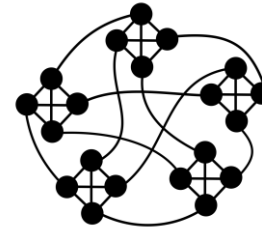
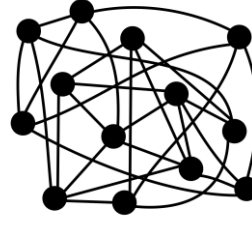
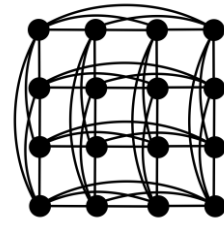
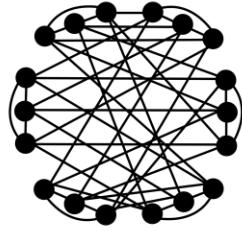


## EVALUATION



We pick Ethernet as a setting  
to accelerate as many Top500  
systems as possible

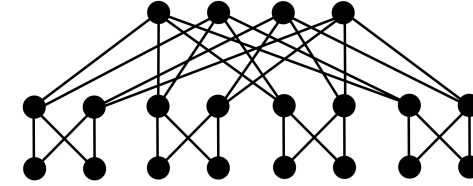
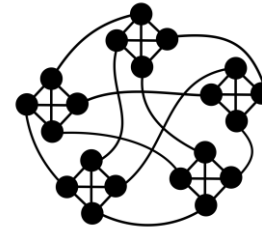
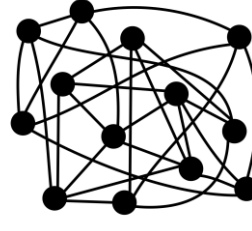
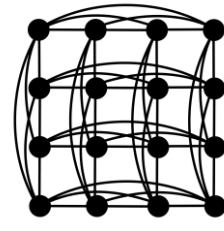
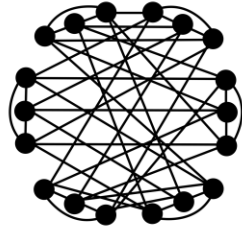
## EVALUATION



We pick Ethernet as a setting to accelerate as many Top500 systems as possible

(more than 50% of systems in Top500 use Ethernet. However, they are not as efficient as InfiniBand and others (details in the paper))

## EVALUATION



We pick Ethernet as a setting to accelerate as many Top500 systems as possible

(more than 50% of systems in Top500 use Ethernet. However, they are not as efficient as InfiniBand and others (details in the paper))

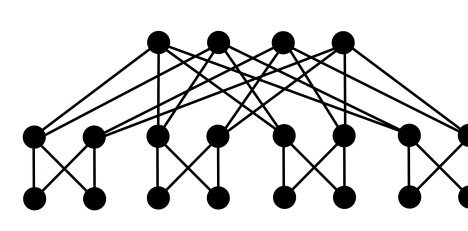
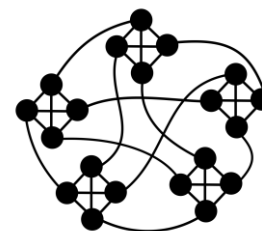
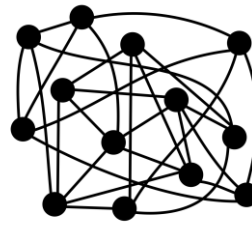
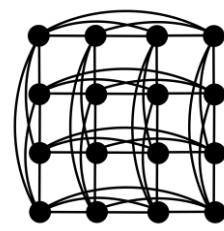
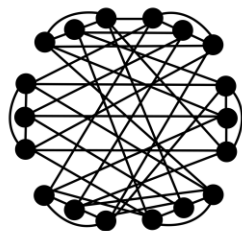
### Setting 1 „Bare Ethernet”

Ethernet, but no TCP.  
**Simulator:** htsim [1].

**NDP [1]:** a very recent baseline, originally for fat trees.

[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

# EVALUATION



We pick Ethernet as a setting to accelerate as many Top500 systems as possible

(more than 50% of systems in Top500 use Ethernet. However, they are not as efficient as InfiniBand and others (details in the paper))

## Setting 1 „Bare Ethernet”

Ethernet, but no TCP.  
**Simulator:** htsim [1].

**NDP [1]:** a very recent baseline, originally for fat trees.

## Setting 2 „Full TCP”

Standard TCP and related.  
**Simulator:** OMNeT++ [2].

**ECMP [3]:** traditional static load balancing.  
**LetFlow [4]:** recent adaptive load balancing.

[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

[2] A. Varga et al. The OMNeT++ discrete event simulation system. ESM'01.

[3] C. Hopps. Analysis of an Equal-Cost Multi-Path Algorithm. RFC2992, 2000.

[4] E. Vanini et al. LetIt Flow: Resilient Asymmetric Load Balancing with Flowlet Switching. NSDI'17.

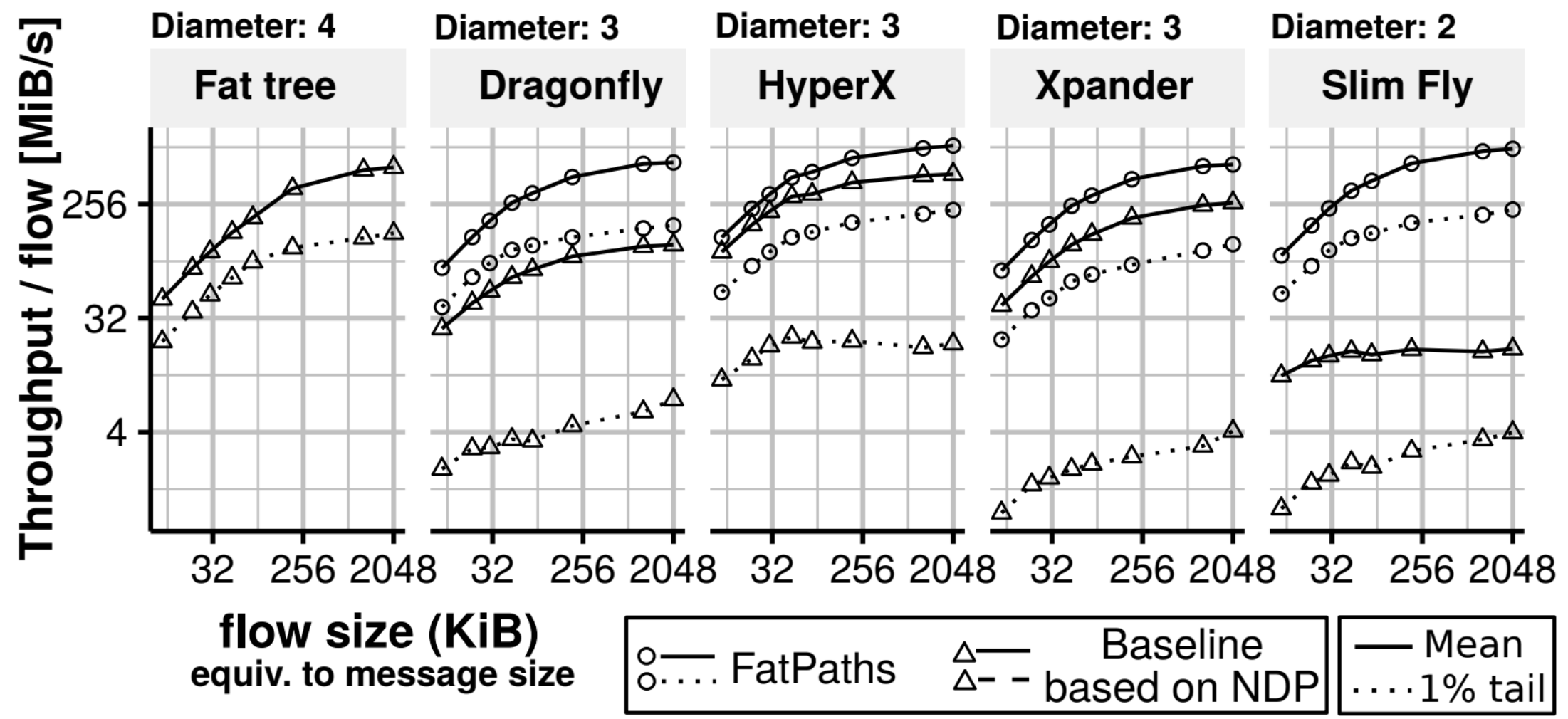
## EVALUATION

$N \approx 10,000$ ; comparable cost; stencil (no remapping); “bare Ethernet” setting

**NDP**: a very recent  
baseline for fat trees [47]

N ≈ 10,000; comparable cost; stencil (no remapping); “bare Ethernet” setting

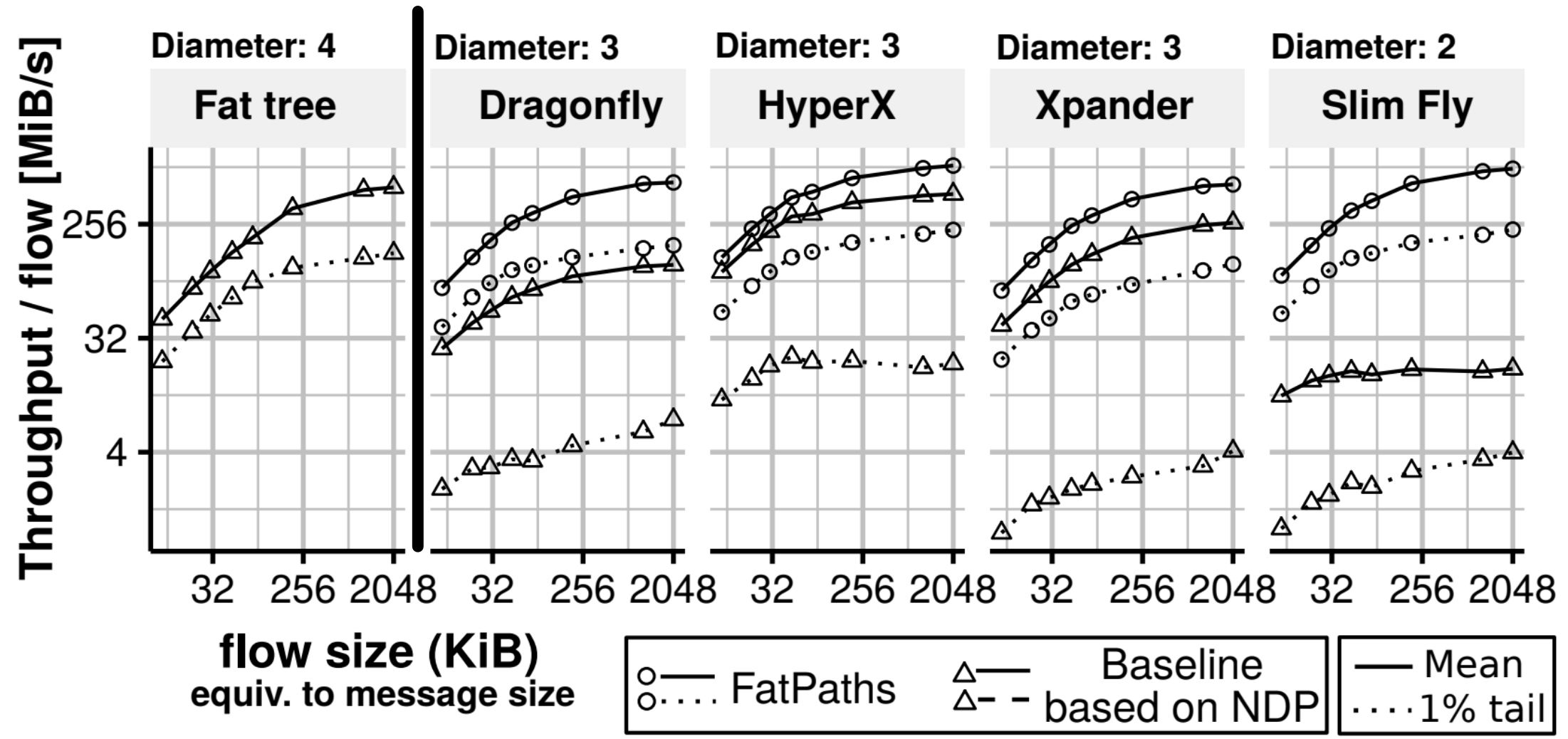
# EVALUATION



**NDP:** a very recent baseline for fat trees [47]

N ≈ 10,000; comparable cost; stencil (no remapping); “bare Ethernet” setting

# EVALUATION

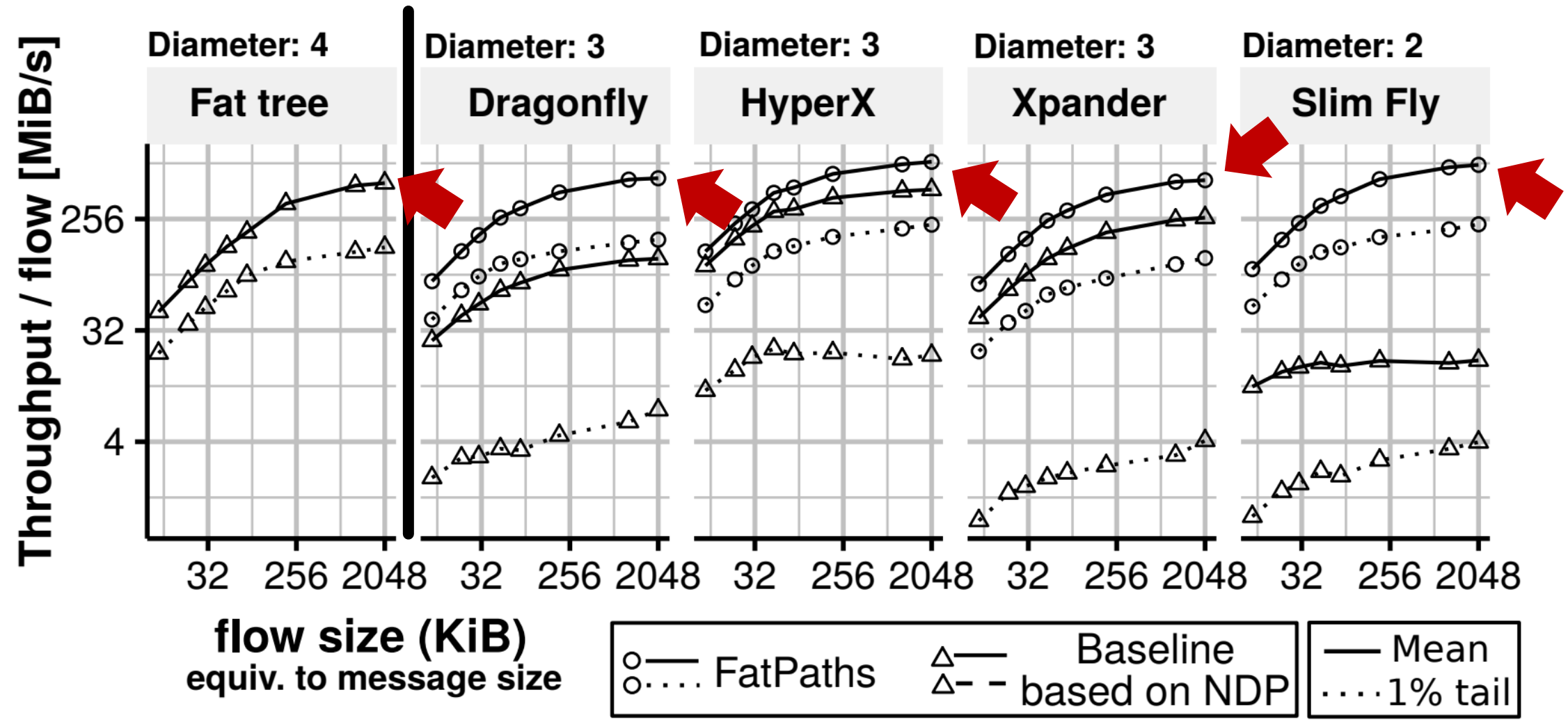


**NDP:** a very recent baseline for fat trees [47]



N ≈ 10,000; comparable cost; stencil (no remapping); “bare Ethernet” setting

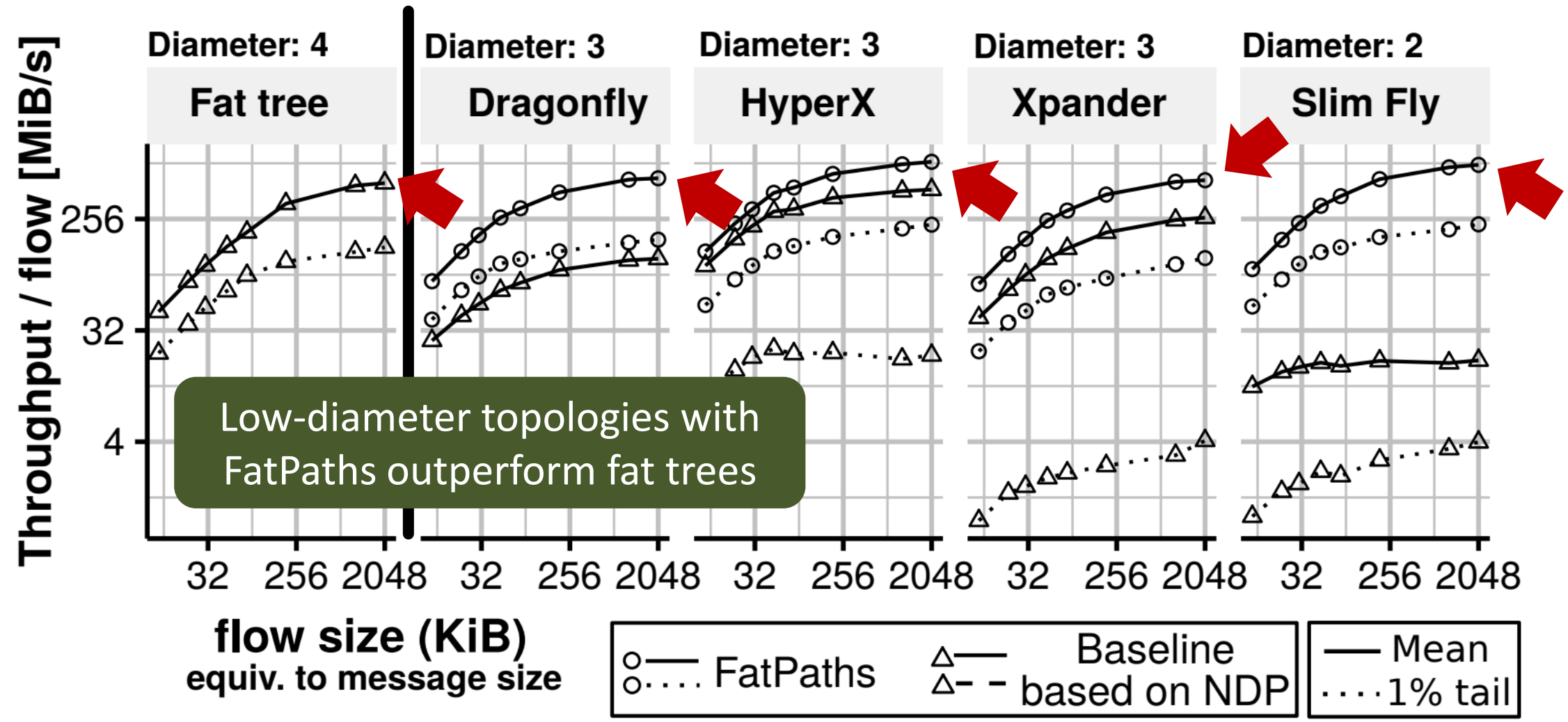
# EVALUATION



**NDP:** a very recent baseline for fat trees [47]

N ≈ 10,000; comparable cost; stencil (no remapping); “bare Ethernet” setting

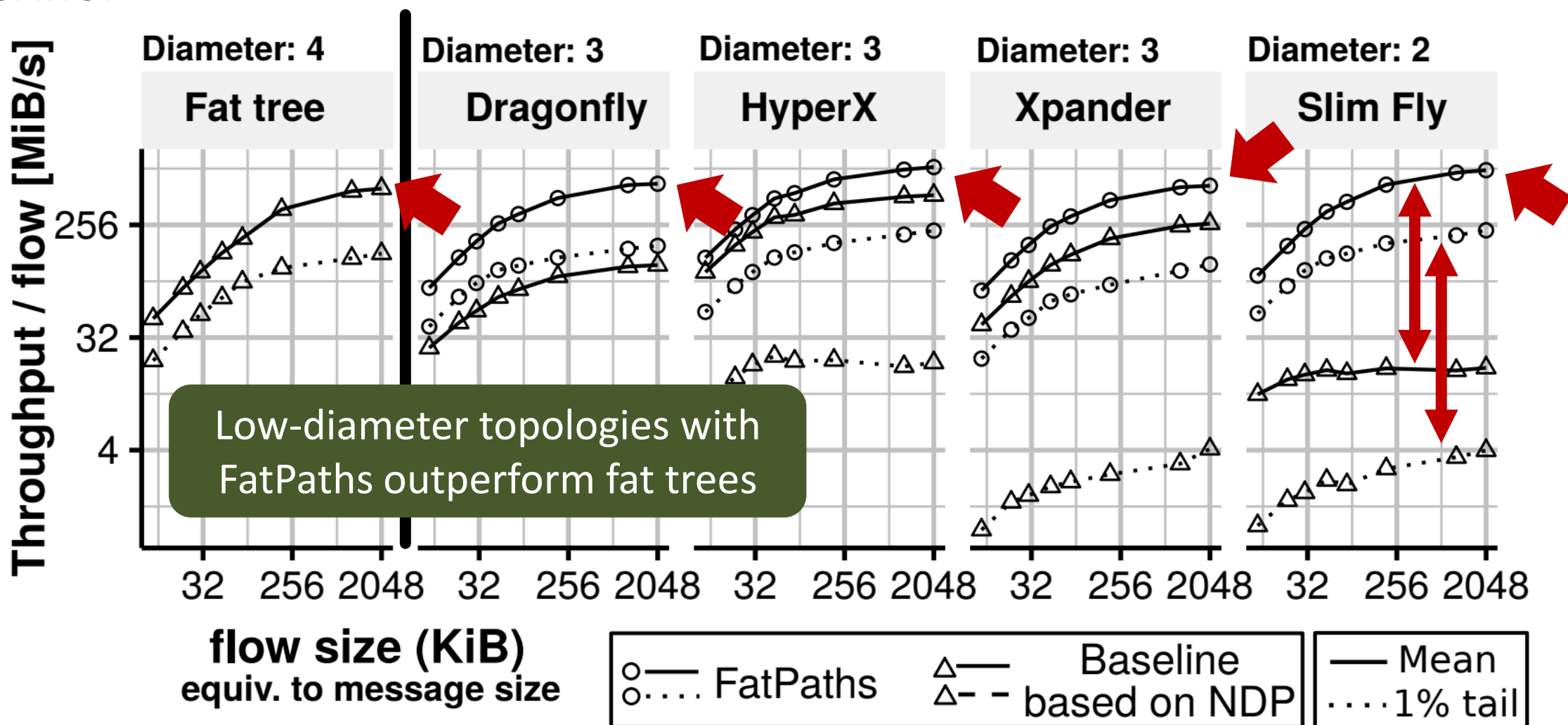
# EVALUATION



**NDP:** a very recent baseline for fat trees [47]

# EVALUATION

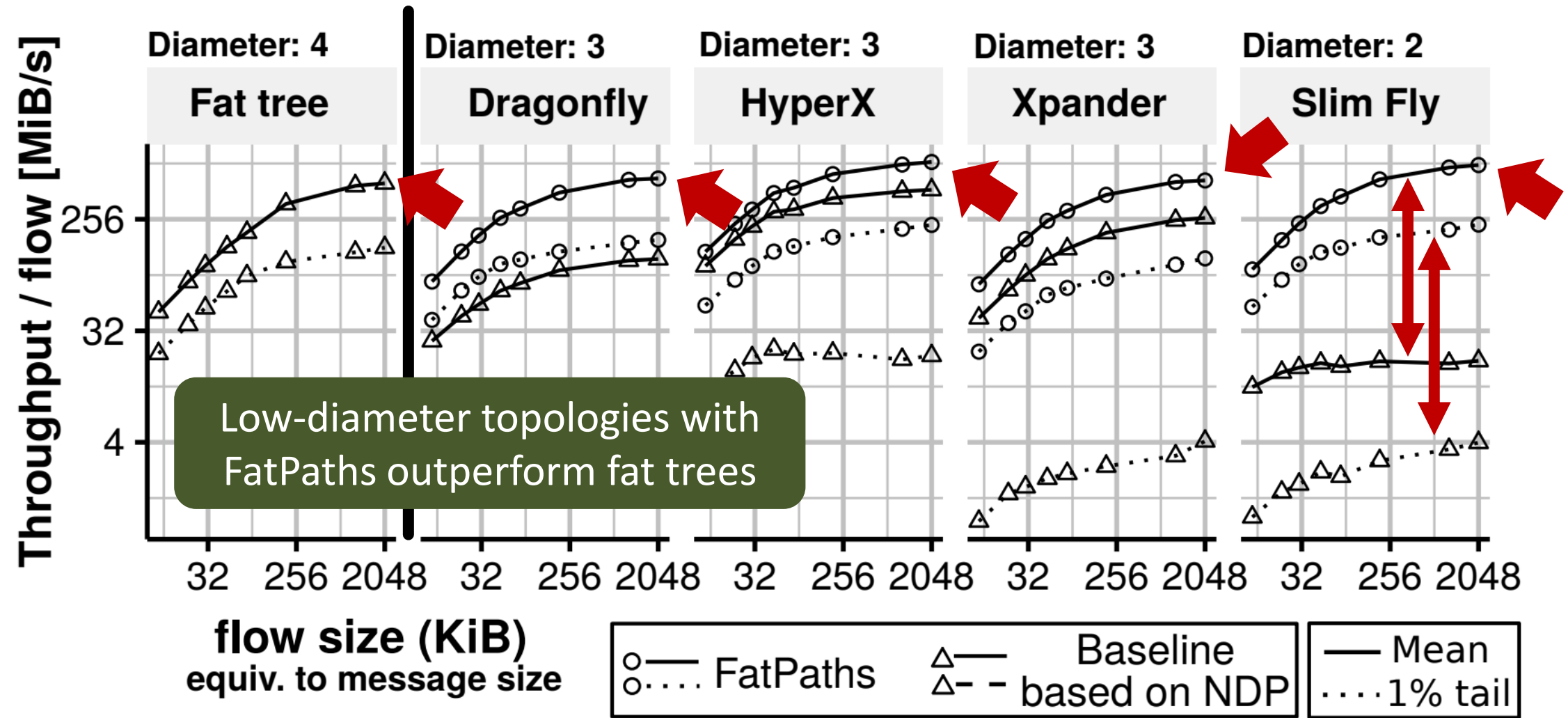
$N \approx 10,000$ ; comparable cost; stencil (no remapping); “bare Ethernet” setting



**NDP:** a very recent baseline for fat trees [47]

N ≈ 10,000; comparable cost; stencil (no remapping); “bare Ethernet” setting

# EVALUATION



**NDP:** a very recent baseline for fat trees [47]

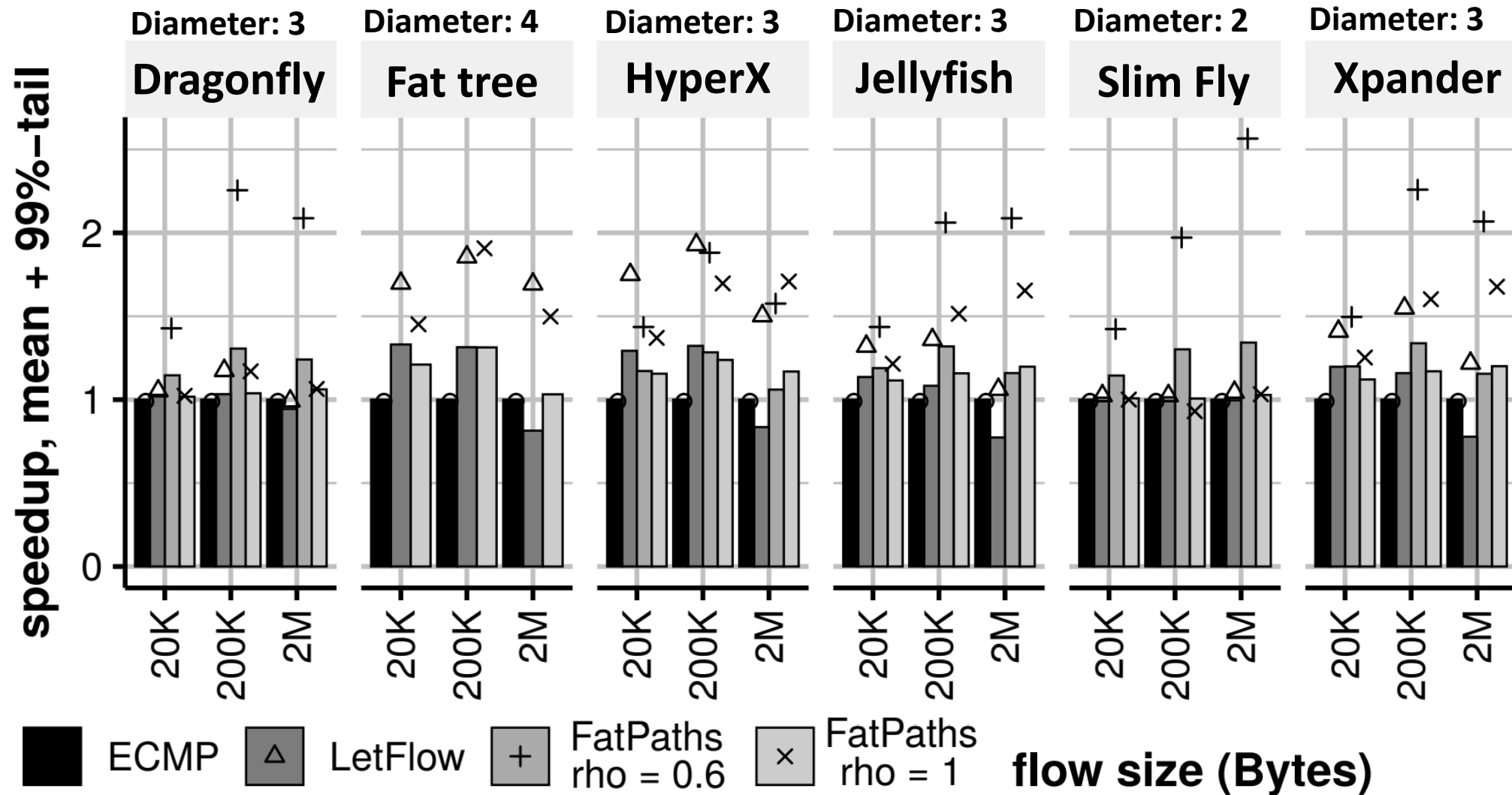
FatPaths outperforms simple NDP executed obliviously on low-diameter networks (up to 30x improvement)

## EVALUATION

$N \approx 10,000$ ; comparable cost; random uniform traffic; “full TCP” setting

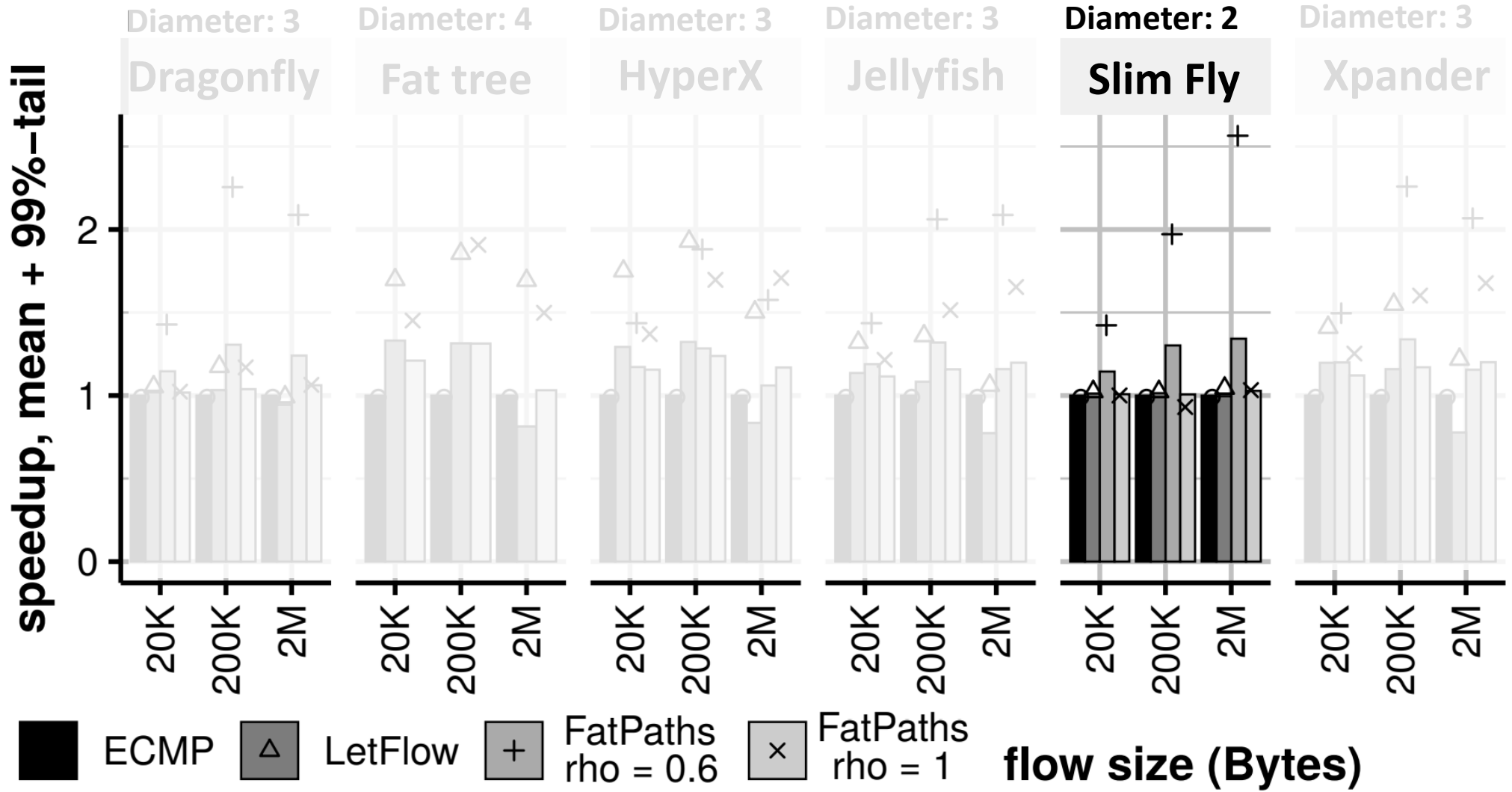
# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting



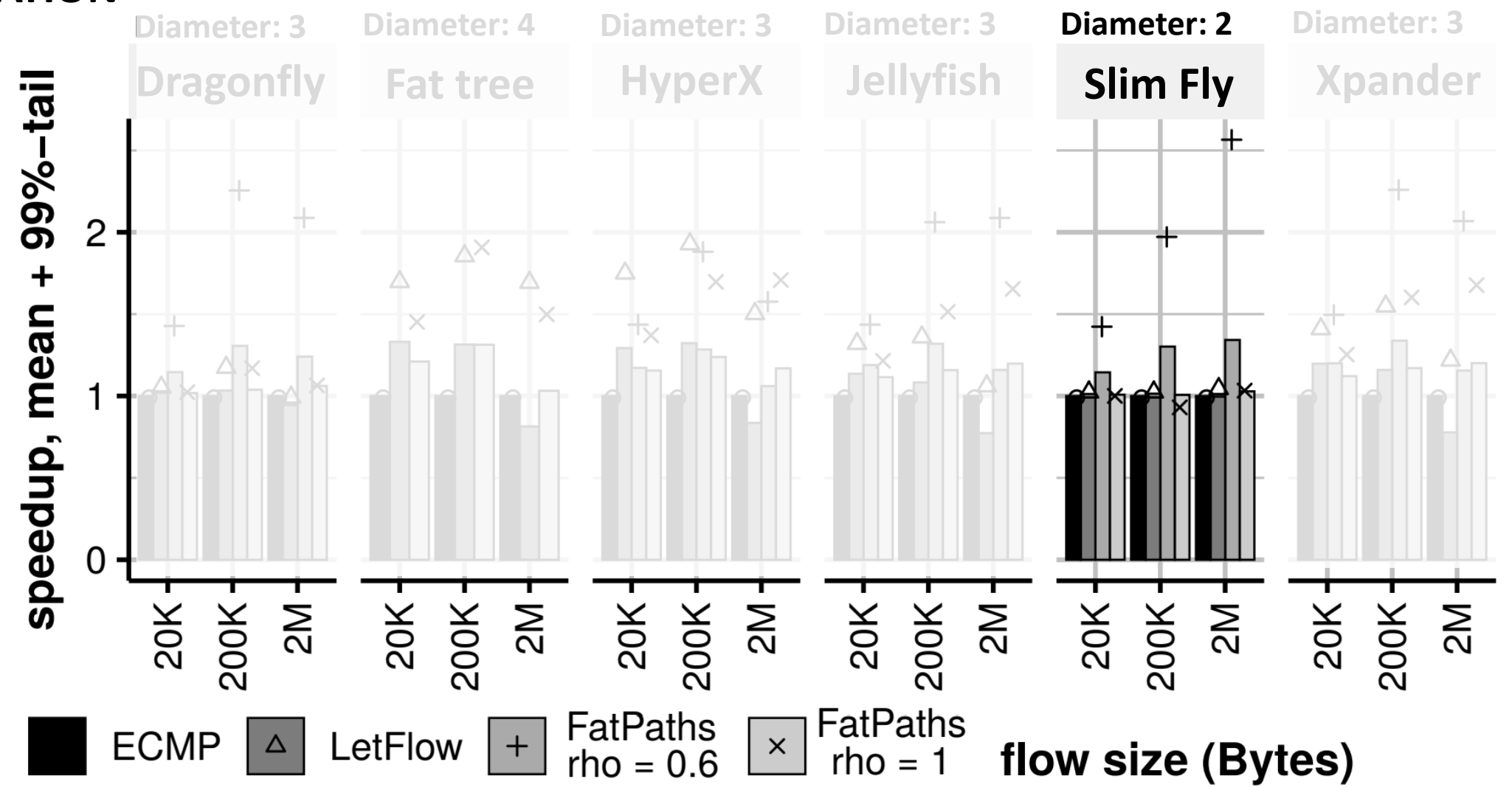
# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting



N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

# EVALUATION



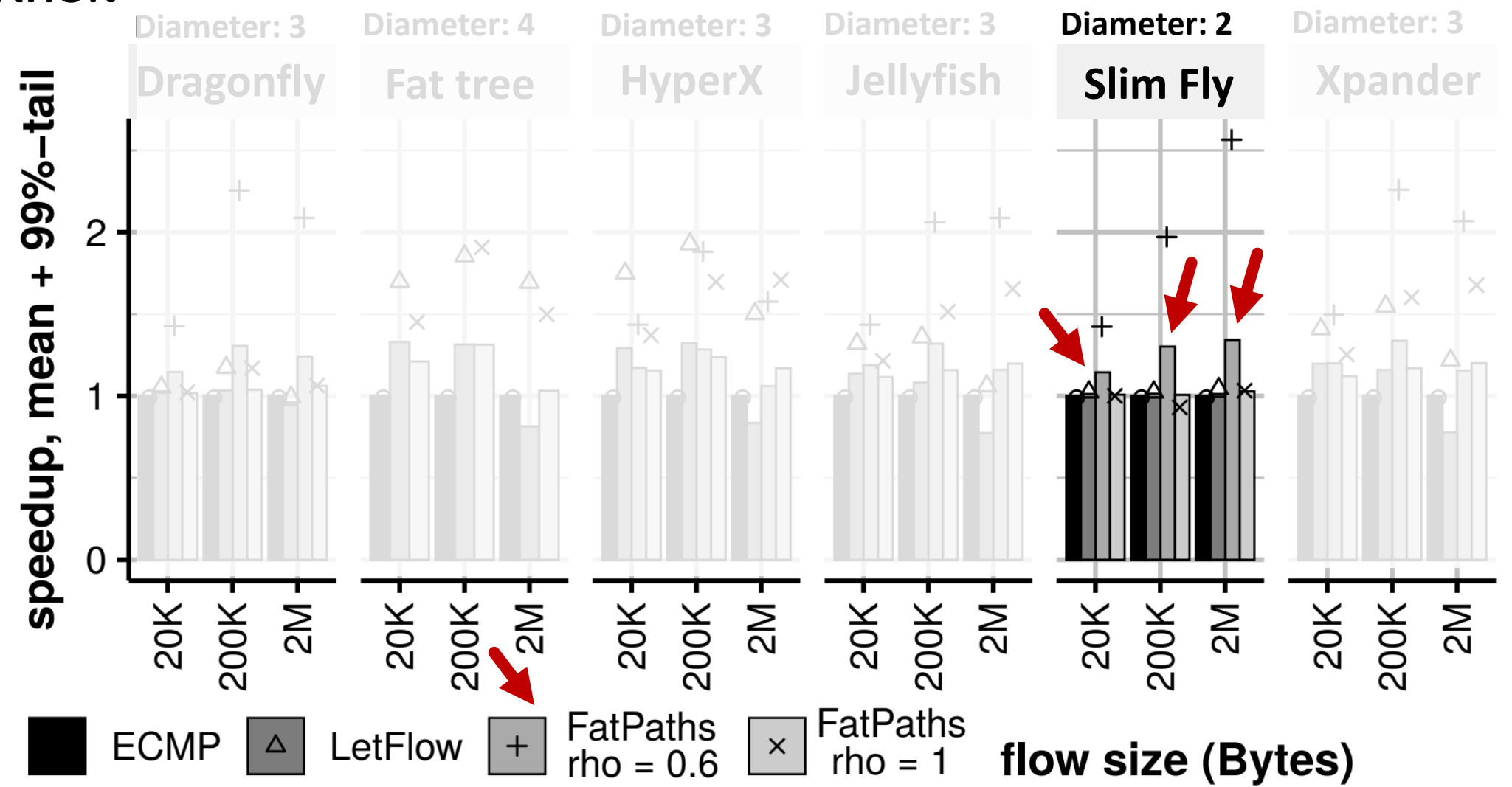
**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer



N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

# EVALUATION

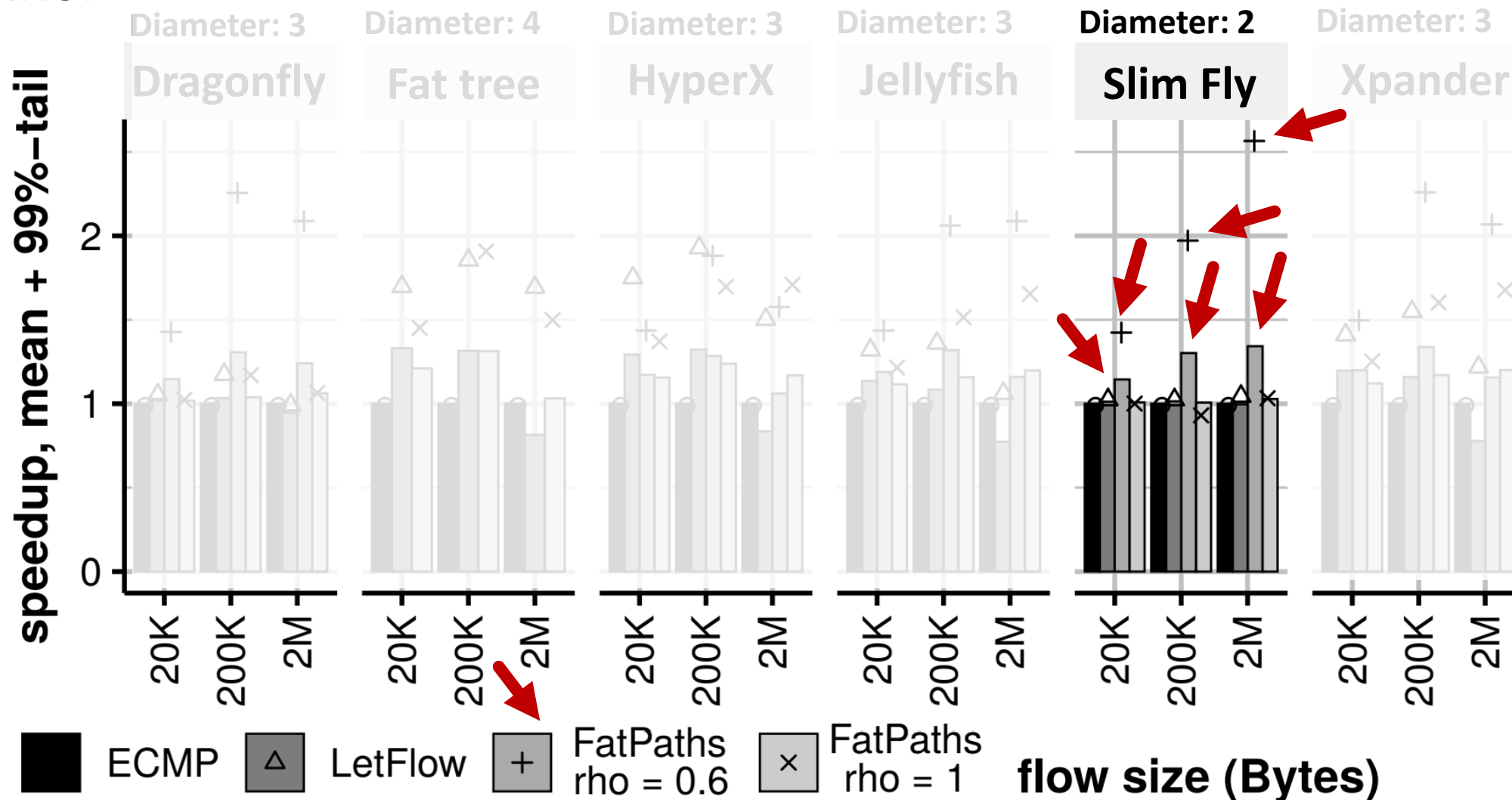


**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

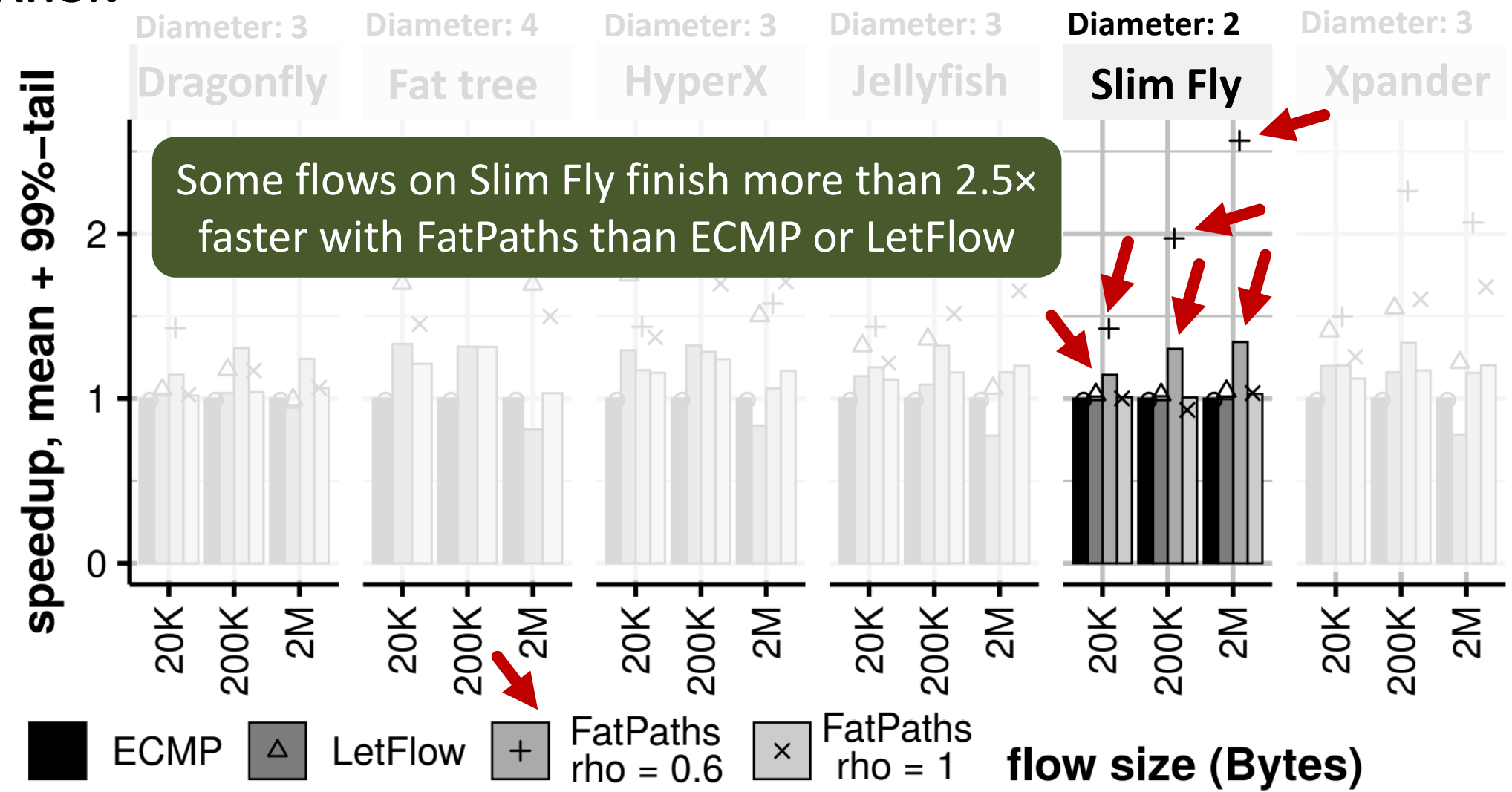


**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

# EVALUATION



**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

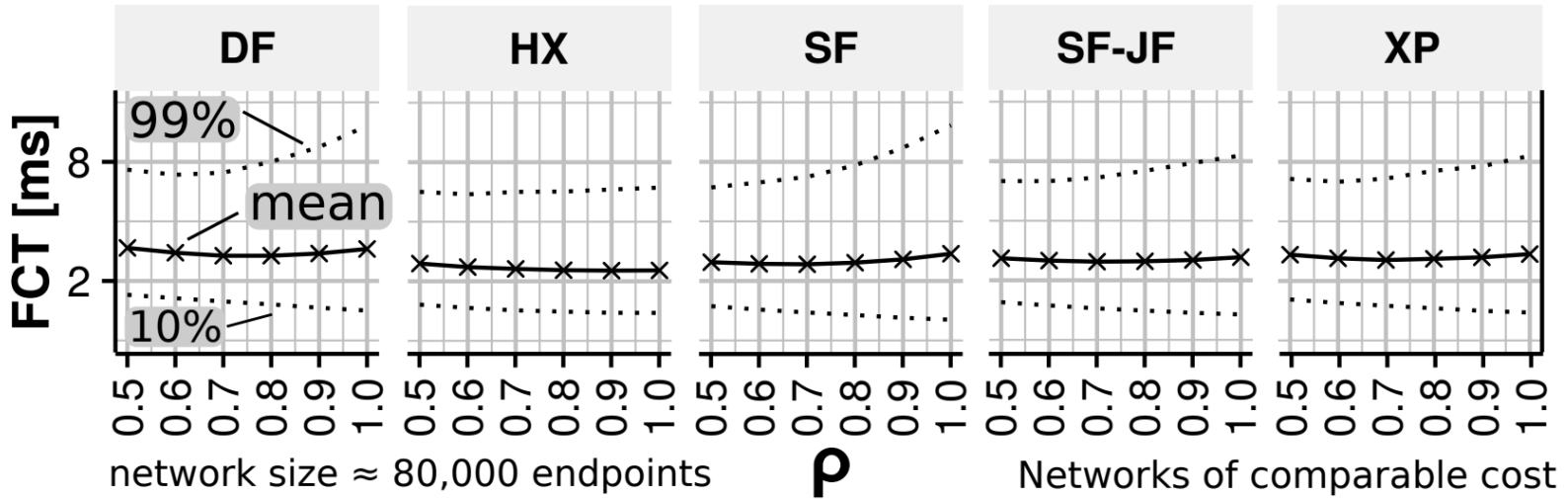
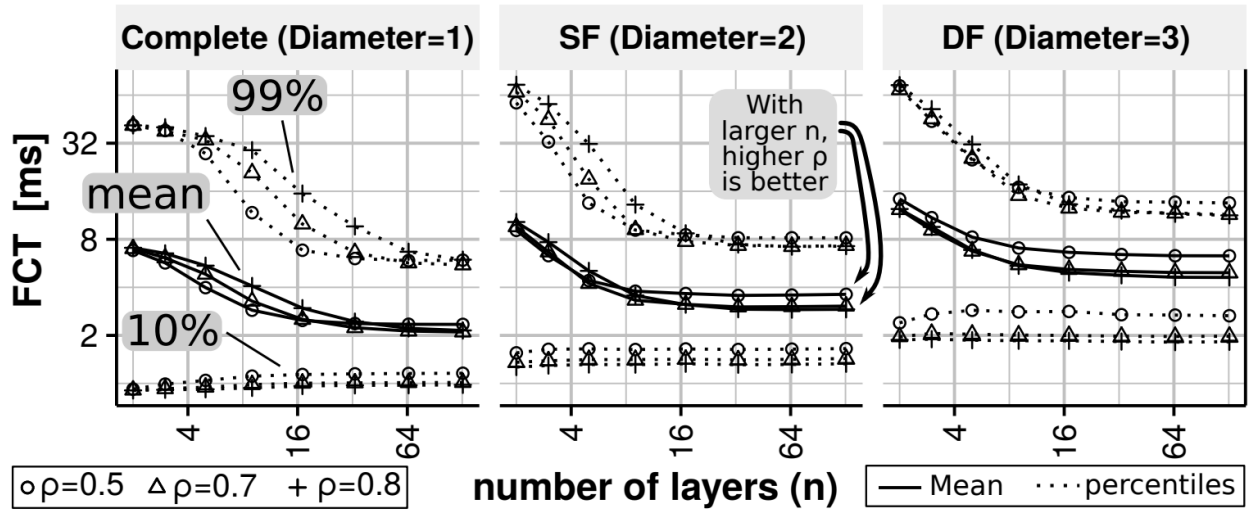
**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

# EVALUATION

Other analyzes

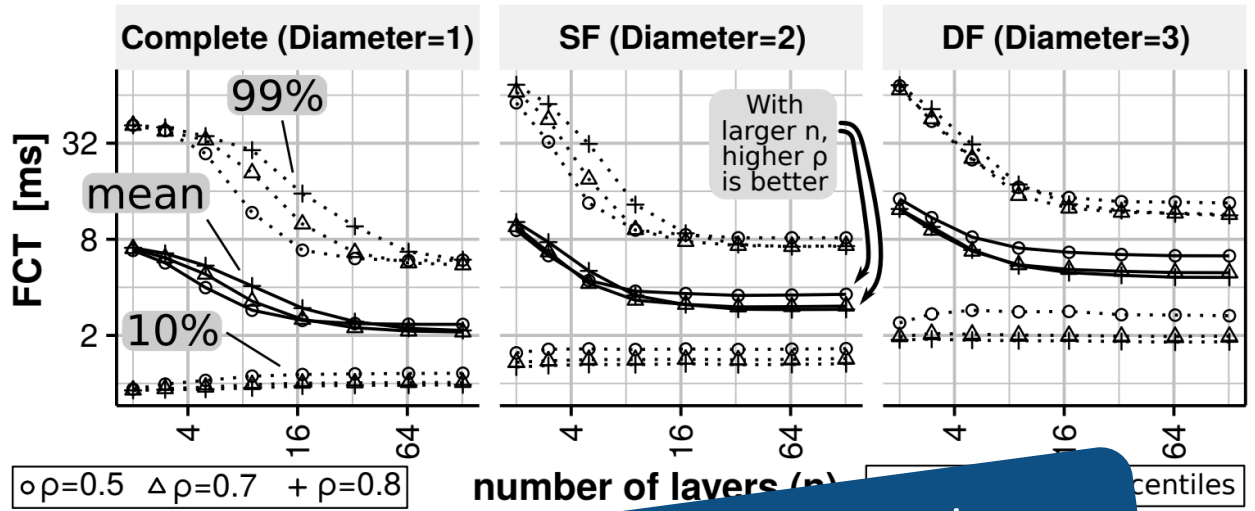
# EVALUATION

## Other analyzes

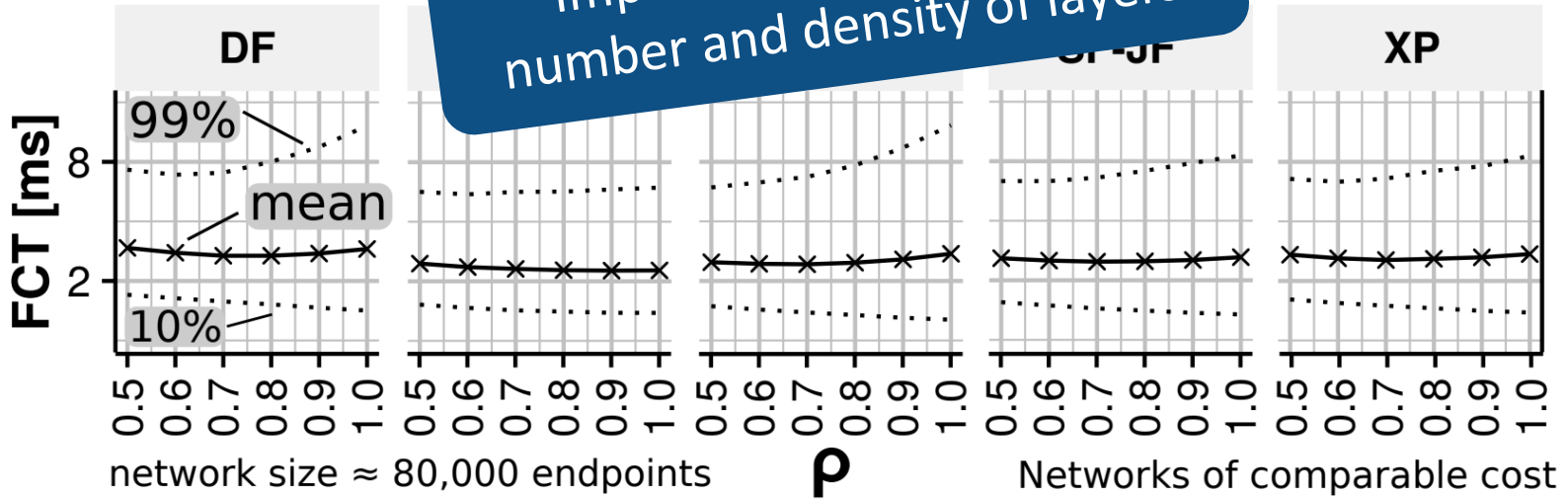


# EVALUATION

Other analyzes

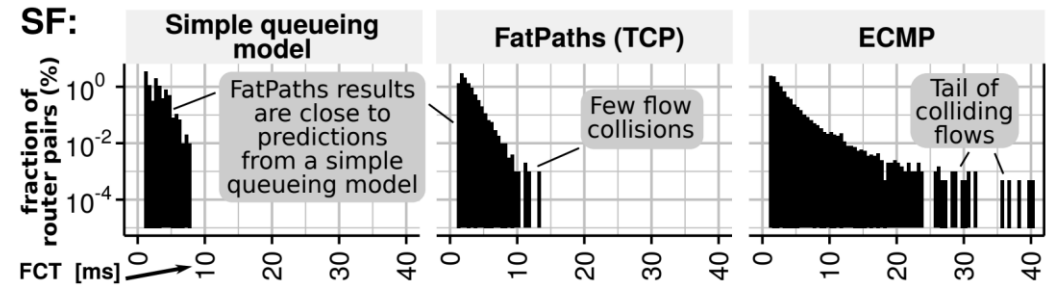
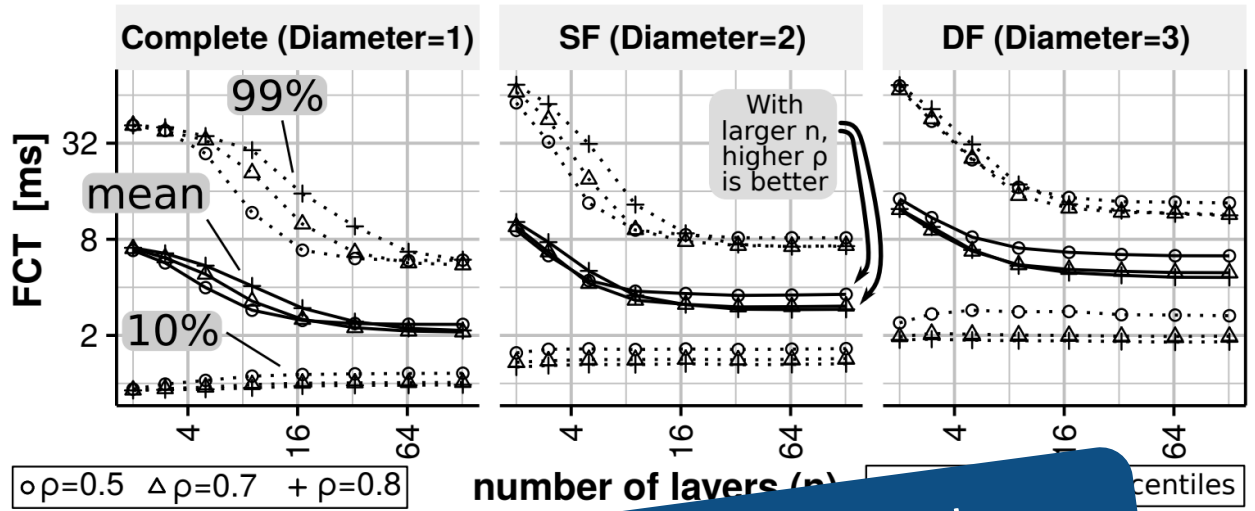


Impact from varying the number and density of layers

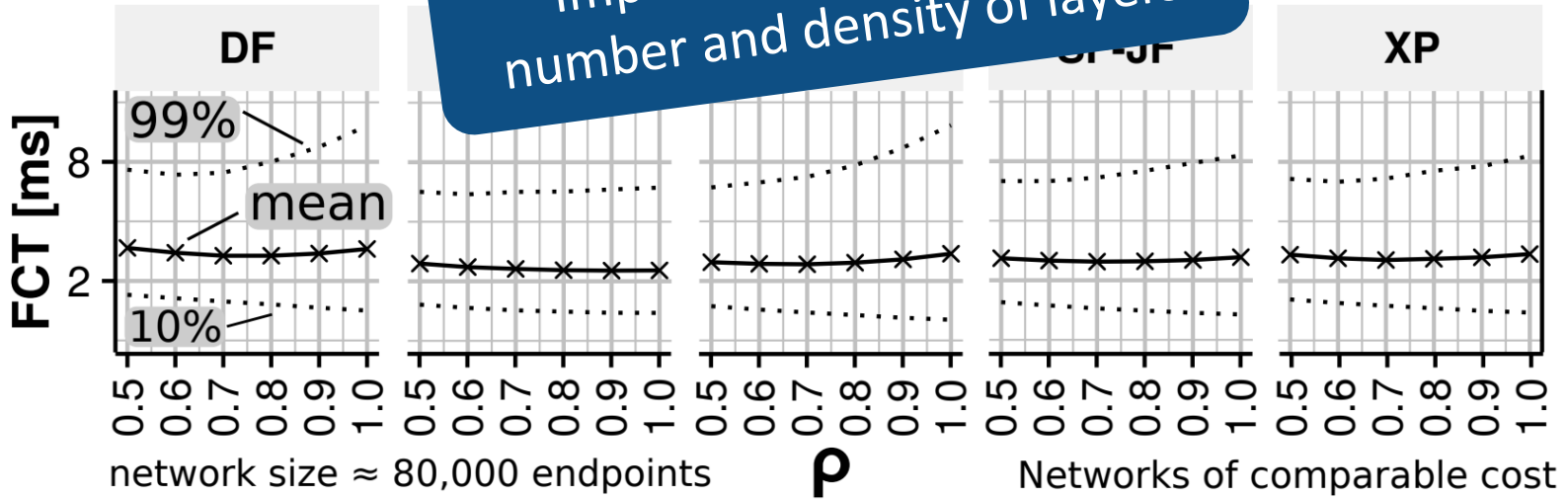


# EVALUATION

## Other analyzes

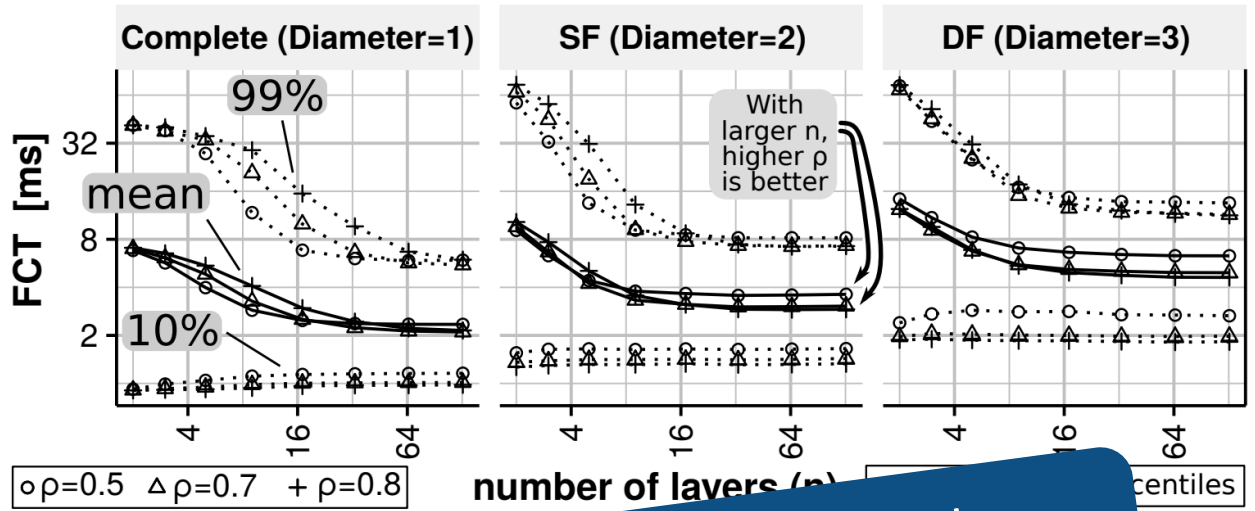


Impact from varying the number and density of layers



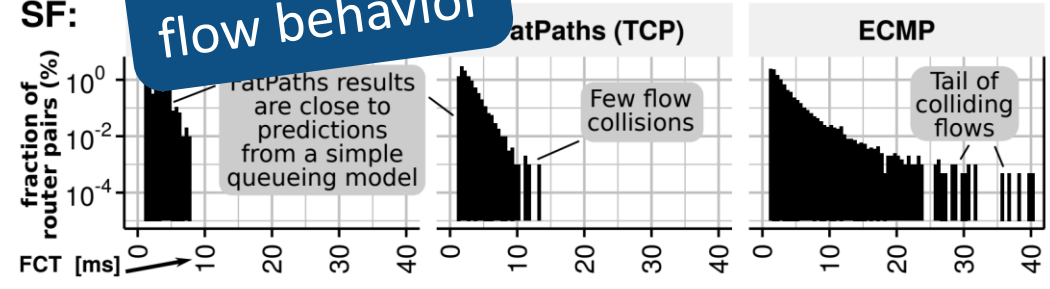
# EVALUATION

## Other analyzes

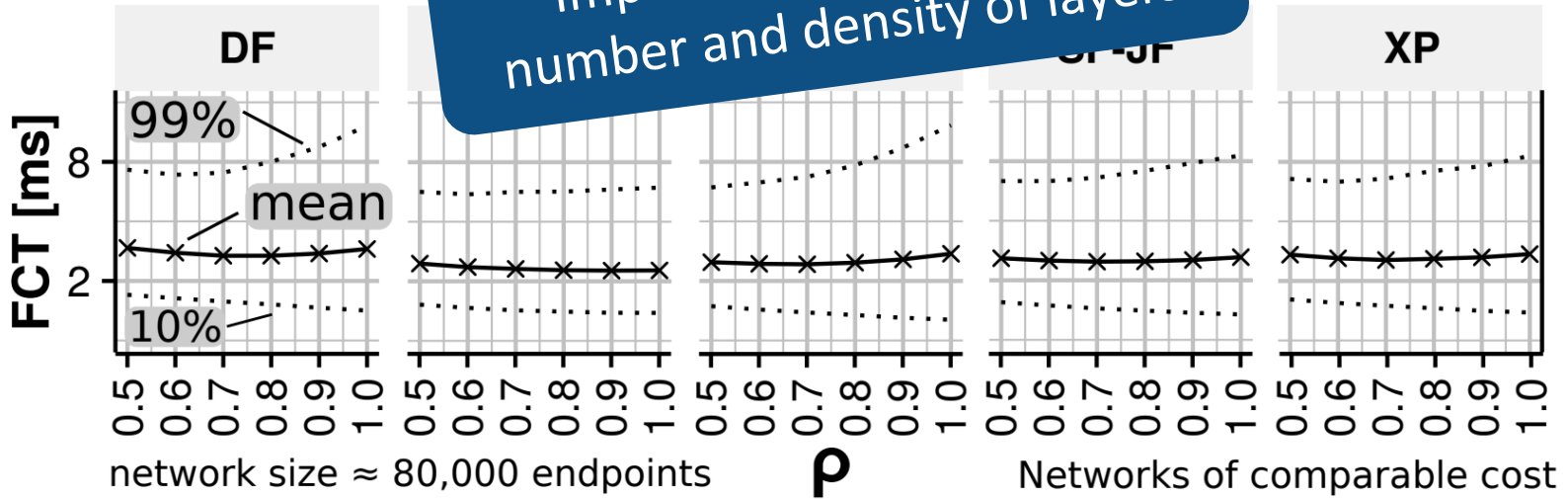


SF:

Details of flow behavior

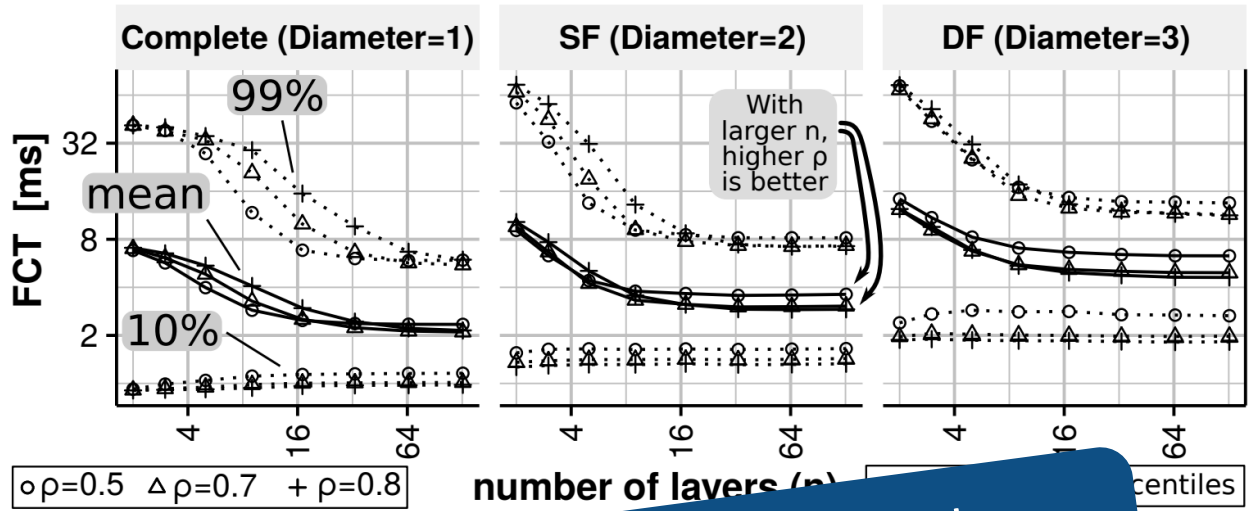


Impact from varying the number and density of layers

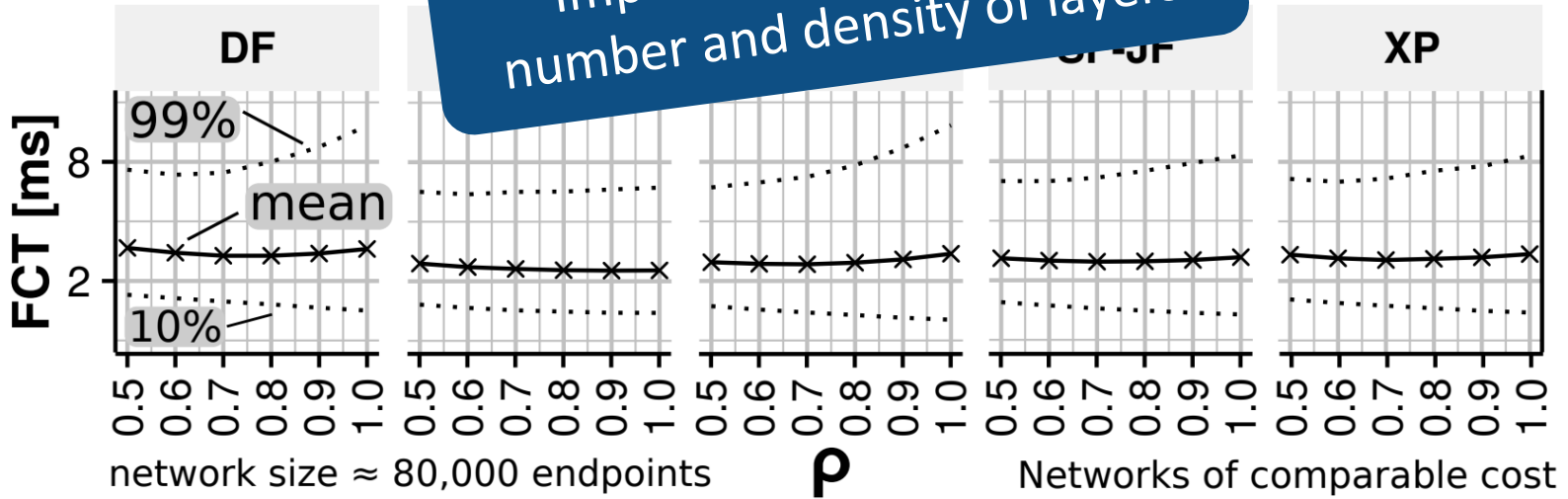




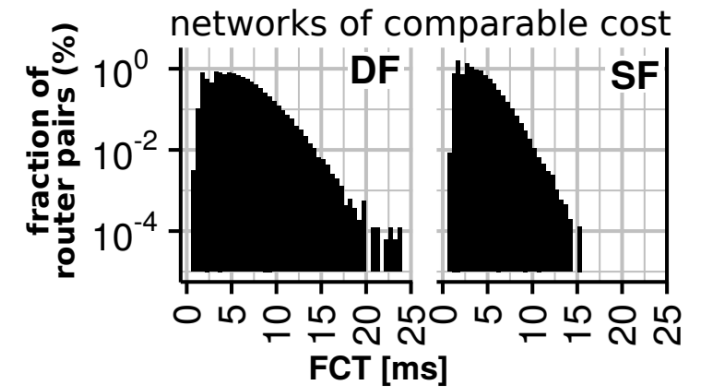
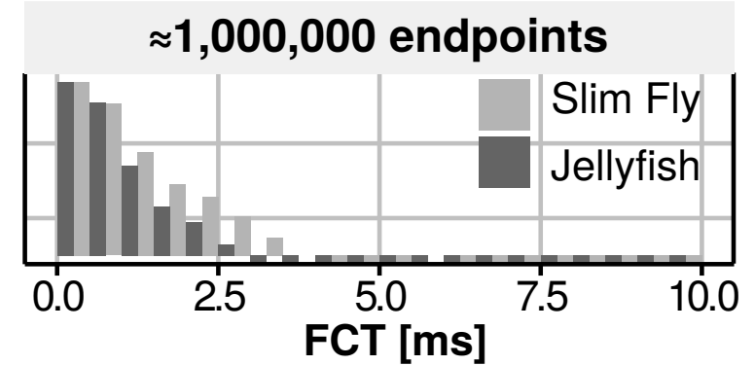
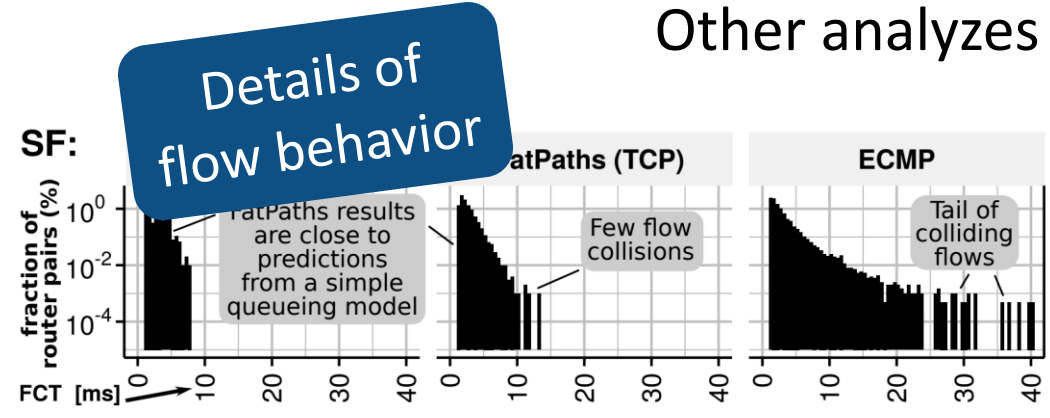
# EVALUATION



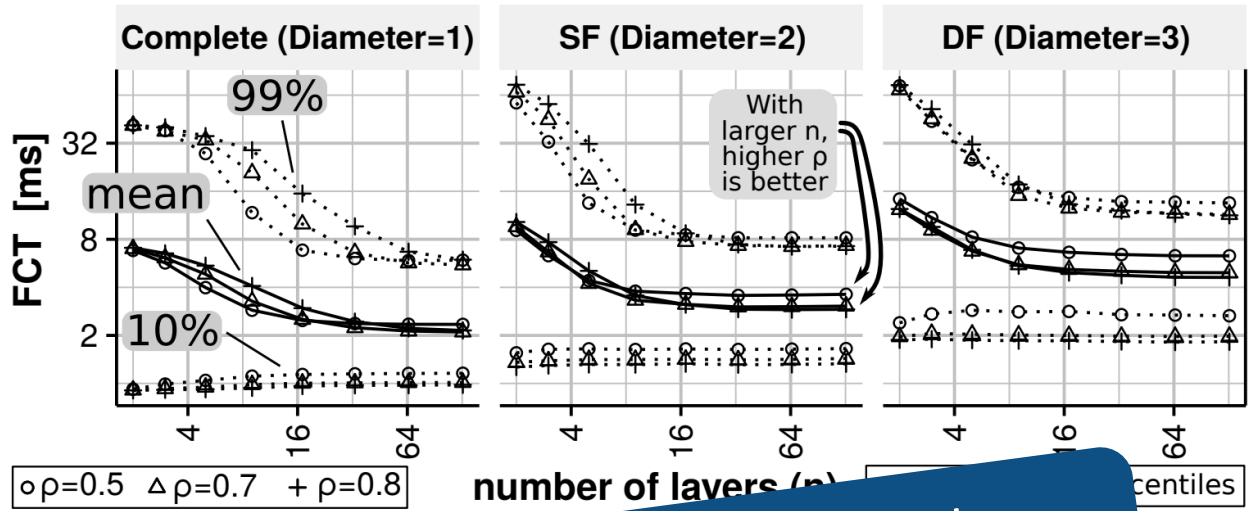
Impact from varying the number and density of layers



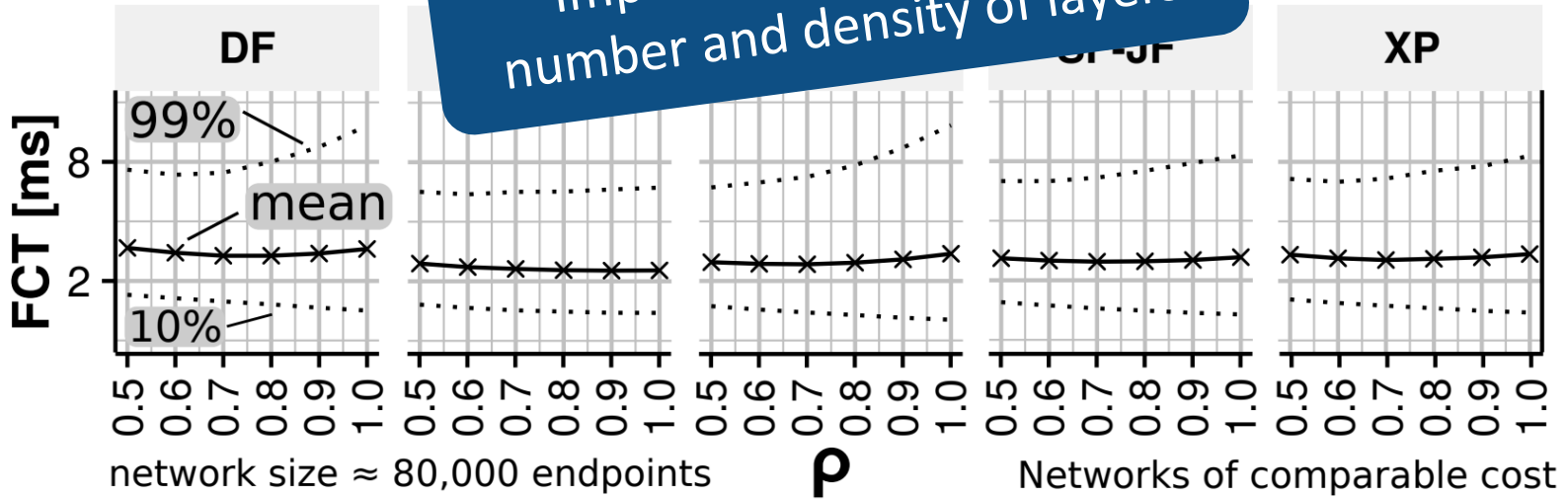
## Other analyzes



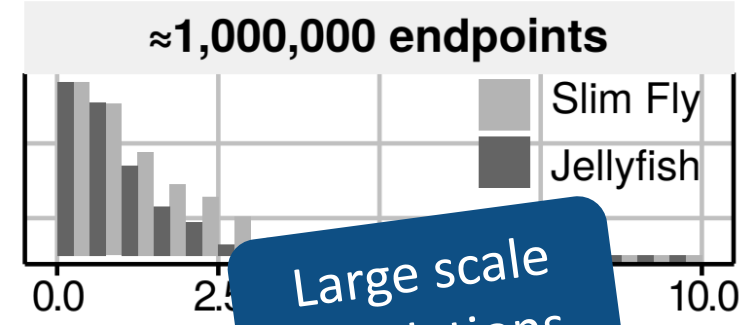
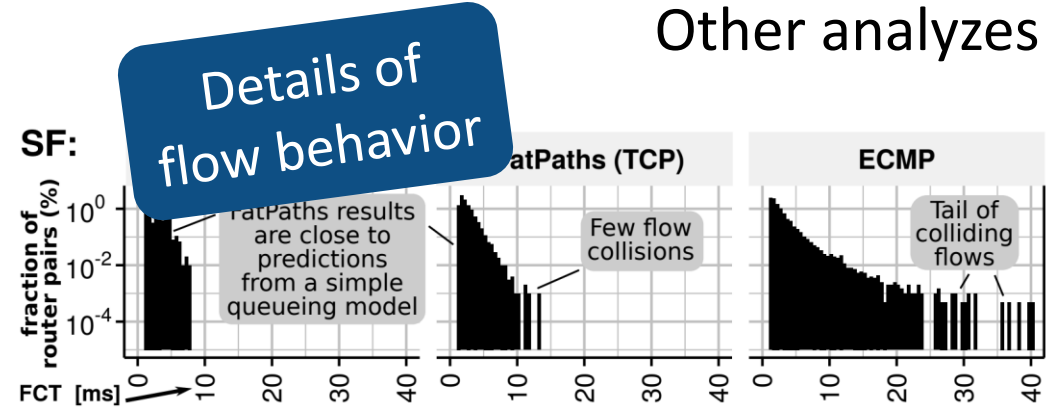
# EVALUATION



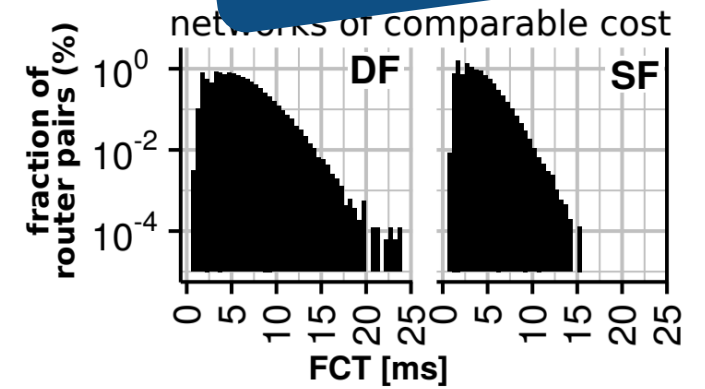
Impact from varying the number and density of layers



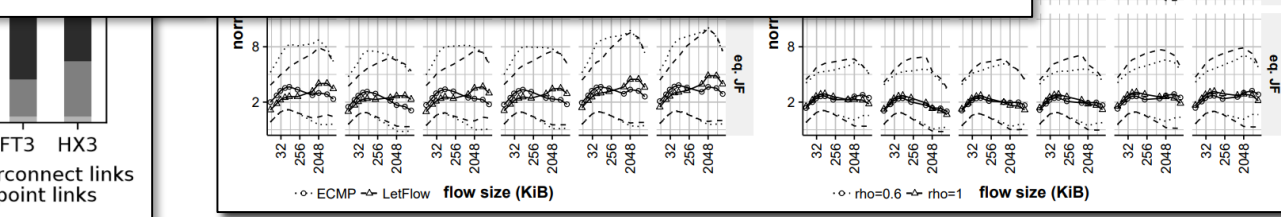
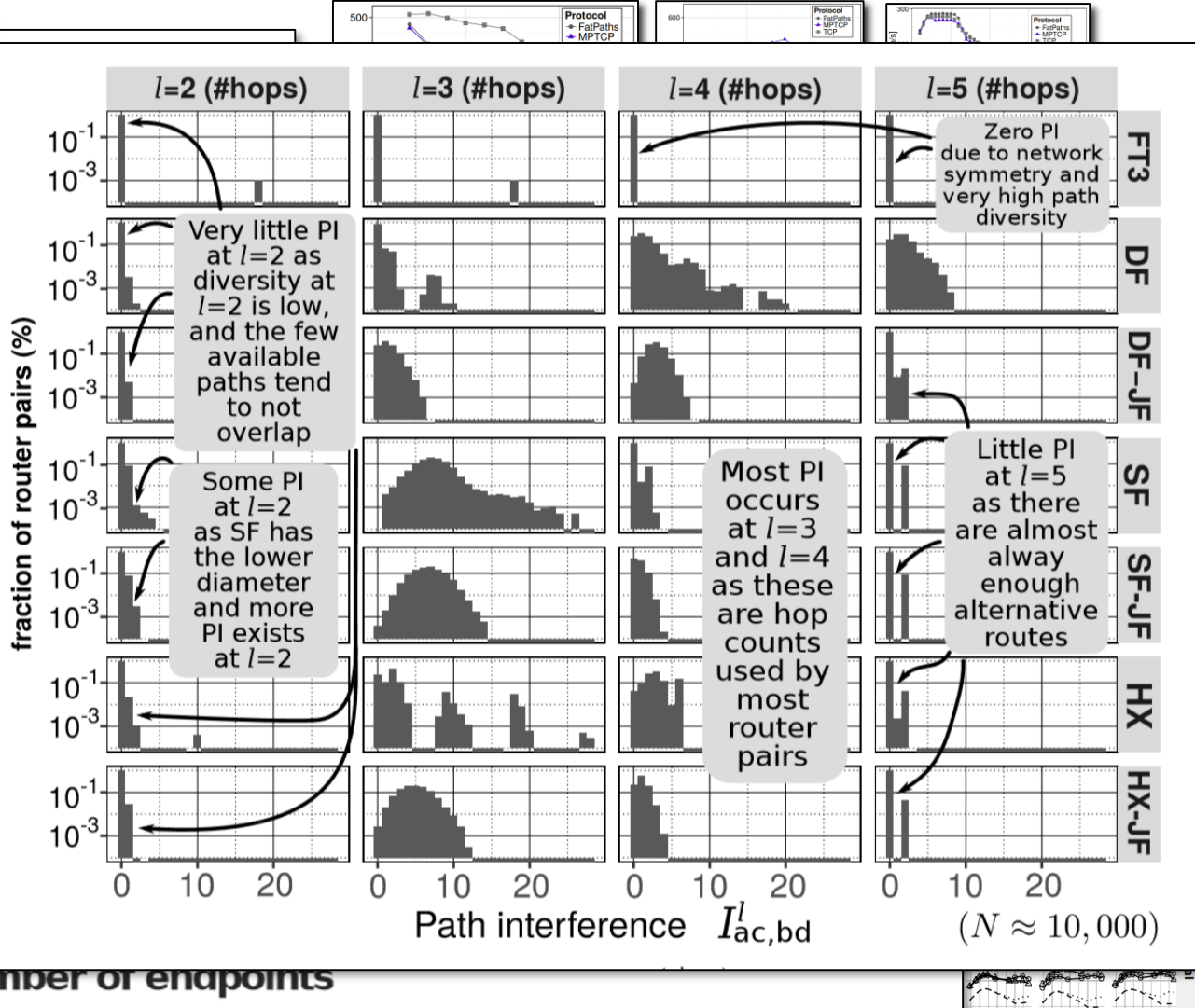
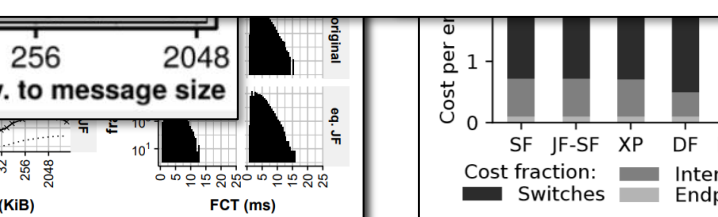
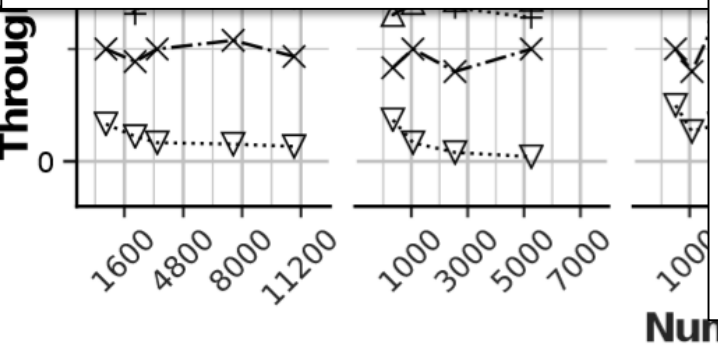
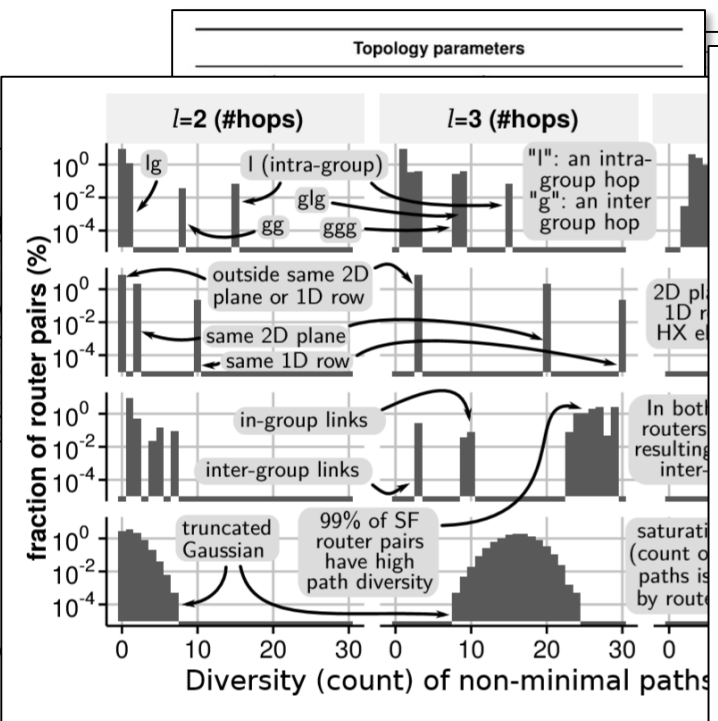
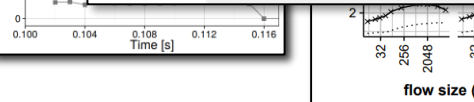
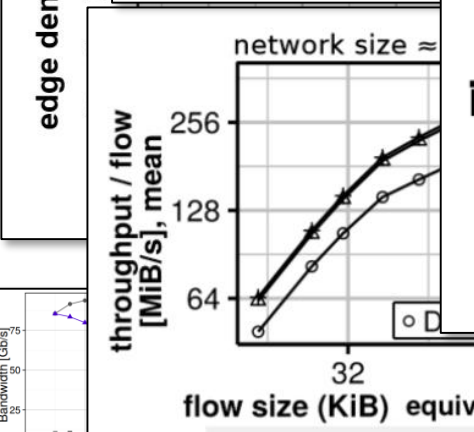
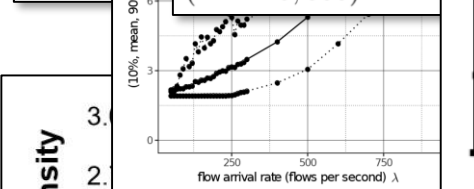
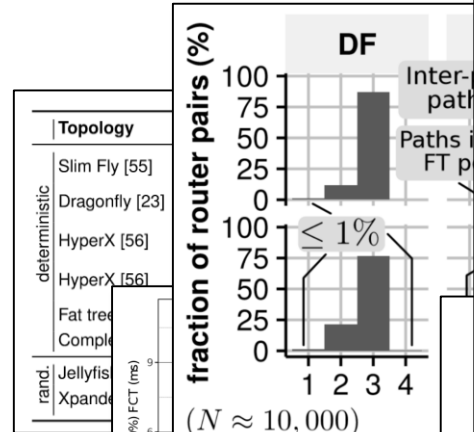
## Other analyzes



Large scale simulations

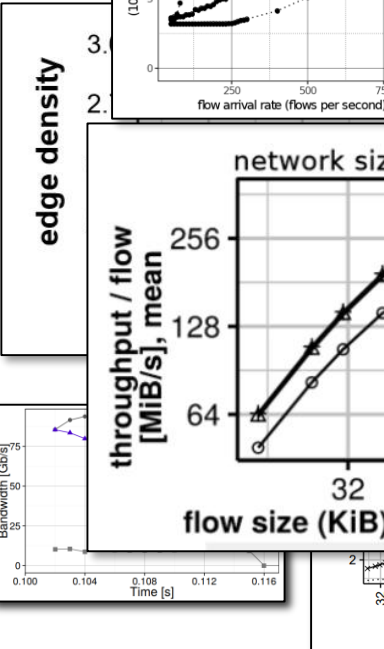
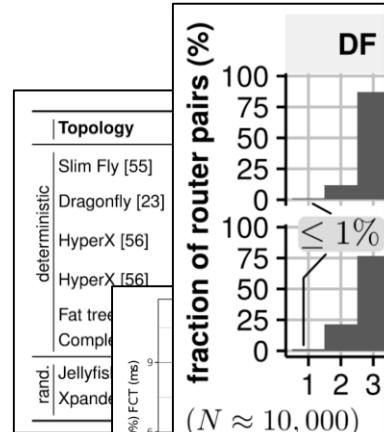


# EVALUATION



Vertical text on the right edge of the slide, partially cut off.

# EVALUATION



Topology parameters

$l=2$  (#hops)     $l=3$  (#hops)     $l=4$  (#hops)     $l=5$  (#hops)

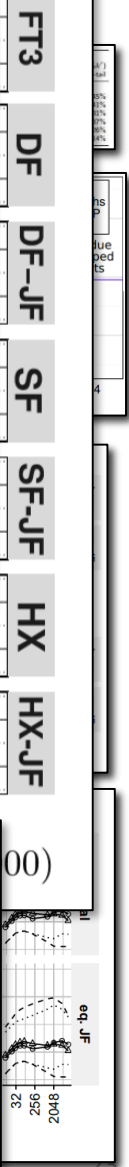
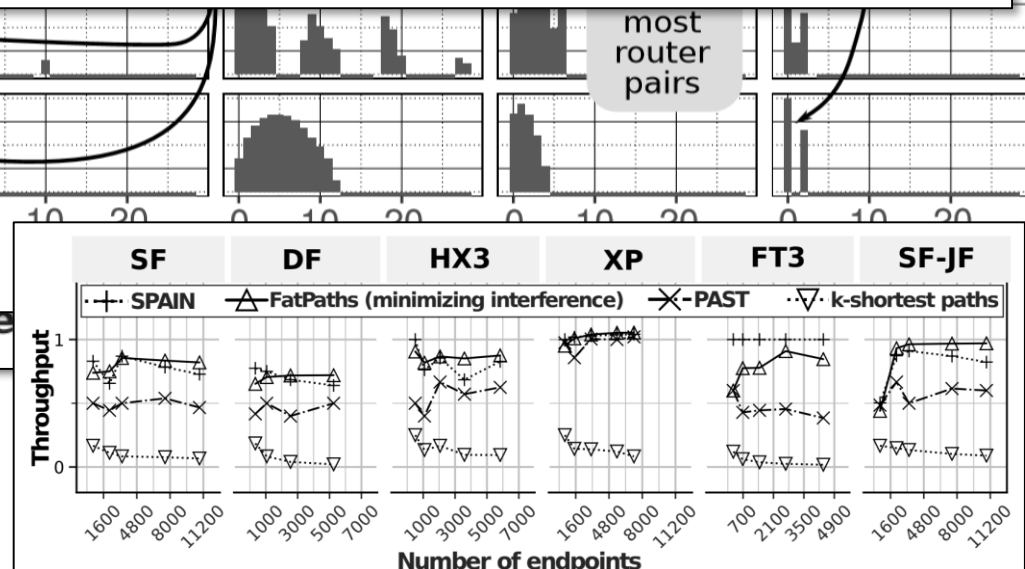
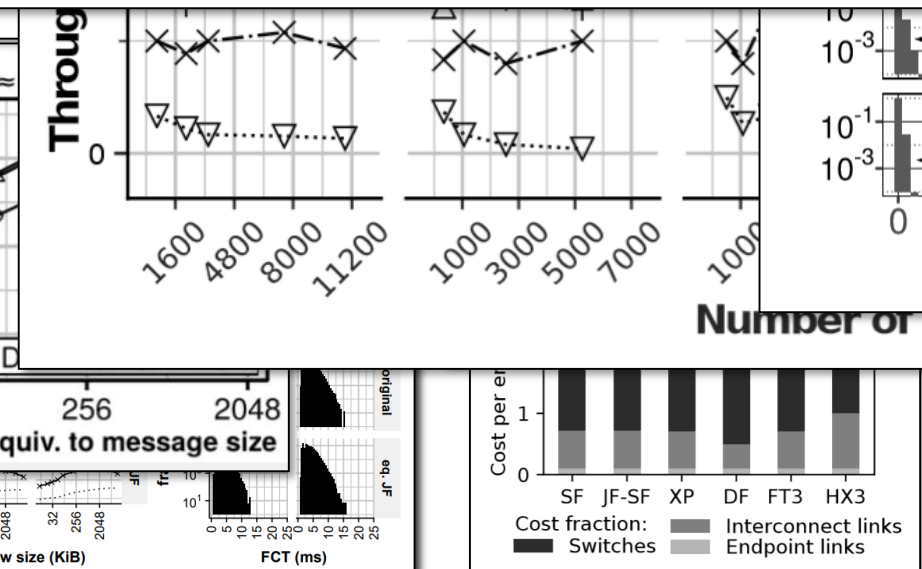
$$-\sum_{v \in V} \sum_{l=1, \dots, n} f_{is_i v l} \cdot \delta_{v, \sigma_l(s_i, t_i)} + T(s_i, t_i) \cdot \mathcal{T} \leq 0, \quad i = 1, 2, \dots, k \quad (5)$$

$$\sum_{i=1, \dots, k} \sum_{l=1, \dots, n} f_{i u v l} \cdot \delta_{v, \sigma_l(u, t_i)} \leq c(u, v), \quad \forall (u, v) \in E \quad (6)$$

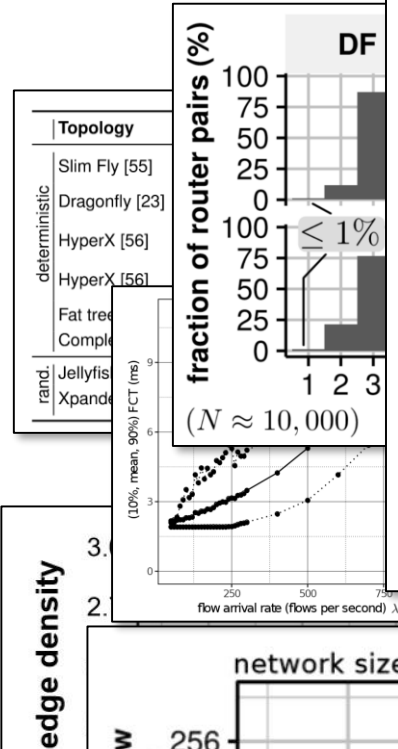
$$\sum_{v \in V} f_{i u v l} \cdot \delta_{v, \sigma_l(u, t_i)} - \sum_{v \in V} f_{i v u l} \cdot \delta_{u, \sigma_l(v, t_i)} = 0, \quad i = 1, \dots, k, \quad l = 1, \dots, n, \quad \forall u \in V \setminus \{s_i, t_i\} \quad (7)$$

$$\sum_{v \in V} \sum_{l=1, \dots, n} f_{is_i v l} \cdot \delta_{v, \sigma_l(s_i, t_i)} \leq \mathcal{T}_{upperbound} \cdot T(s_i, t_i), \quad i = 1, \dots, k \quad (8)$$

$$\sum_{v \in V} \sum_{l=1, \dots, n} f_{i v s_i l} \cdot \delta_{s_i, \sigma_l(v, t_i)} = 0, \quad i = 1, \dots, k \quad (9)$$



# EVALUATION



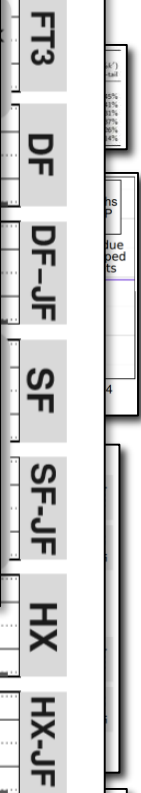
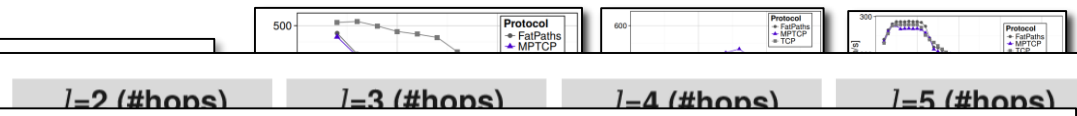
$$-\sum_{v \in V} \sum_{l=1, \dots, n} f_{is_i v l} \cdot \delta_{v, \sigma_l(s_i, t_i)} + T(s_i, t_i) \cdot \mathcal{T} \leq 0, \quad i = 1, 2, \dots, k \quad (5)$$

$$\sum_{i=1, \dots, k} \sum_{l=1, \dots, n} f_{i u v l} \cdot \delta_{v, \sigma_l(u, t_i)} \leq c(u, v), \quad \forall (u, v) \in E \quad (6)$$

$$\sum_{v \in V} f_{i u v l} \cdot \delta_{v, \sigma_l(u, t_i)} - \sum_{v \in V} f_{i v u l} \cdot \delta_{u, \sigma_l(v, t_i)} = 0, \quad i = 1, \dots, k, \quad l = 1, \dots, n, \quad \forall u \in V \setminus \{s_i, t_i\} \quad (7)$$

$$\sum_{v \in V} \sum_{l=1, \dots, n} f_{is_i v l} \cdot \delta_{v, \sigma_l(s_i, t_i)} \leq \mathcal{T}_{upperbound} \cdot T(s_i, t_i), \quad i = 1, \dots, k \quad (8)$$

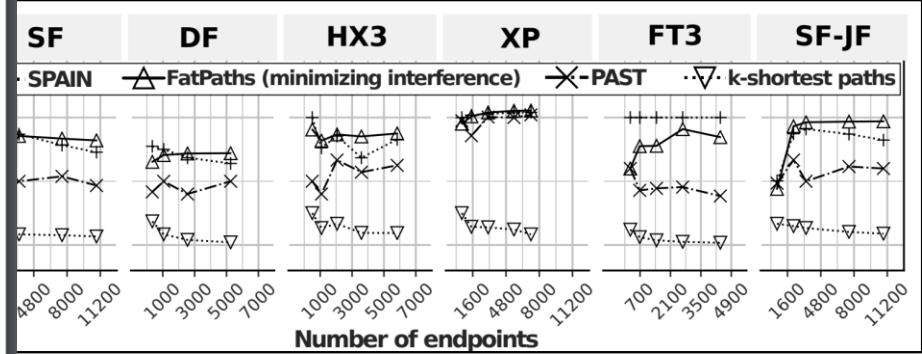
$$\sum_{v \in V} \sum_{l=1, \dots, n} f_{i v s_i l} \cdot \delta_{s_i, \sigma_l(v, t_i)} = 0, \quad i = 1, \dots, k \quad (9)$$

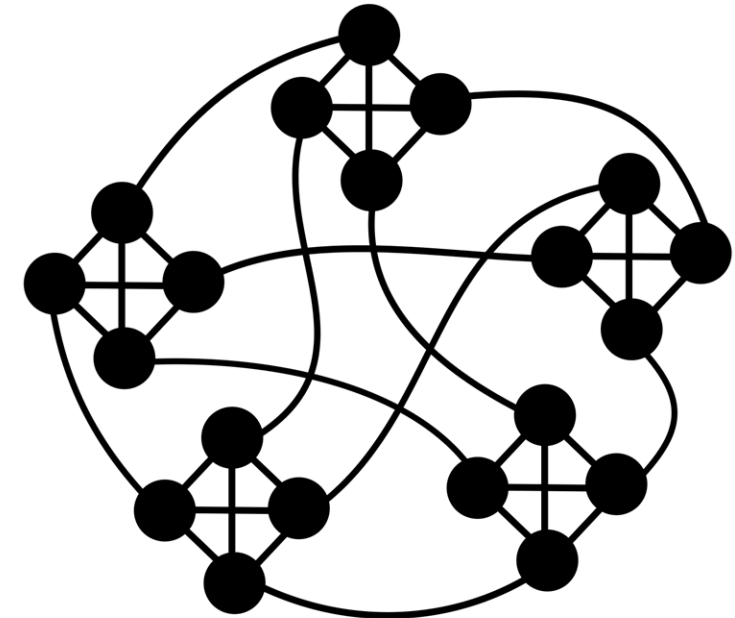
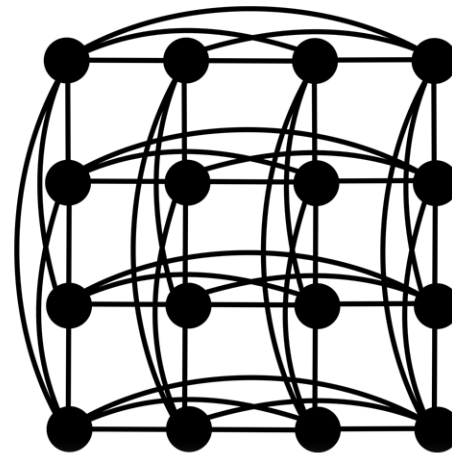
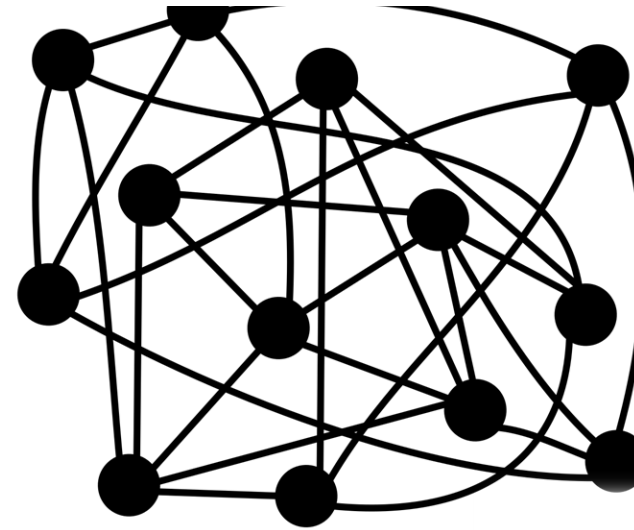
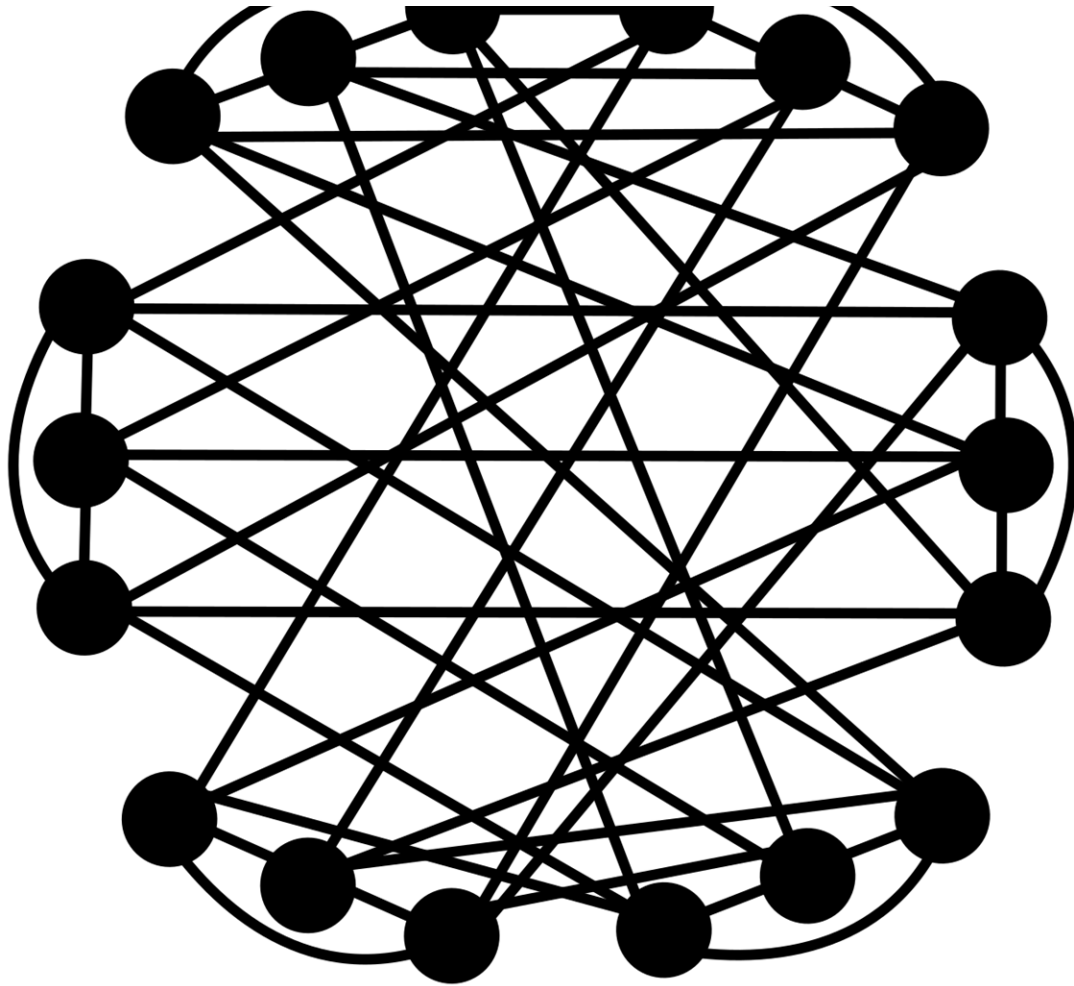


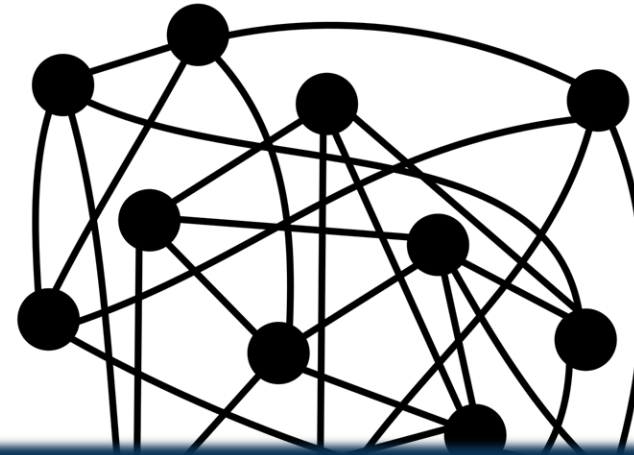
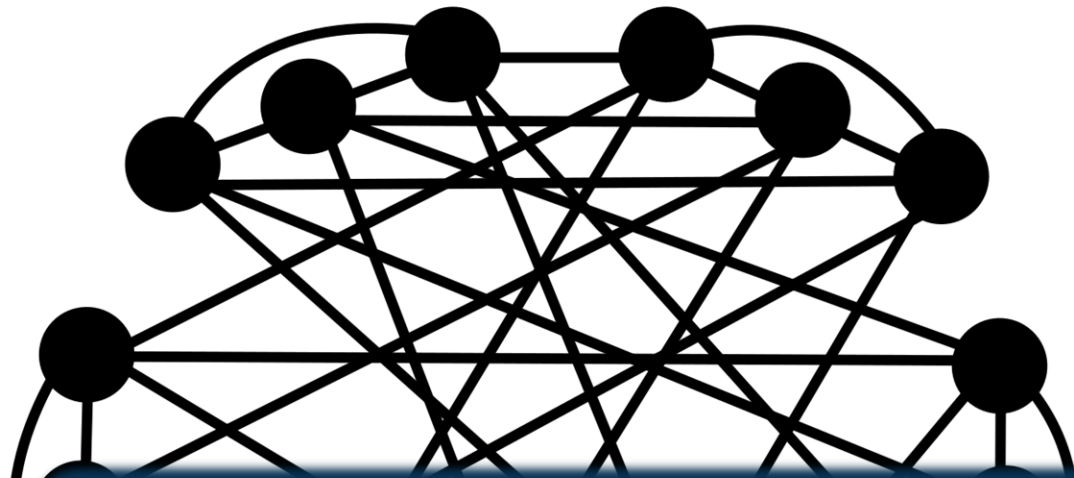
## High-Performance Routing with Multipathing and Path Diversity in Supercomputers and Data Centers

Maciej Besta<sup>1</sup>, Jens Domke<sup>2</sup>, Marcel Schneider<sup>1</sup>, Marek Konieczny<sup>3</sup>,  
 Salvatore Di Girolamo<sup>1</sup>, Timo Schneider<sup>1</sup>, Ankit Singla<sup>1</sup>, Torsten Hoefler<sup>1</sup>

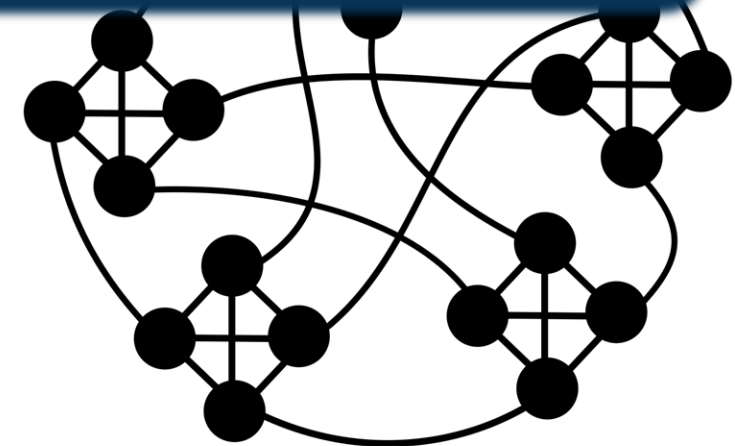
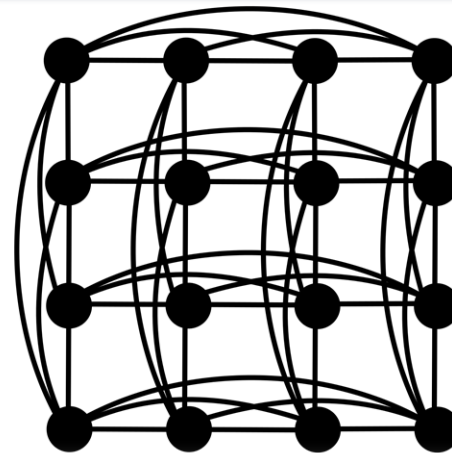
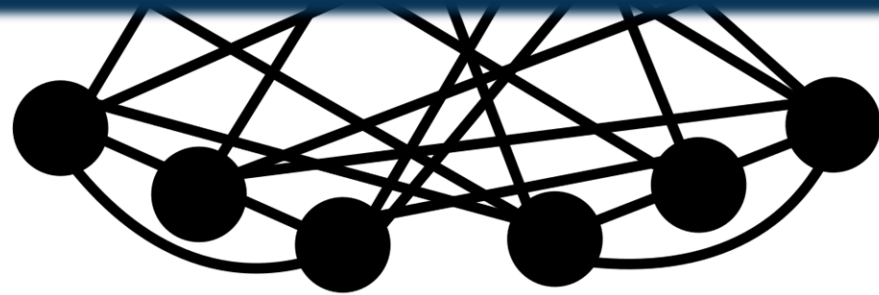
<sup>1</sup>Department of Computer Science, ETH Zurich; <sup>2</sup>RIKEN Center for Computational Science (R-CCS)  
<sup>3</sup>Department of Computer Science, Electronics and Telecommunications; AGH-UST







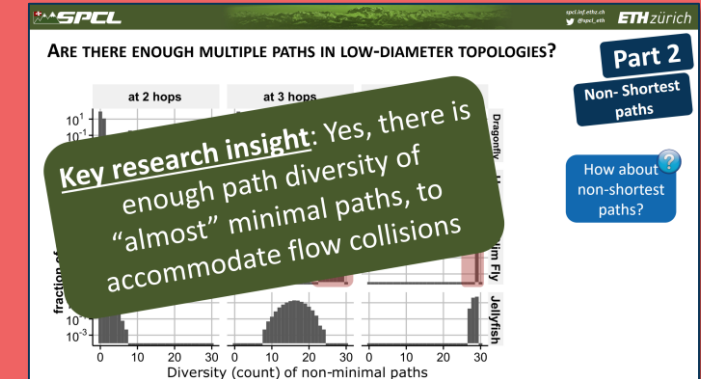
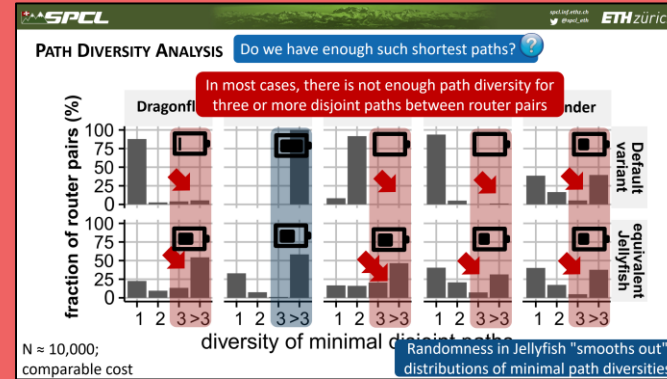
TAKEAWAY MESSAGE: FATPATHS ENABLES HIGH-PERFORMANCE  
MULTIPATH ROUTING IN LOW-DIAMETER NETWORKS



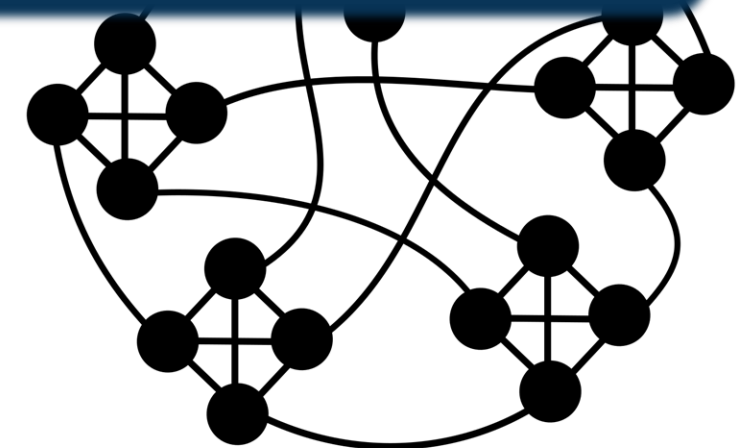
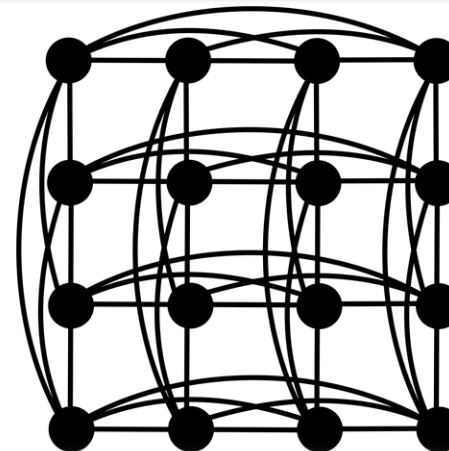
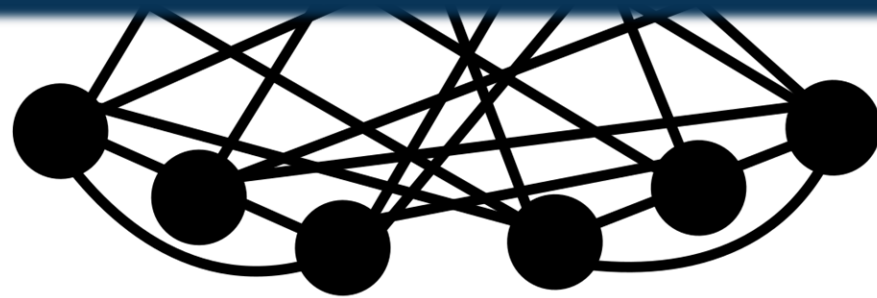
## POTENTIAL FOR MULTIPATH ROUTING IN MODERN LOW-DIAMETER NETWORKS

LOW-DIAMETER NETWORK TOPOLOGIES

How to route modern low diameter topologies?



**TAKEAWAY MESSAGE: FATPATHS ENABLES HIGH-PERFORMANCE MULTIPATH ROUTING IN LOW-DIAMETER NETWORKS**

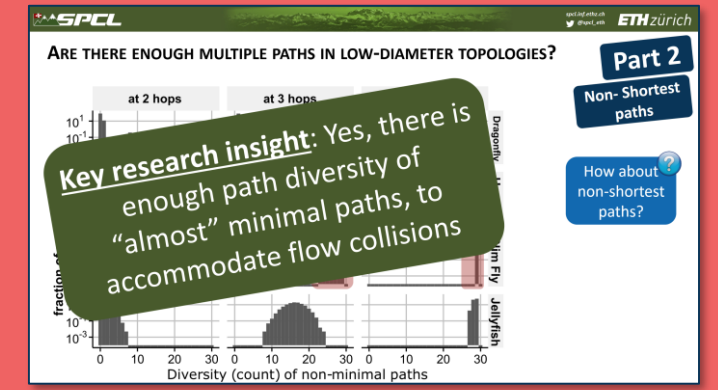
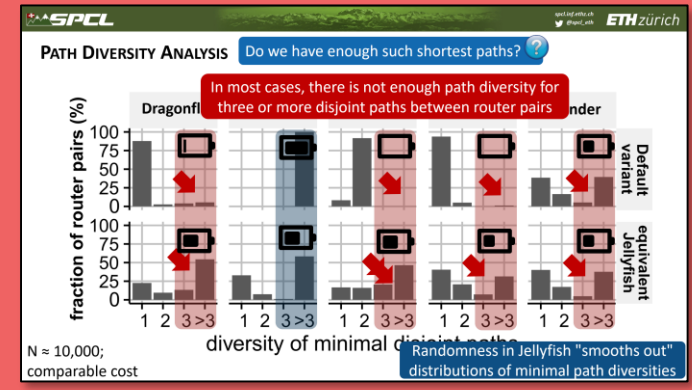




# POTENTIAL FOR MULTIPATH ROUTING IN MODERN LOW-DIAMETER NETWORKS

**LOW-DIAMETER NETWORK TOPOLOGIES**

How to route modern low diameter topologies?



# TAKEAWAY MESSAGE: FATPATHS ENABLES HIGH-PERFORMANCE MULTIPATH ROUTING IN LOW-DIAMETER NETWORKS

## HIGH-PERFORMANCE ROUTING PROTOCOL & ARCHITECTURE

**LAYERED ROUTING** How to encode & use diversity of shortest and non-minimal paths?

(a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

(b) Divide one flow into subflows and send subflows across different layers

(c) Route minimally in each layer

Default topology Shortest path: 2 hops

Layer 1: "Almost"-shortest path: 3 hops

Layer 2: "Almost"-shortest path: 3 hops

**FATPATHS ARCHITECTURE**

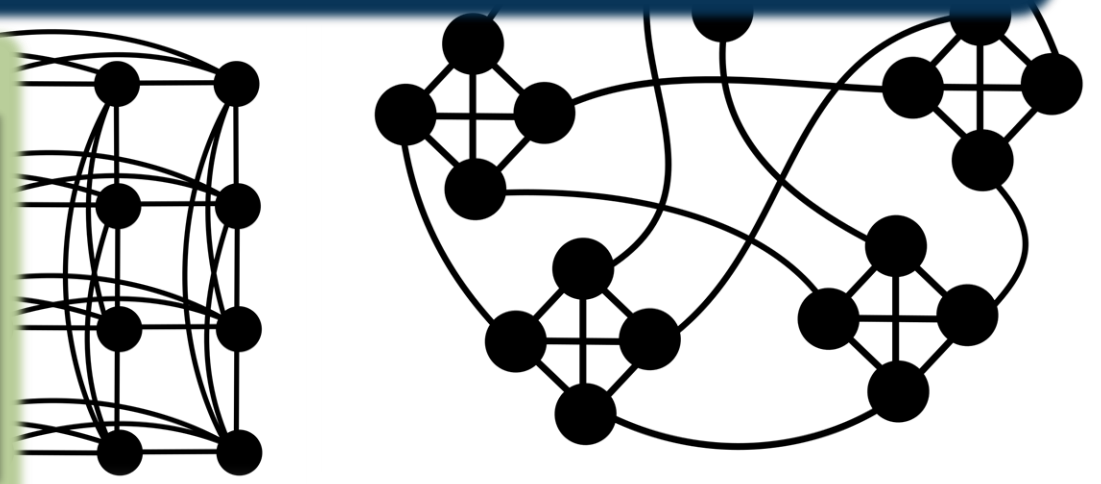
Design details

Check the

Key design insight: Layered routing enables (1) easy encoding of the diversity of multiple paths, and (2) enables simple and robust load balancing

Deploying layers (VLAN, ...)

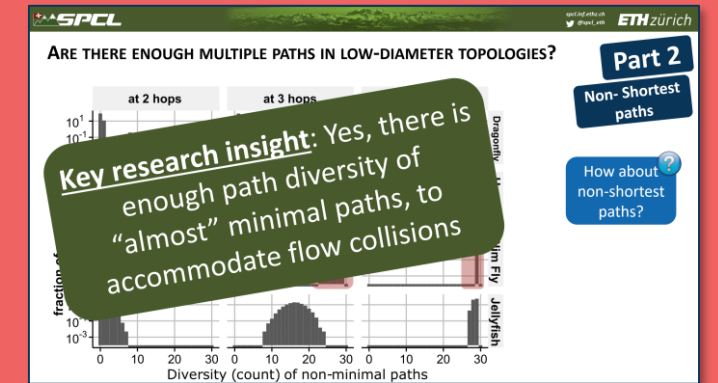
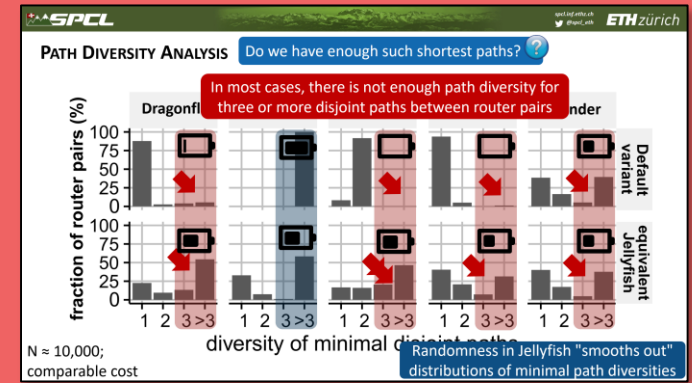
Populating forwarding entries



# POTENTIAL FOR MULTIPATH ROUTING IN MODERN LOW-DIAMETER NETWORKS

**LOW-DIAMETER NETWORK TOPOLOGIES**

How to route modern low diameter topologies?



# TAKEAWAY MESSAGE: FATPATHS ENABLES HIGH-PERFORMANCE MULTIPATH ROUTING IN LOW-DIAMETER NETWORKS

## HIGH-PERFORMANCE ROUTING PROTOCOL & ARCHITECTURE

**LAYERED ROUTING** How to encode & use diversity of shortest and non-minimal paths?

- (a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links
- (a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)
- (b) Divide one flow into subflows and send subflows across different layers
- (c) Route minimally in each layer

**Default topology**  
Shortest path: 2 hops

**Layer 1: "Almost"-shortest path: 3 hops**

**Layer 2: "Almost"-shortest path: 3 hops**

**FATPATHS ARCHITECTURE**

Design details

Check the

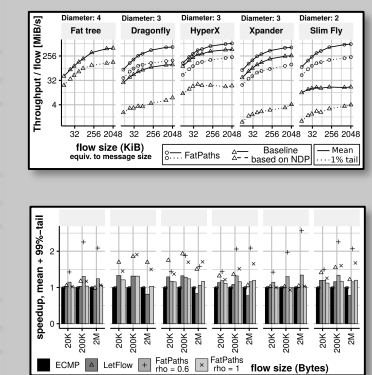
Purified Transport based on NDP

Key design insight: Layered routing enables (1) easy encoding of the diversity of multiple paths, and (2) enables simple and robust load balancing

Deploying layers (VLAN, ...)

Populating forwarding entries

## RICH EVALUATION & THEORETICAL ANALYSIS



**EVALUATION**

Legend: SF, DF, HKX, XP, FT3, SF-FP

Legend: SF, DF, HKX, XP, FT3, SF-FP

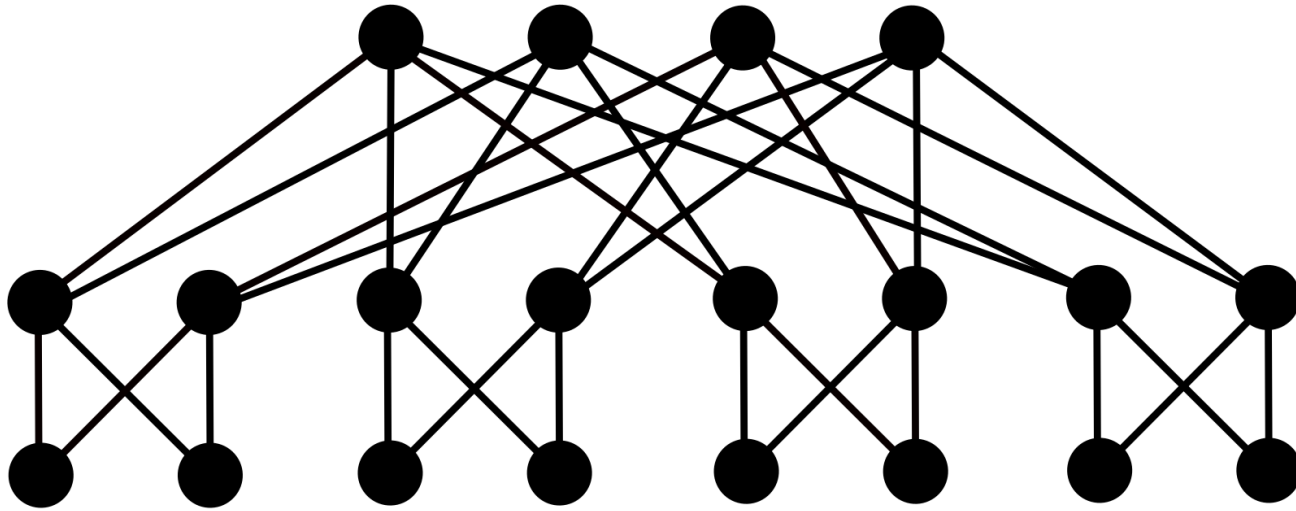
$$-\sum_{i=1}^k \sum_{u,v \in V} f_{u,v} \cdot \delta_{i,n}(u,v) + T(u,v) \cdot T \leq 0, \quad i = 1, 2, \dots, k \quad (5)$$

$$\sum_{i=1}^k \sum_{u,v \in V} f_{u,v} \cdot \delta_{i,n}(u,v) \leq c(u,v), \quad \forall (u,v) \in E \quad (6)$$

$$\sum_{i=1}^k \sum_{u,v \in V} f_{u,v} \cdot \delta_{i,n}(u,v) - \sum_{i=1}^k \sum_{u,v \in V} f_{u,v} \cdot \delta_{i,n}(u,v) = 0, \quad i = 1, \dots, k, \quad \forall u,v \in V \setminus \{s_i, t_i\} \quad (7)$$

## Backup Slides and Slides' Variants

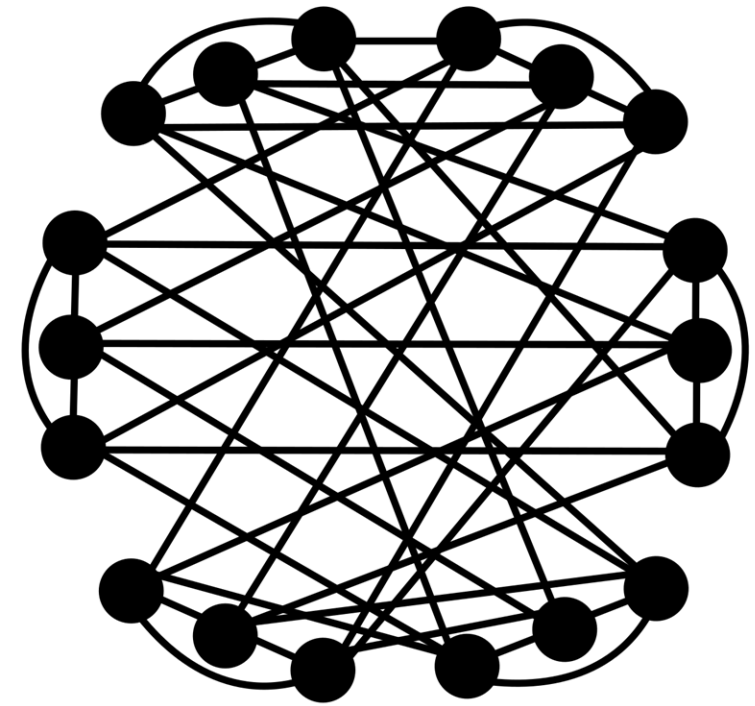
# LOW-DIAMETER NETWORK TOPOLOGIES VS FAT TREES



3-level Fat tree

Diameter = 4

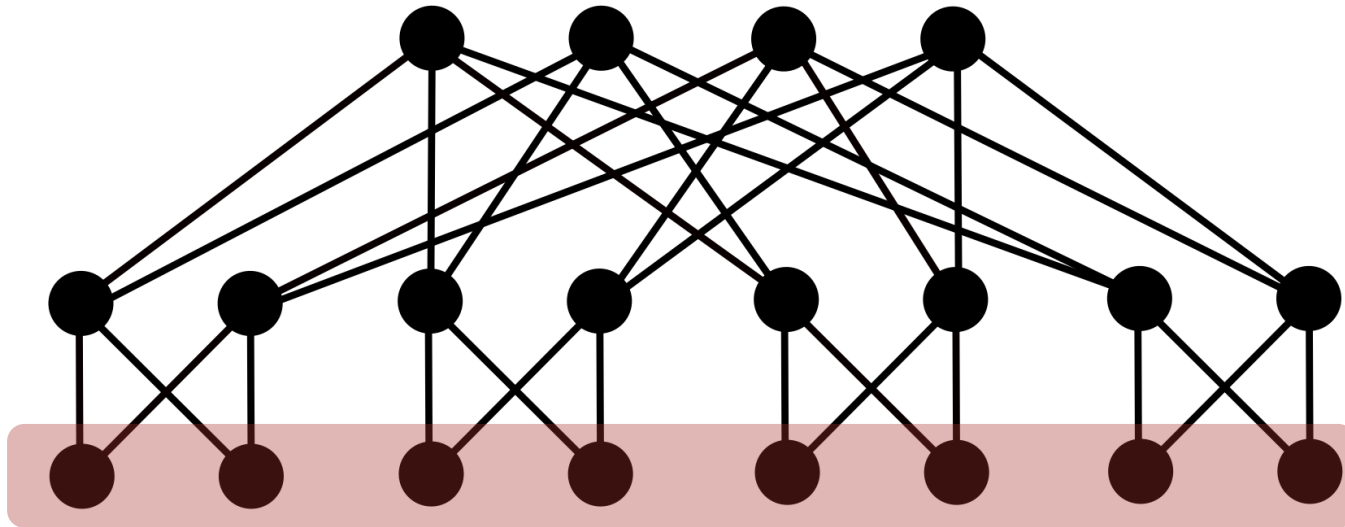
Example Slim Fly (based on the Hoffman-Singleton Graph [1])



Diameter = 2

[1] Hoffman, Alan J.; Singleton, Robert R. (1960), *Moore graphs with diameter 2 and 3*, IBM Journal of Research and Development

# LOW-DIAMETER NETWORK TOPOLOGIES VS FAT TREES

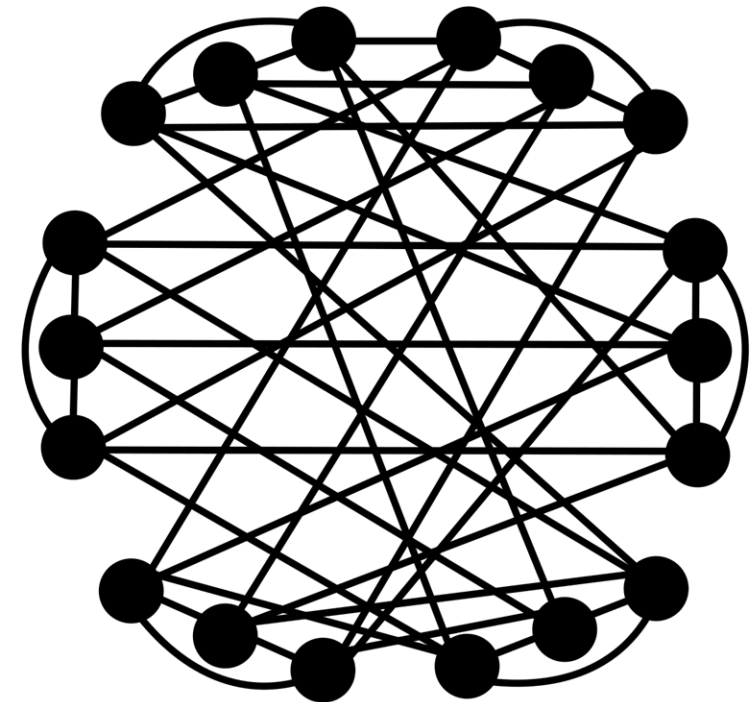


3-level Fat tree

Diameter = 4

Only edge (leaf) routers attach to endpoints

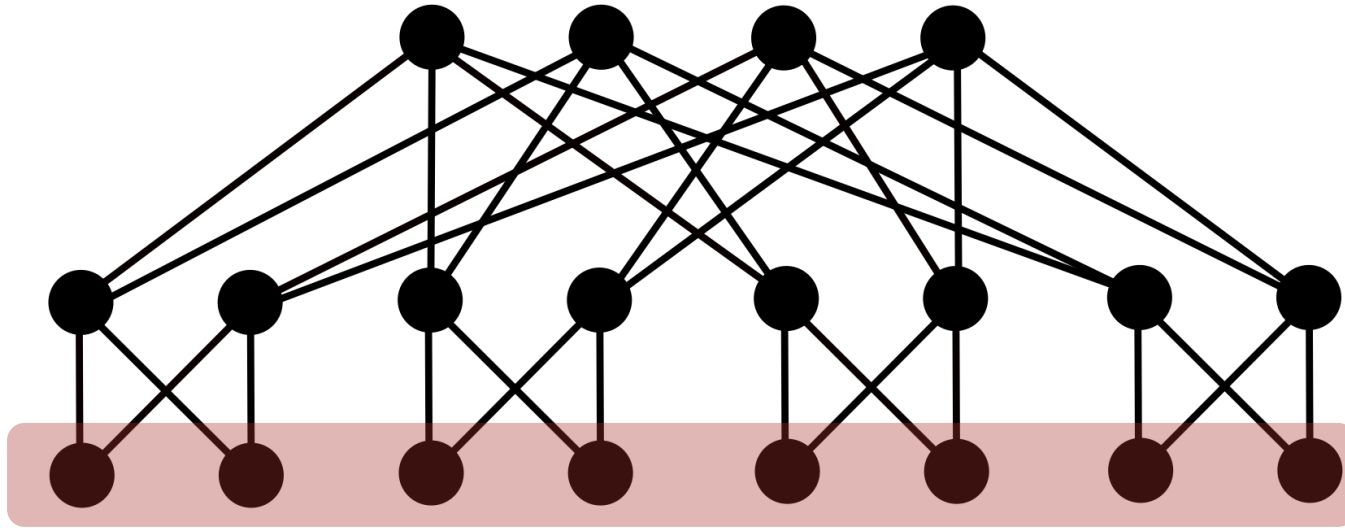
Example Slim Fly (based on the Hoffman-Singleton Graph [1])



Diameter = 2

[1] Hoffman, Alan J.; Singleton, Robert R. (1960), *Moore graphs with diameter 2 and 3*, IBM Journal of Research and Development

# LOW-DIAMETER NETWORK TOPOLOGIES VS FAT TREES



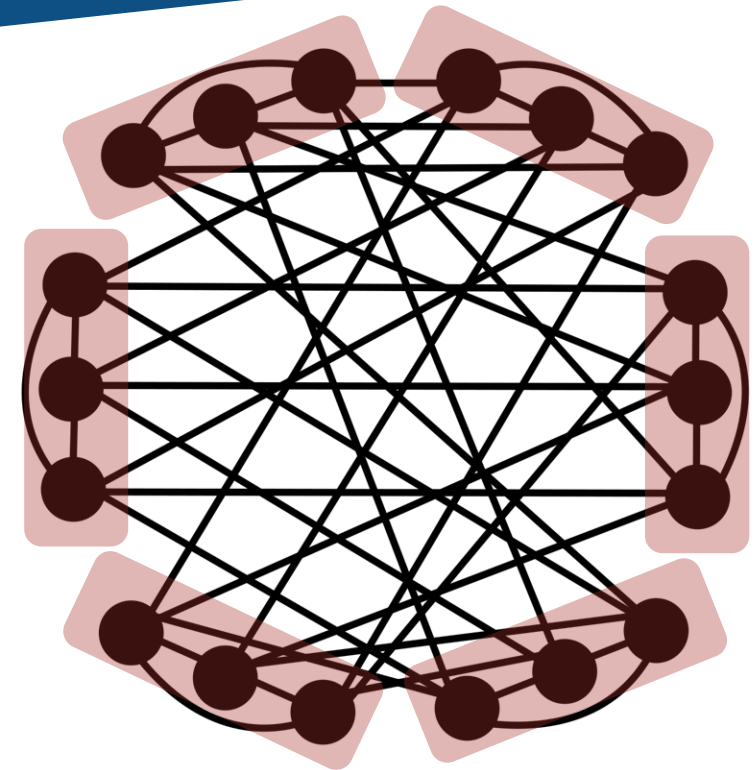
3-level Fat tree

Diameter = 4

Only edge (leaf) routers attach to endpoints

All routers attach to endpoints

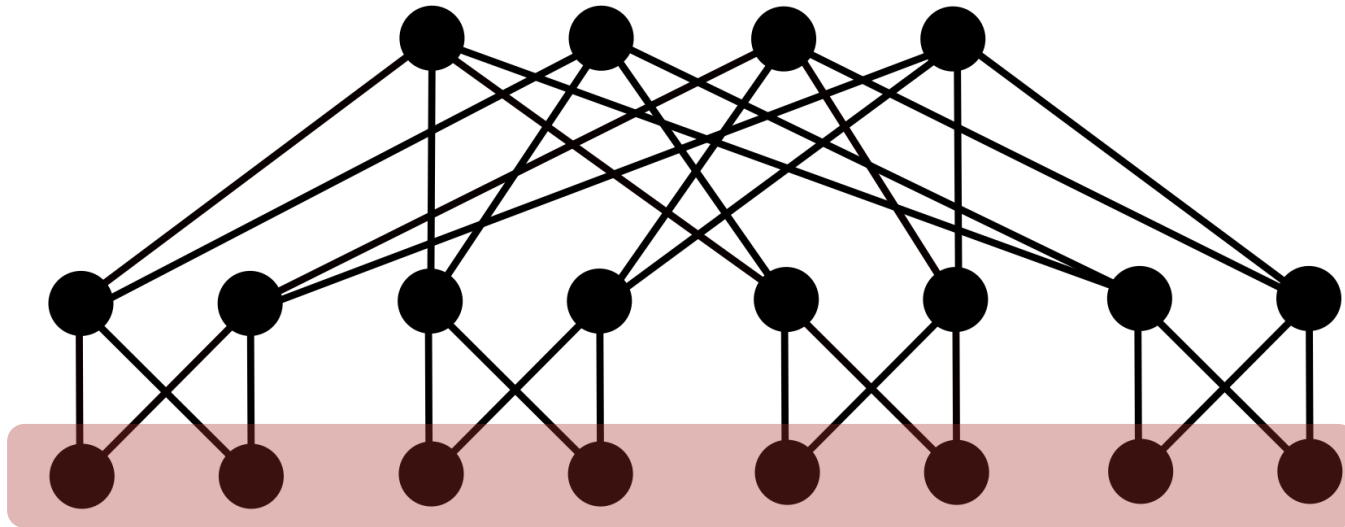
Example Slim Fly (based on the Hoffman-Singleton Graph [1])



Diameter = 2

[1] Hoffman, Alan J.; Singleton, Robert R. (1960), *Moore graphs with diameter 2 and 3*, IBM Journal of Research and Development

# LOW-DIAMETER NETWORK TOPOLOGIES VS FAT TREES

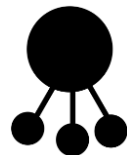


3-level Fat tree

Diameter = 4

Only edge (leaf) routers attach to endpoints

Lower concentration

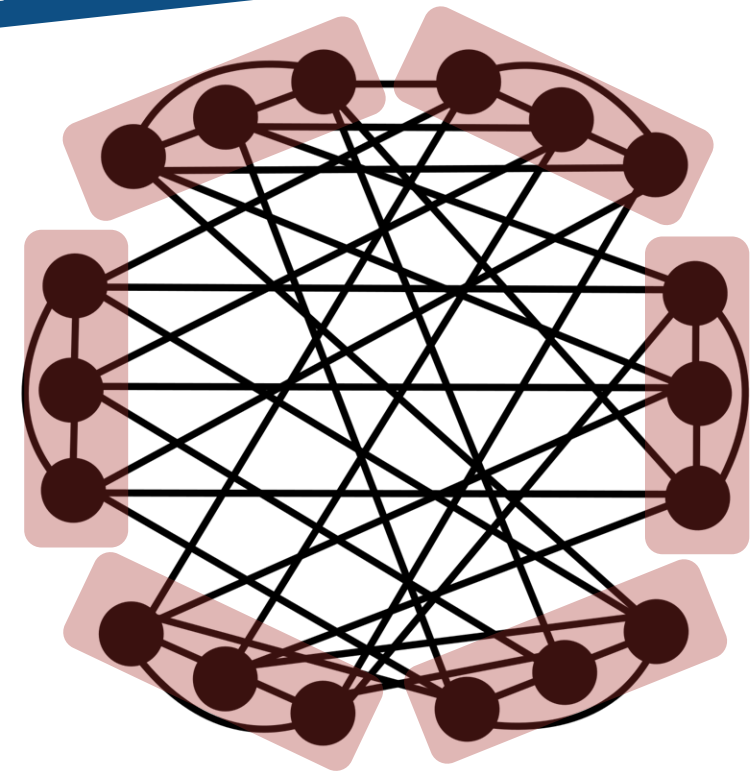


All routers attach to endpoints

Higher concentration



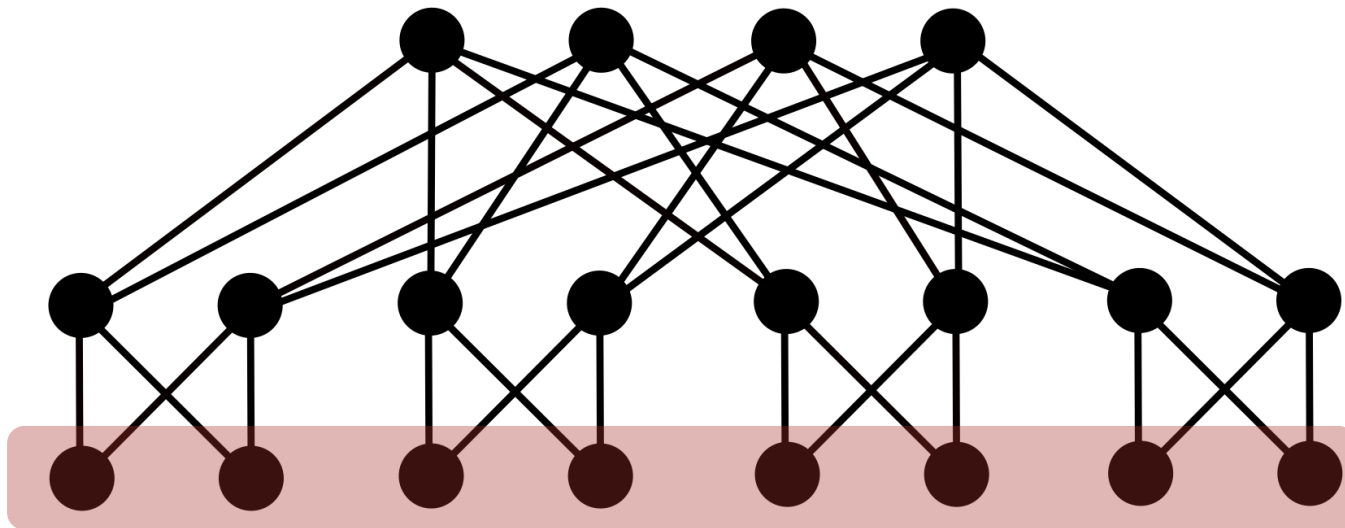
Example Slim Fly (based on the Hoffman-Singleton Graph [1])



Diameter = 2

[1] Hoffman, Alan J.; Singleton, Robert R. (1960), *Moore graphs with diameter 2 and 3*, IBM Journal of Research and Development

# LOW-DIAMETER NETWORK TOPOLOGIES VS FAT TREES

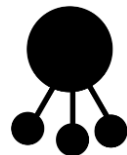


3-level Fat tree

Diameter = 4

Only edge (leaf) routers attach to endpoints

Lower concentration



All routers attach to endpoints

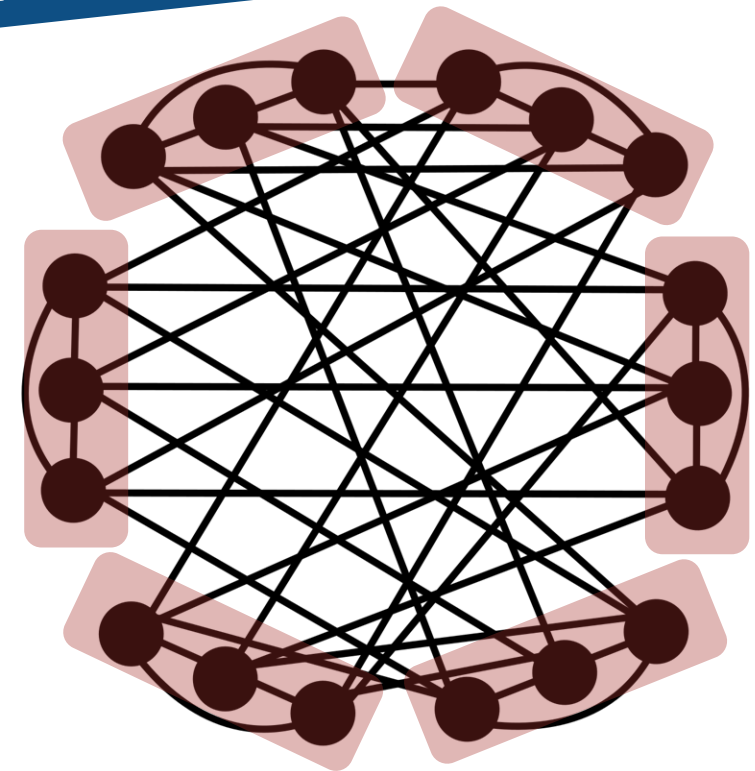
> ≈50% fewer routers

> ≈30% fewer cables

Higher concentration



Example Slim Fly (based on the Hoffman-Singleton Graph [1])



Diameter = 2

[1] Hoffman, Alan J.; Singleton, Robert R. (1960), *Moore graphs with diameter 2 and 3*, IBM Journal of Research and Development



# MULTIPATH ROUTING: MOTIVATION



## MULTIPATH ROUTING: MOTIVATION



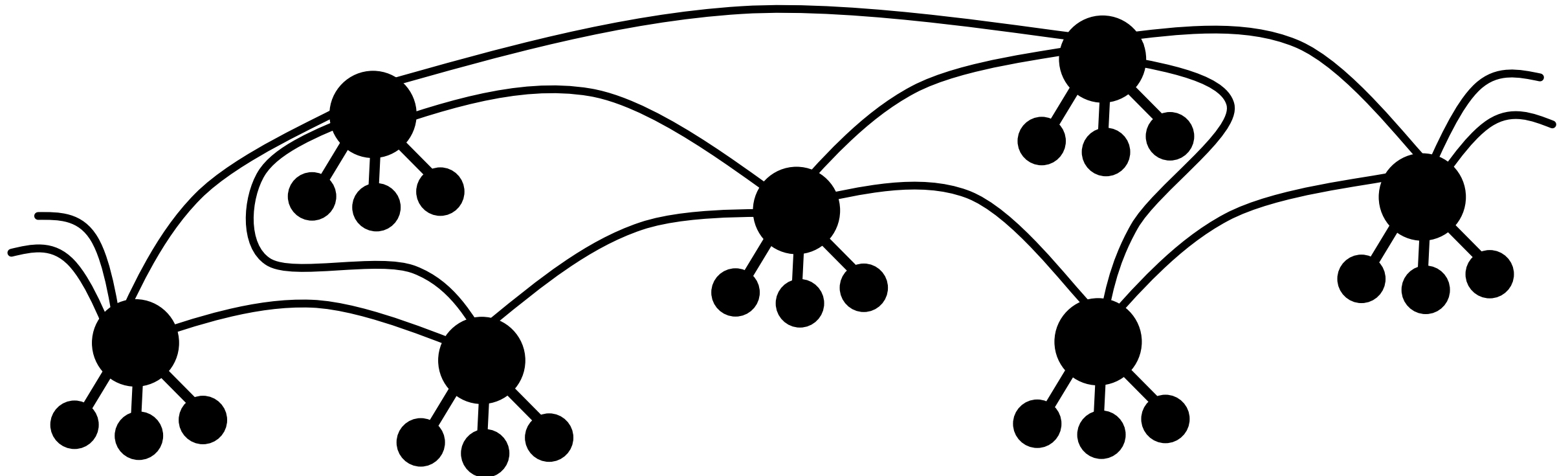
What are the problems  
that we want to tackle  
with multipathing?



## MULTIPATH ROUTING: MOTIVATION



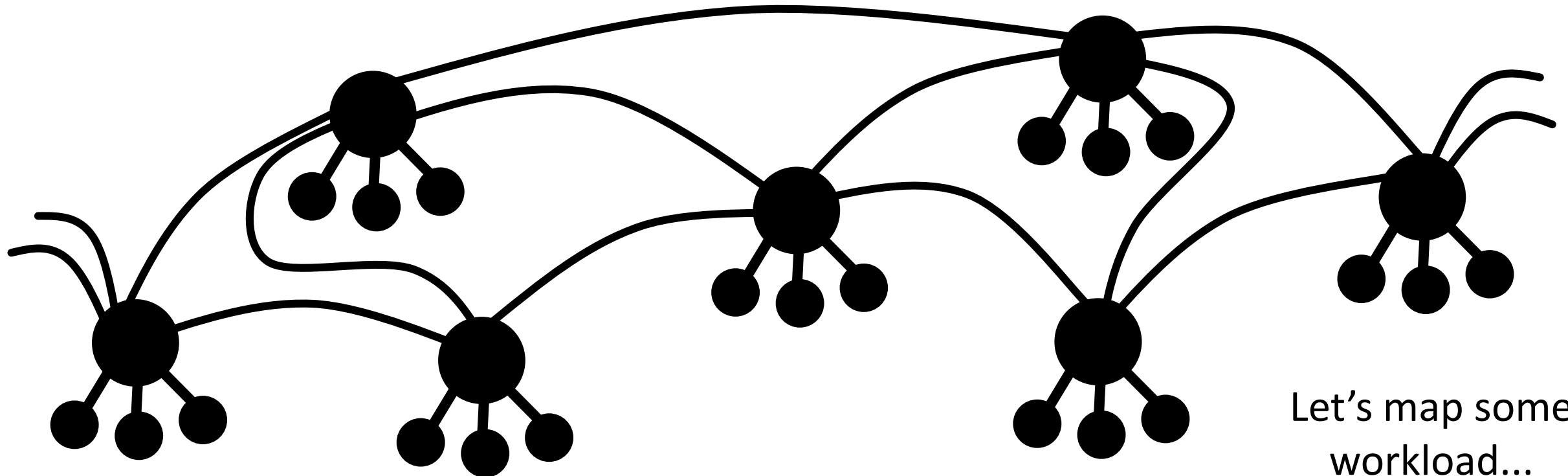
What are the problems that we want to tackle with multipathing?



## MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?

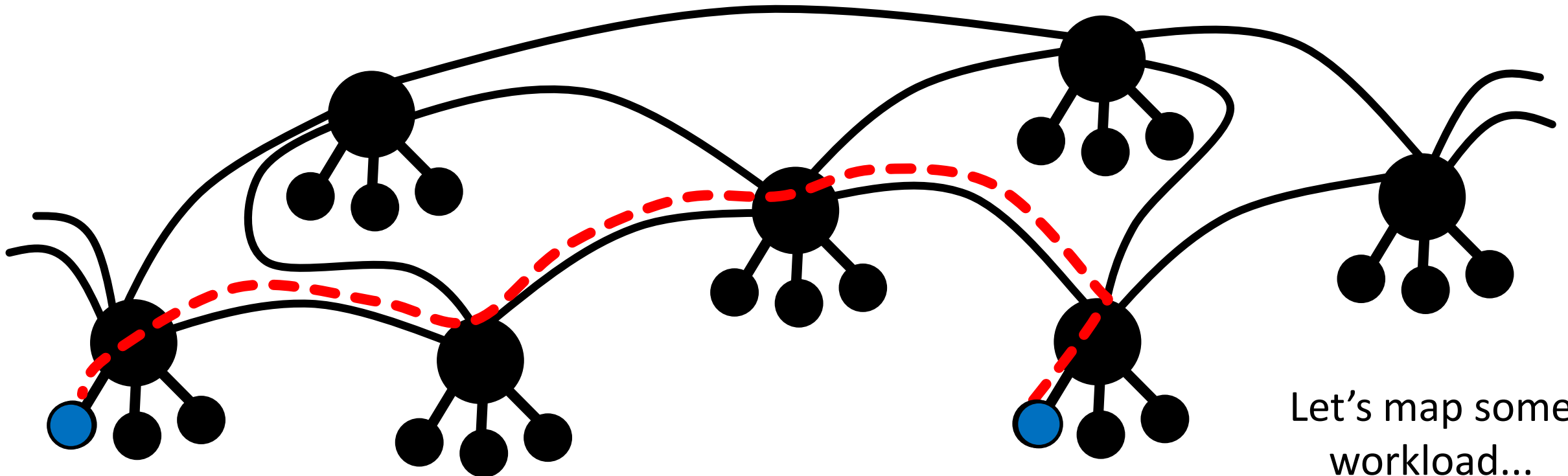


Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?

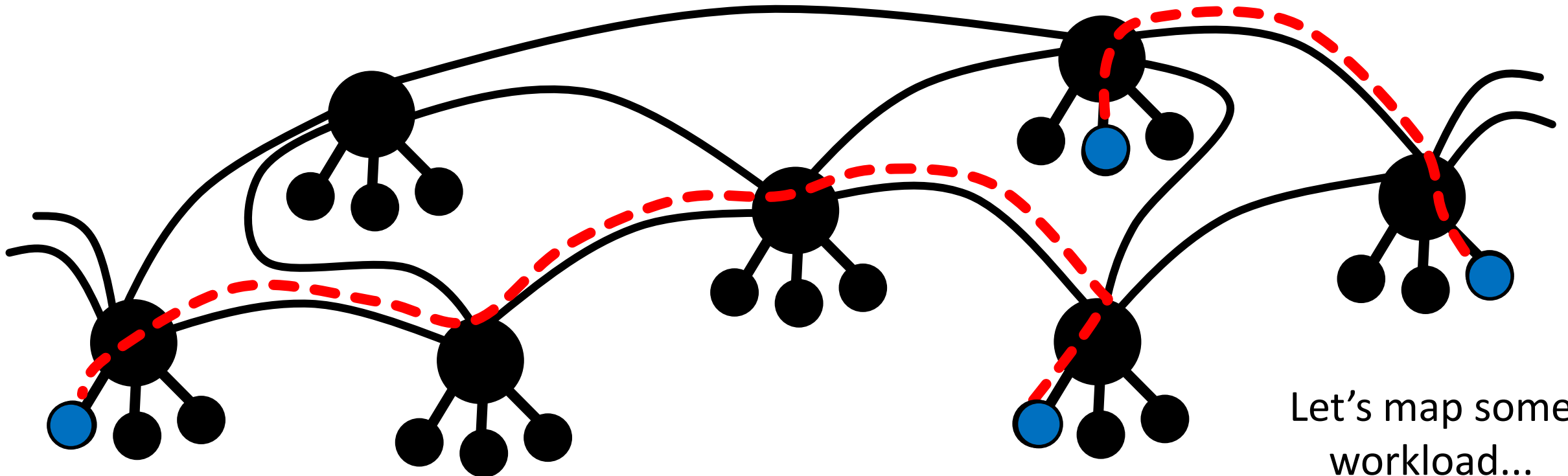


Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?

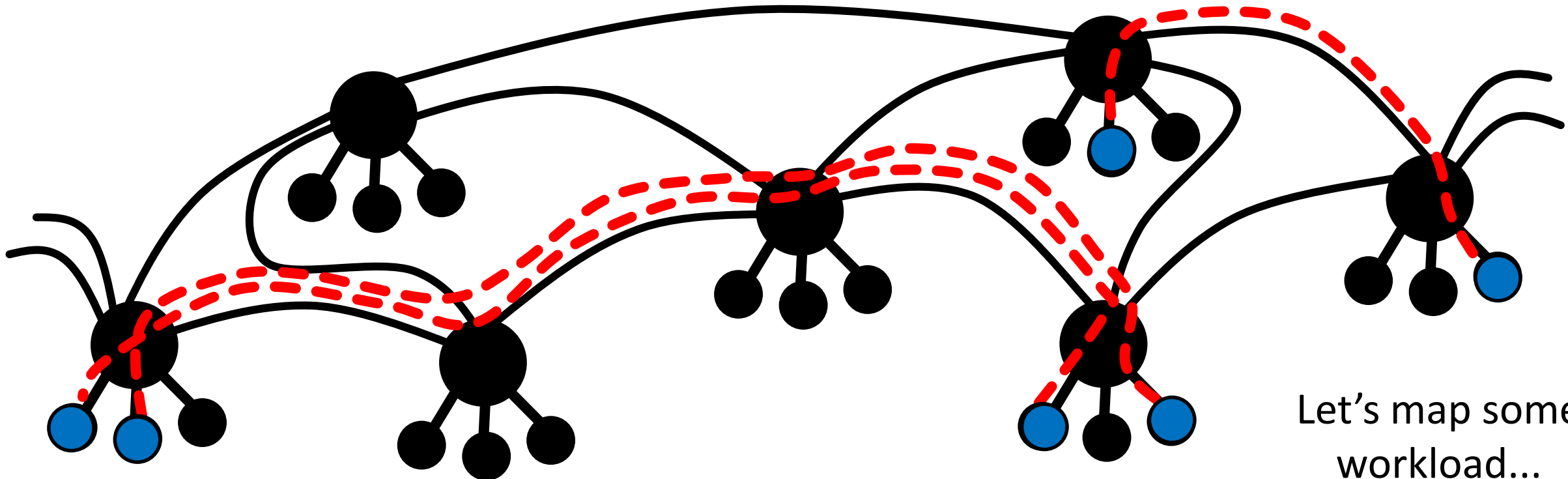


Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



What are the problems that we want to tackle with multipathing?



Let's map some workload...

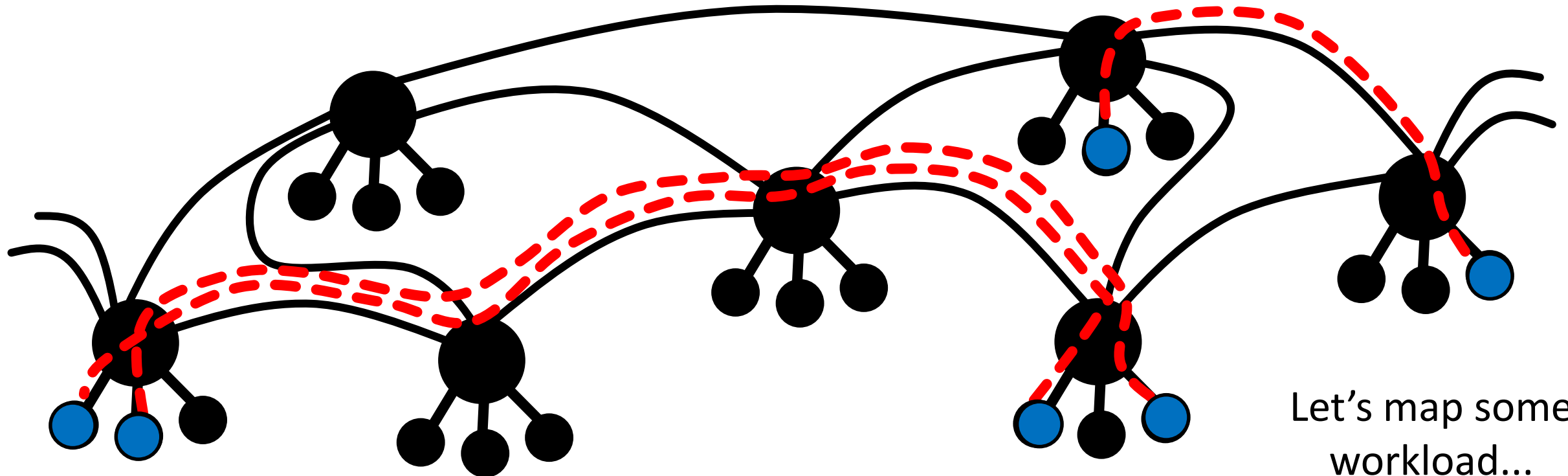
# MULTIPATH ROUTING: MOTIVATION



**Flows collide!**



What are the problems that we want to tackle with multipathing?



Let's map some workload...

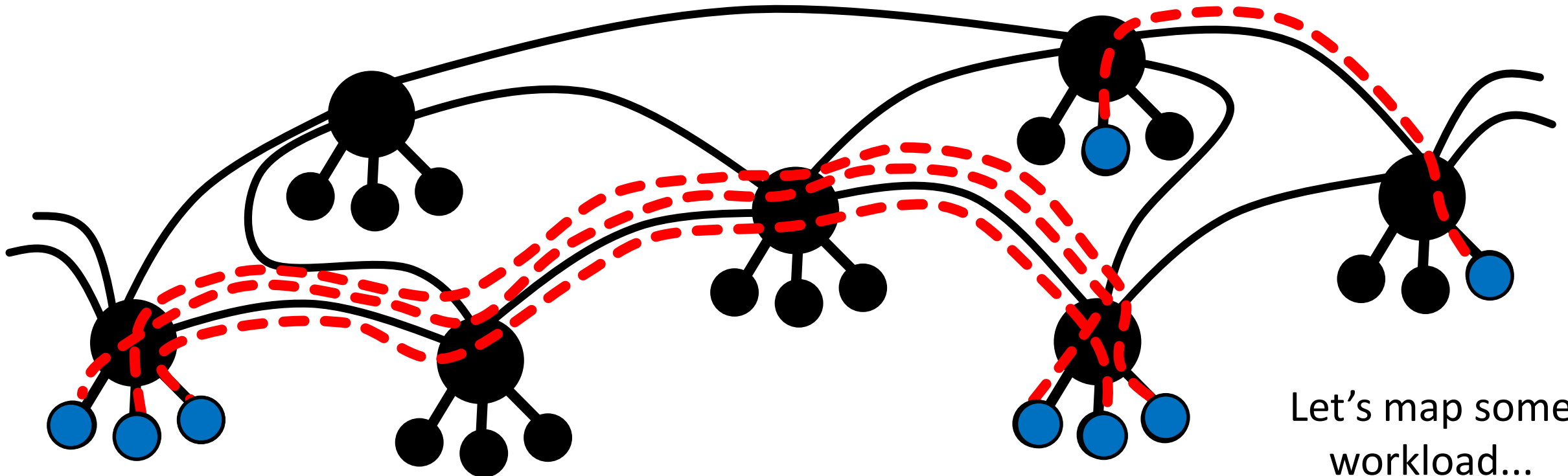


# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



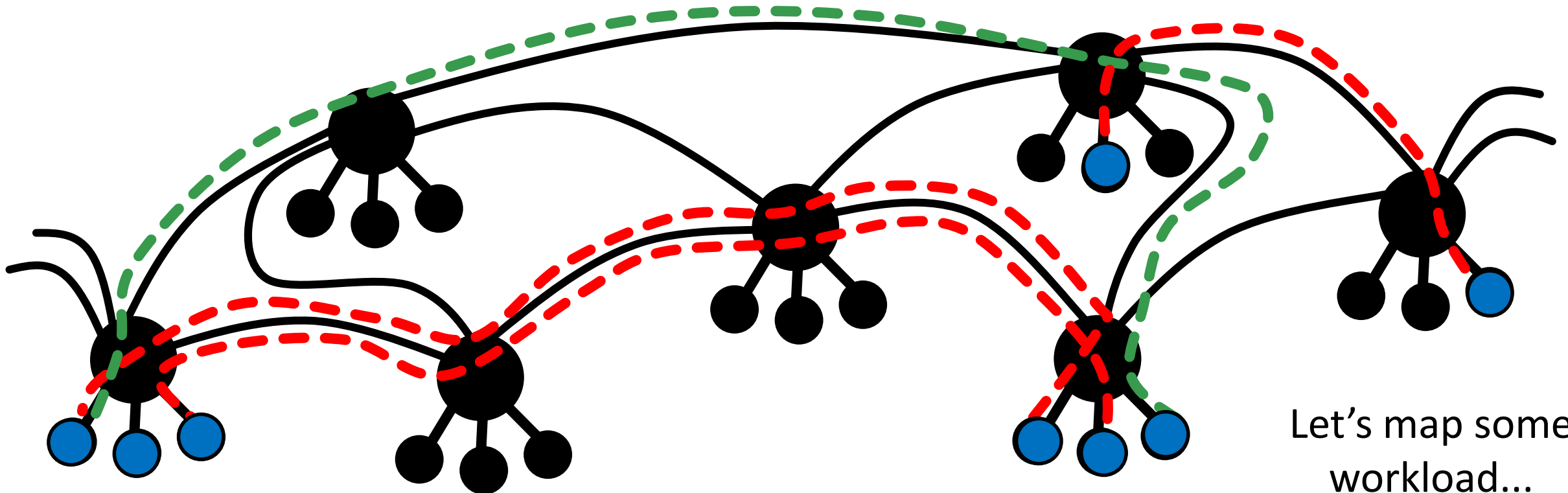
Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

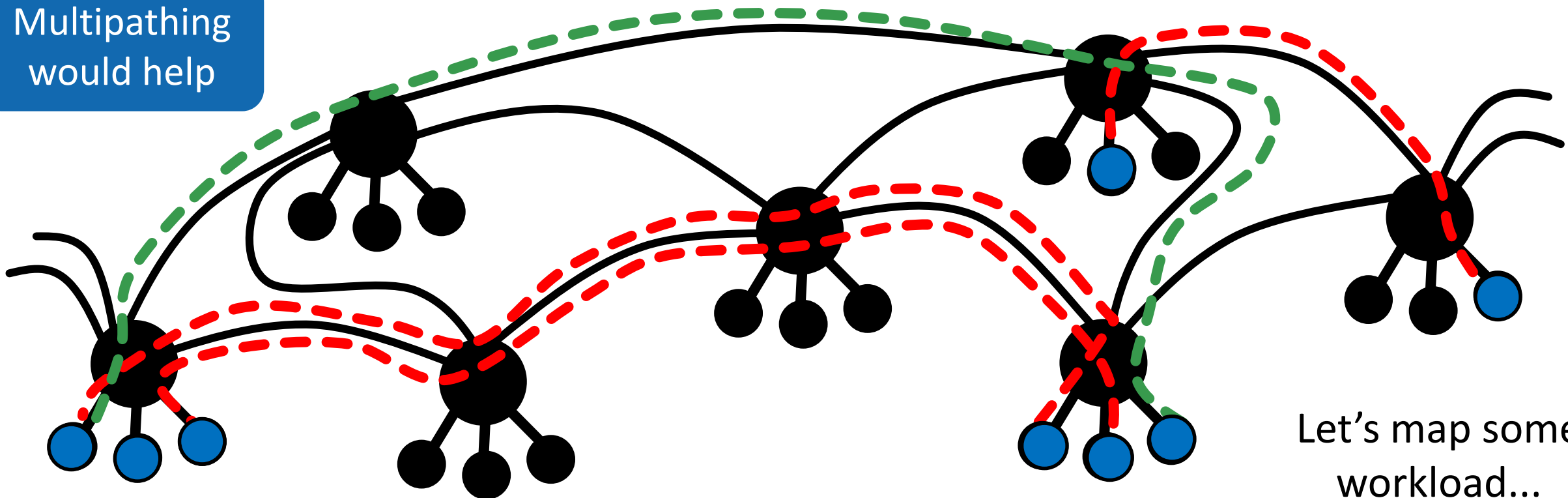


**✗ Flows collide!**

What are the problems that we want to tackle with multipathing?



Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

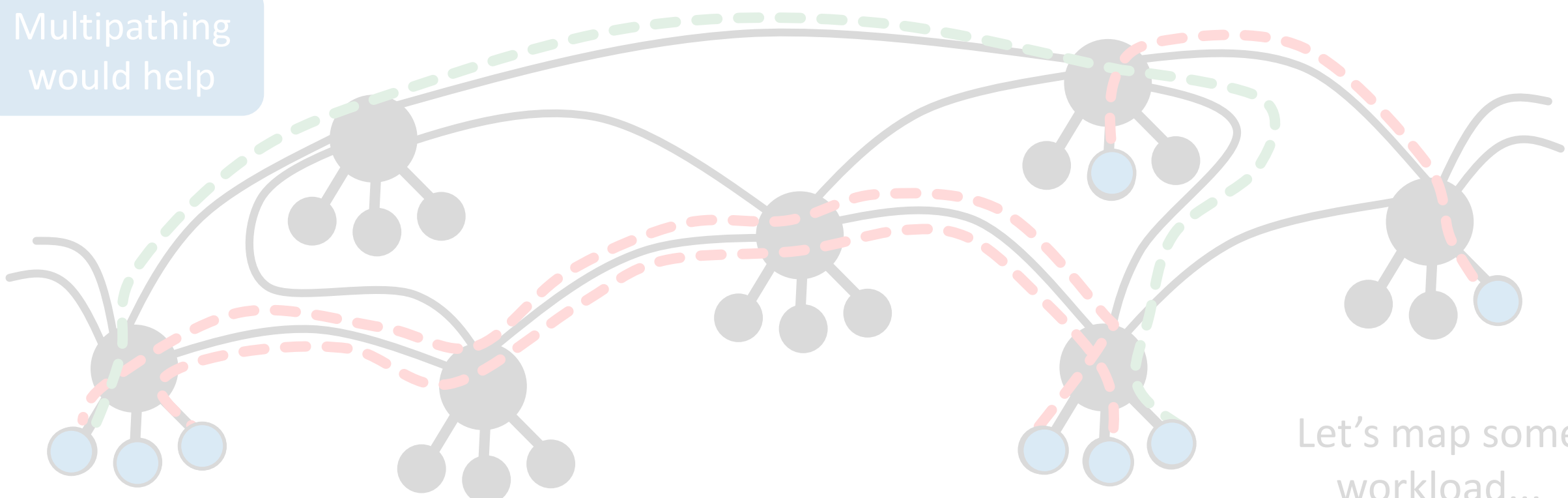


**✗** Flows collide!

What are the problems that we want to tackle with multipathing?



Multipathing would help



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

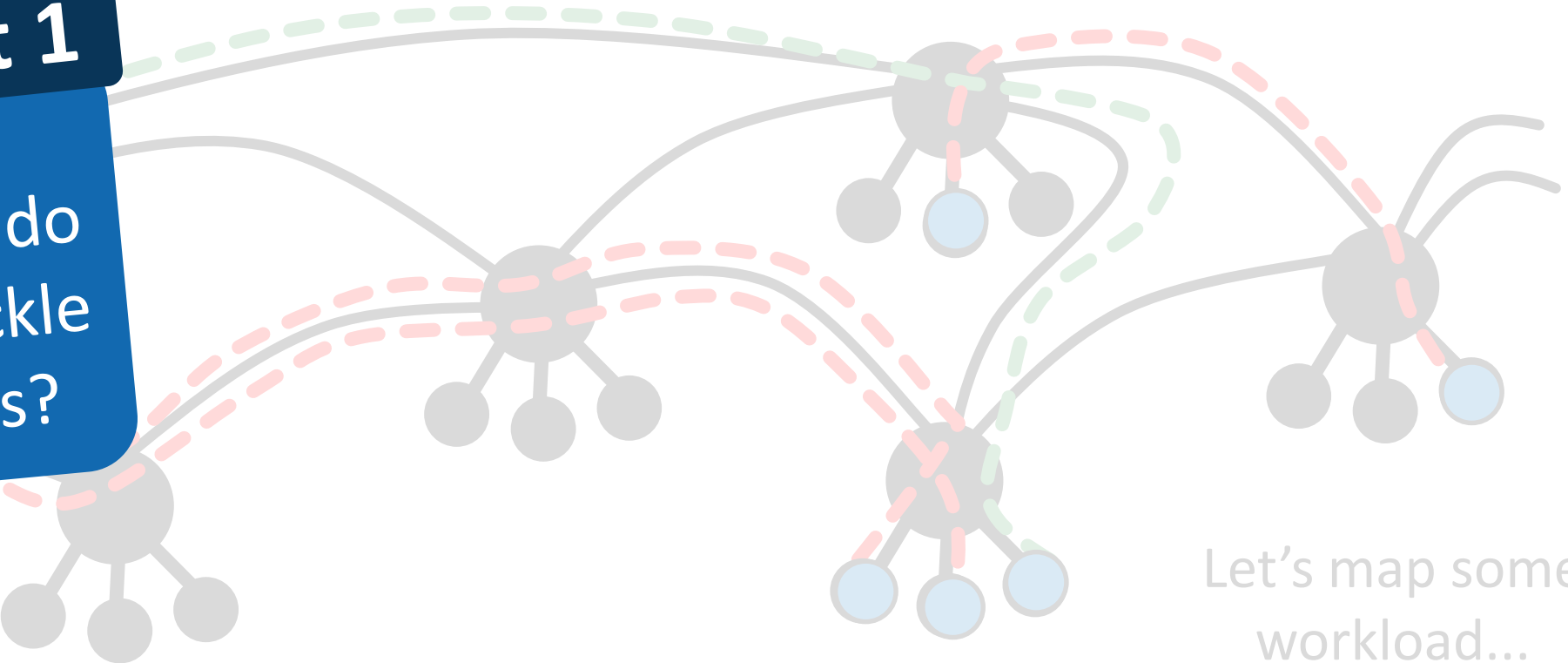


Multipathing

## Part 1



How many multiple paths do we need to tackle flow collisions?



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION

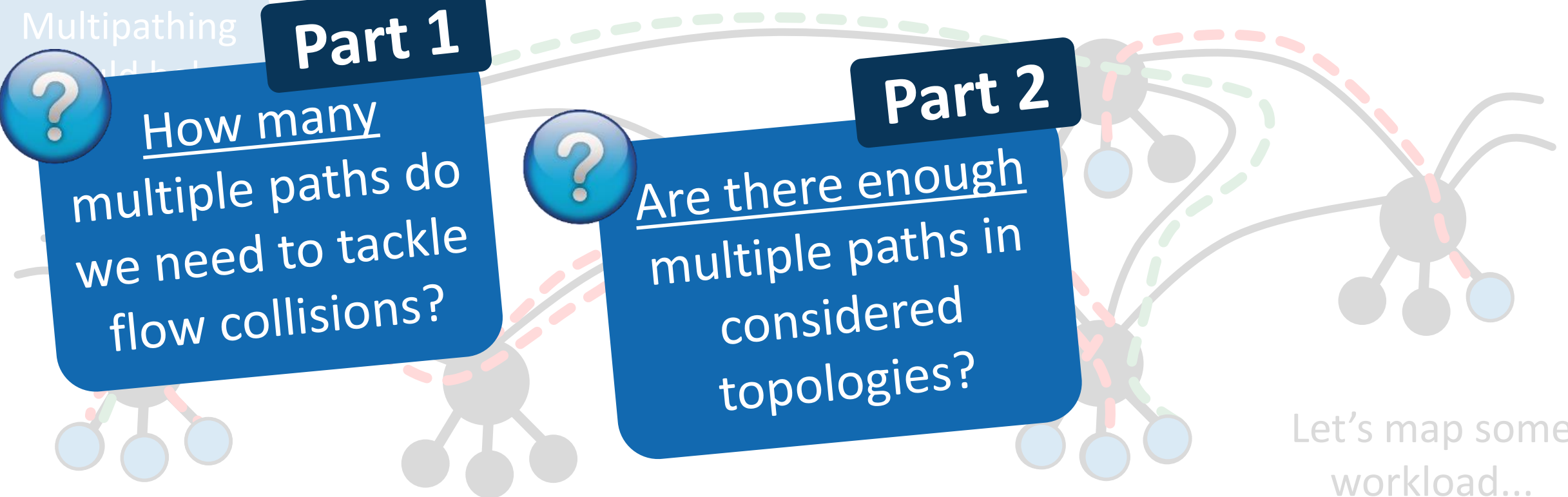


**✗** Flows collide!

What are the problems that we want to tackle with multipathing?

Multipathing  
**Part 1**  
**?** How many multiple paths do we need to tackle flow collisions?

**Part 2**  
**?** Are there enough multiple paths in considered topologies?



Let's map some workload...

# MULTIPATH ROUTING: MOTIVATION



**✗** Flows collide!

What are the problems that we want to tackle with multipathing?



Multipathing

## Part 1



How many multiple paths do we need to tackle flow collisions?

## Part 2



Are there enough multiple paths in considered topologies?

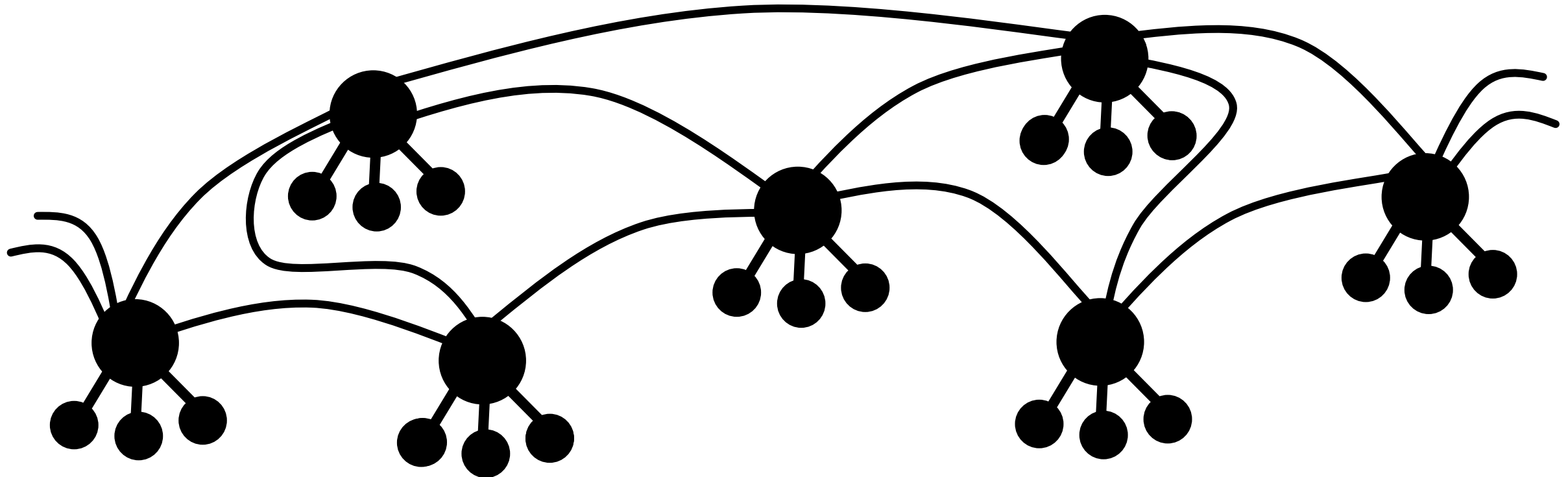
## Part 3



How to use such multiple paths in considered topologies?

# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

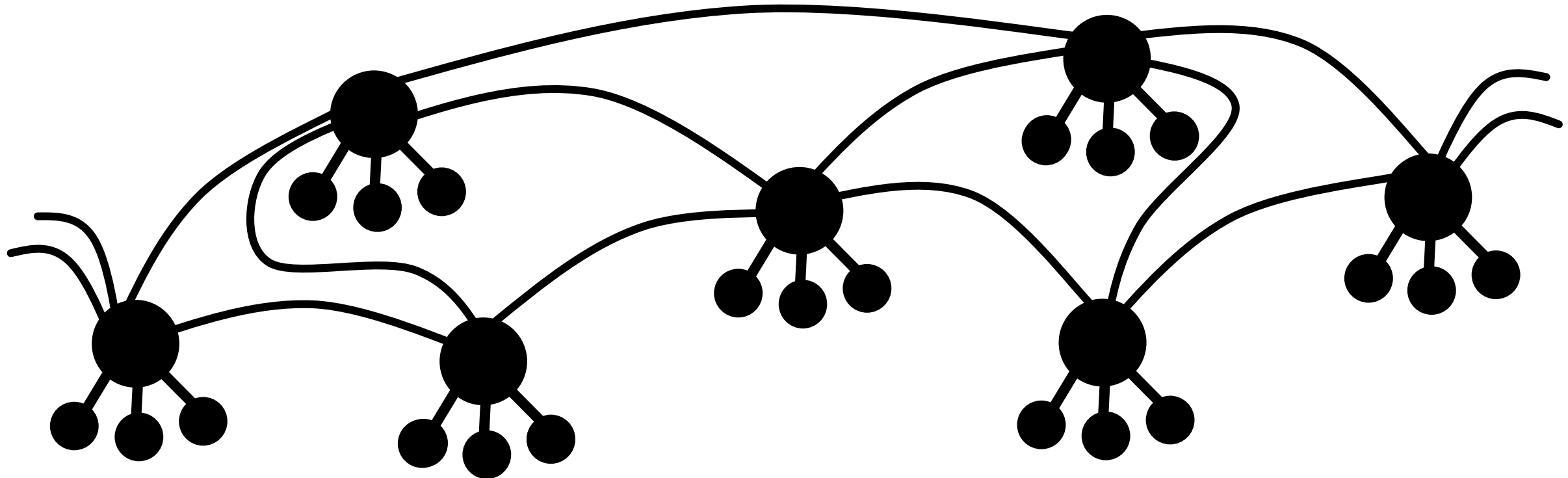
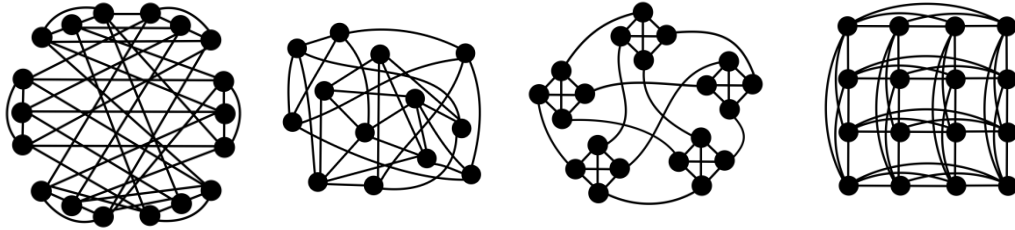




# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

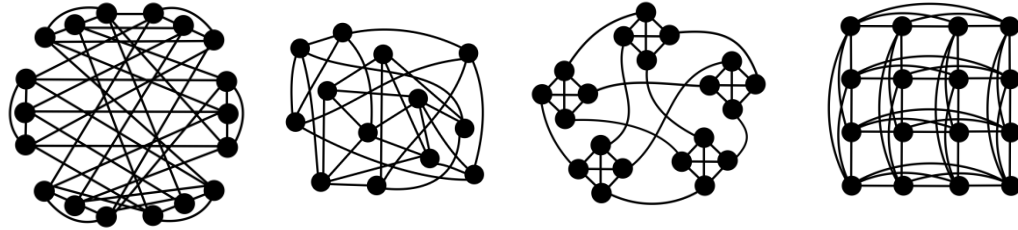
...For different topologies



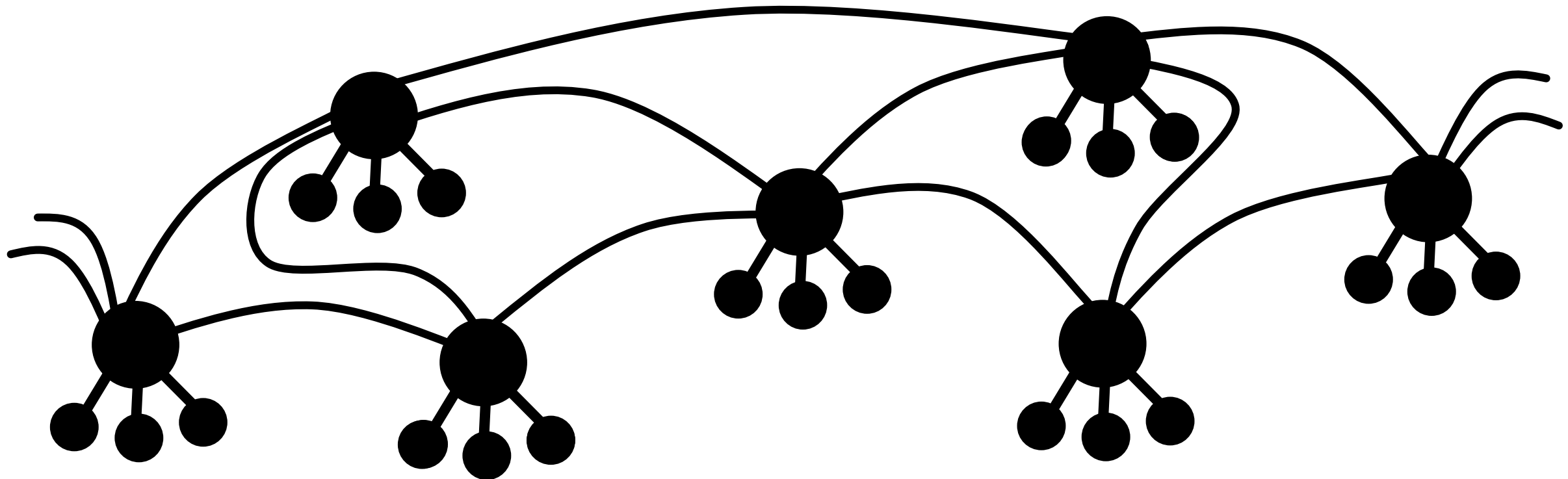
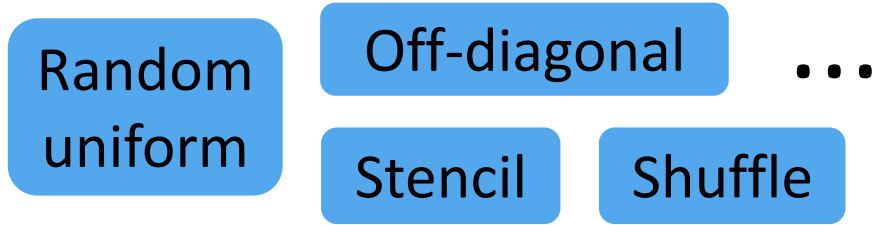
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



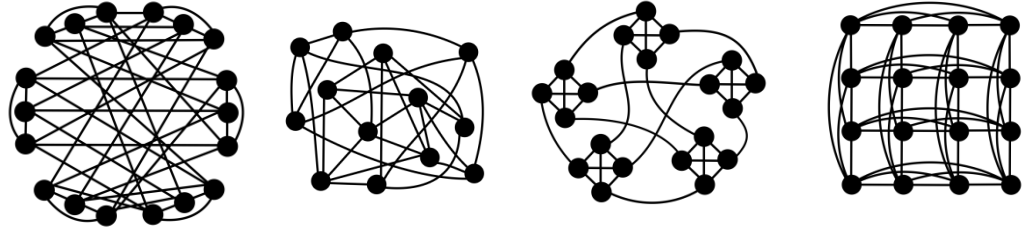
...For different workloads



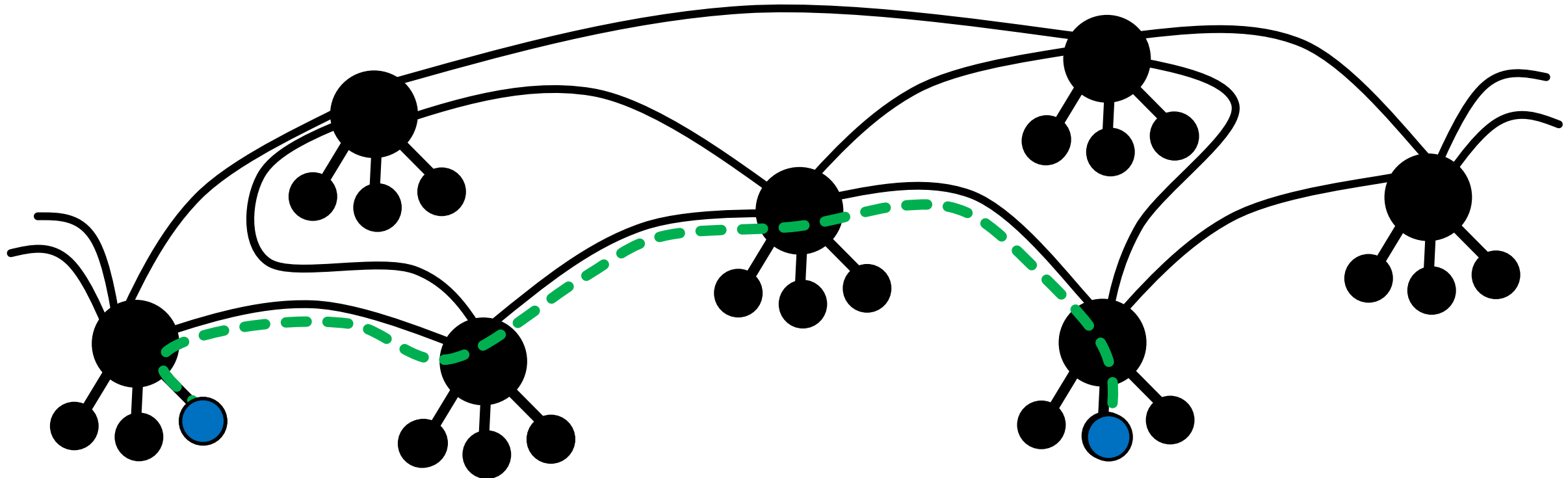
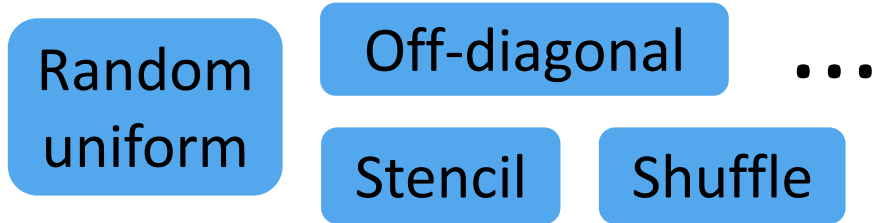
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



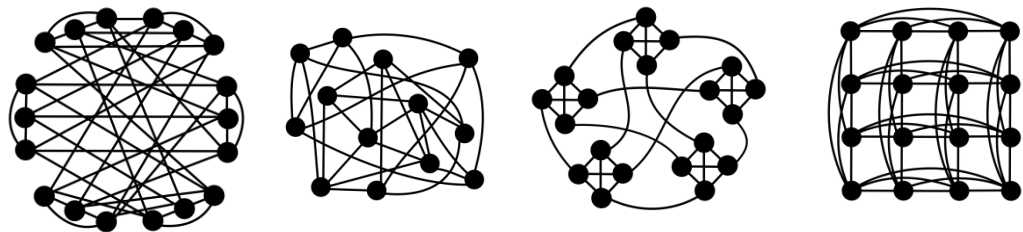
...For different workloads



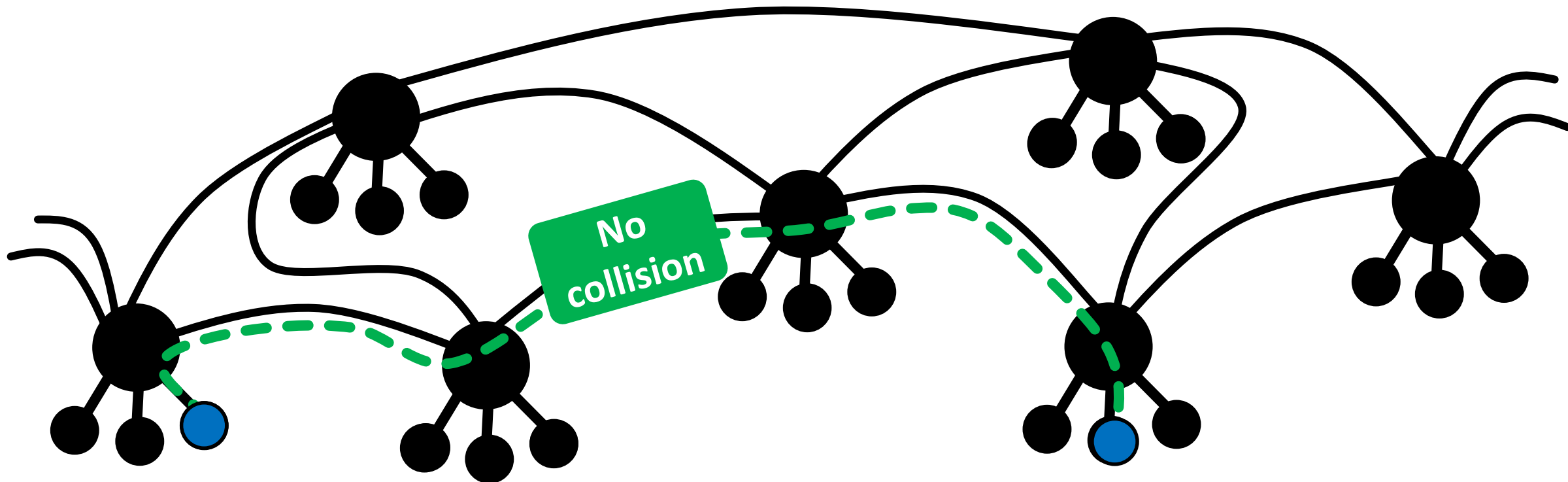
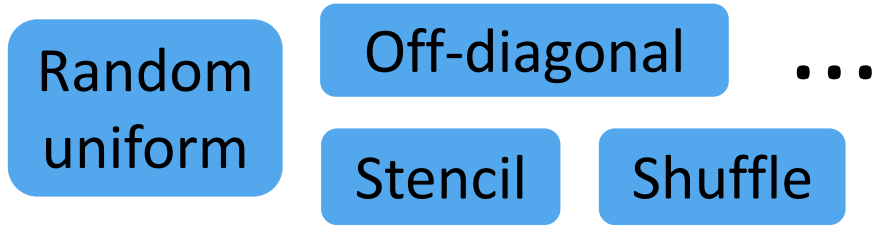
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



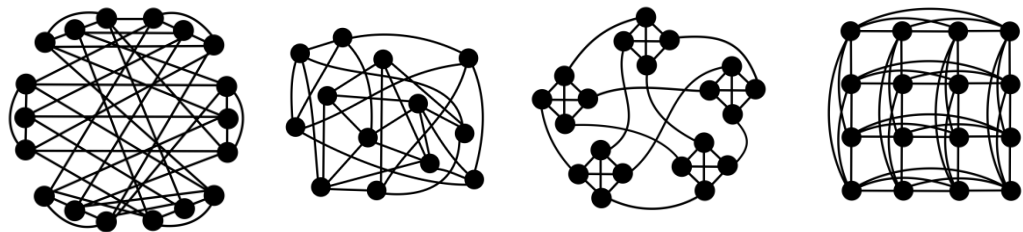
...For different workloads



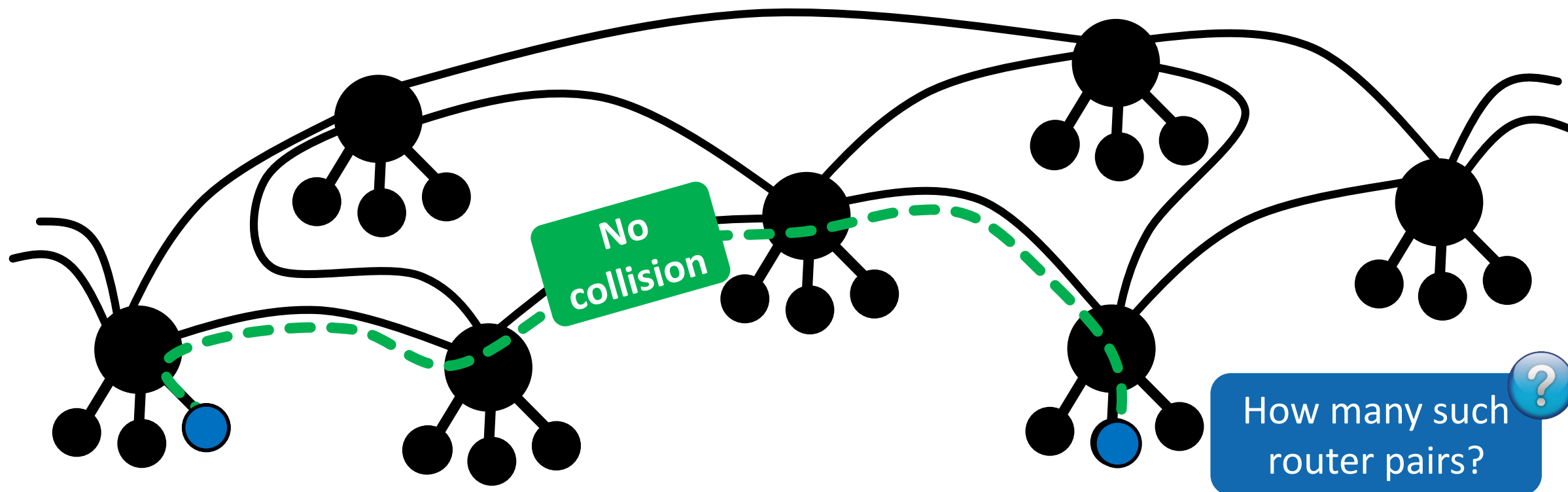
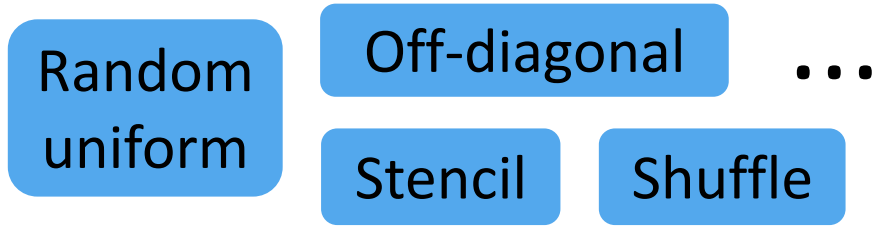
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



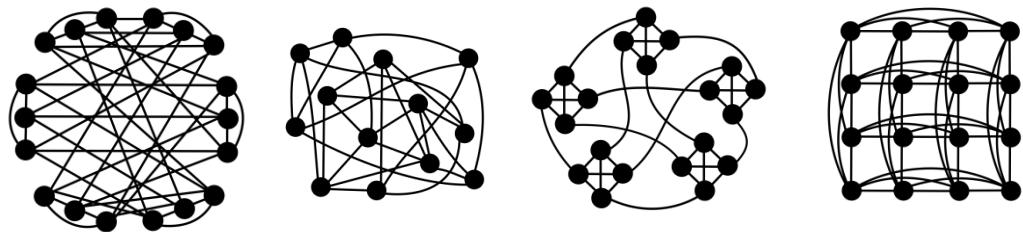
...For different workloads



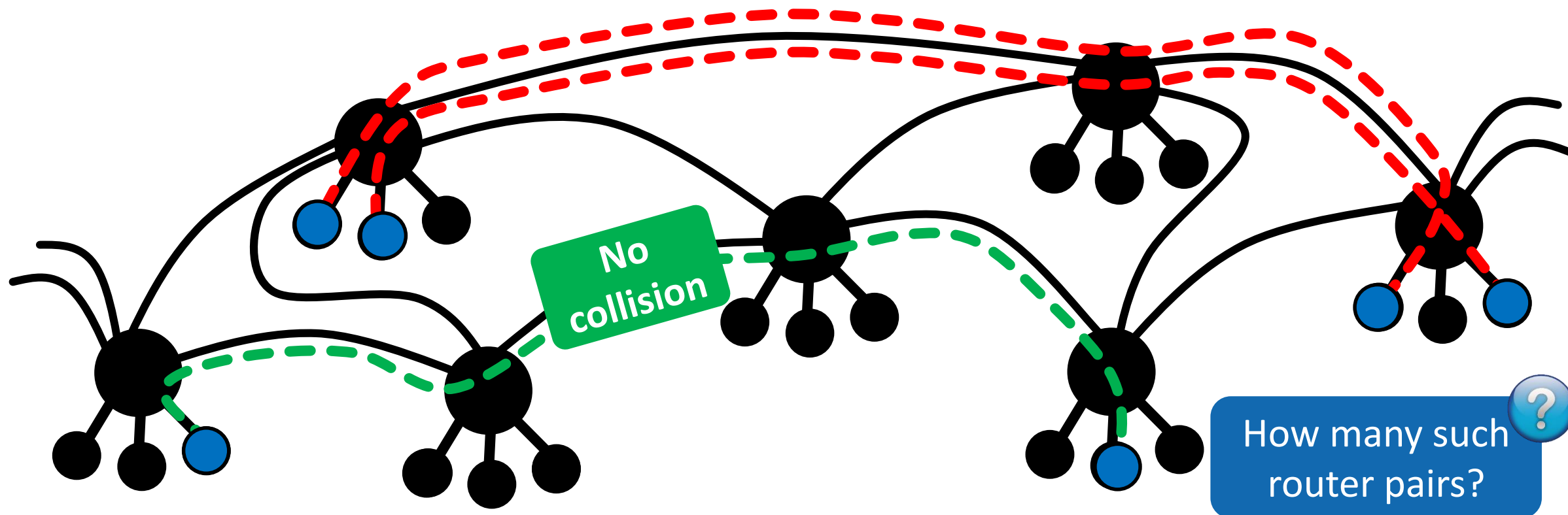
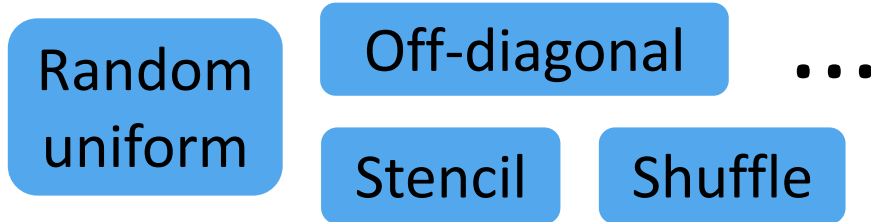
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



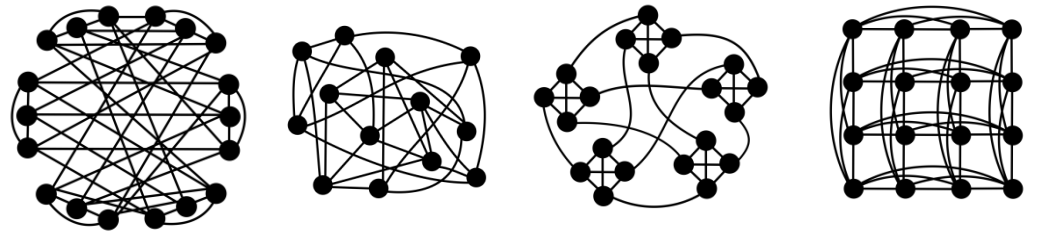
...For different workloads



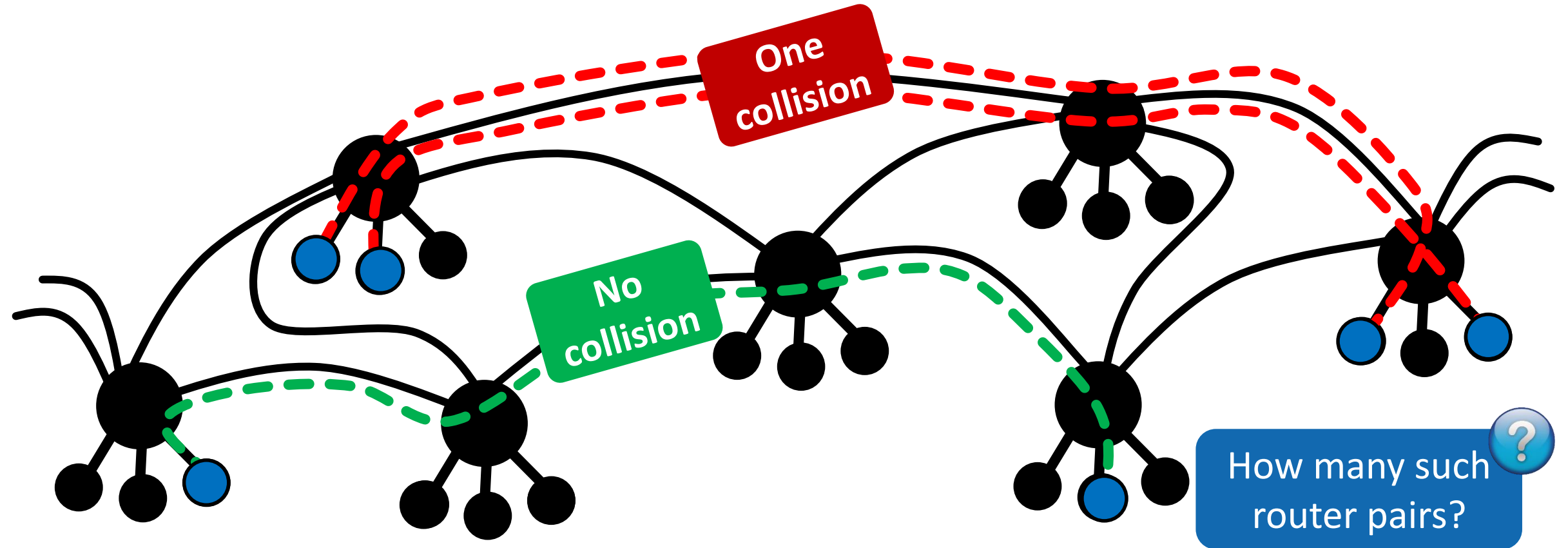
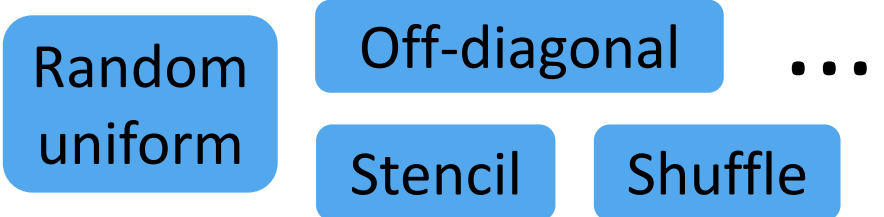
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



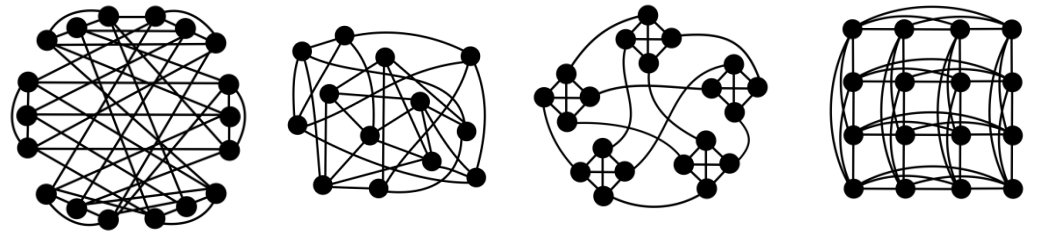
...For different workloads



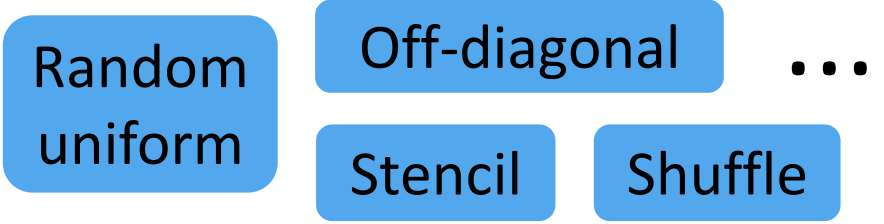
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

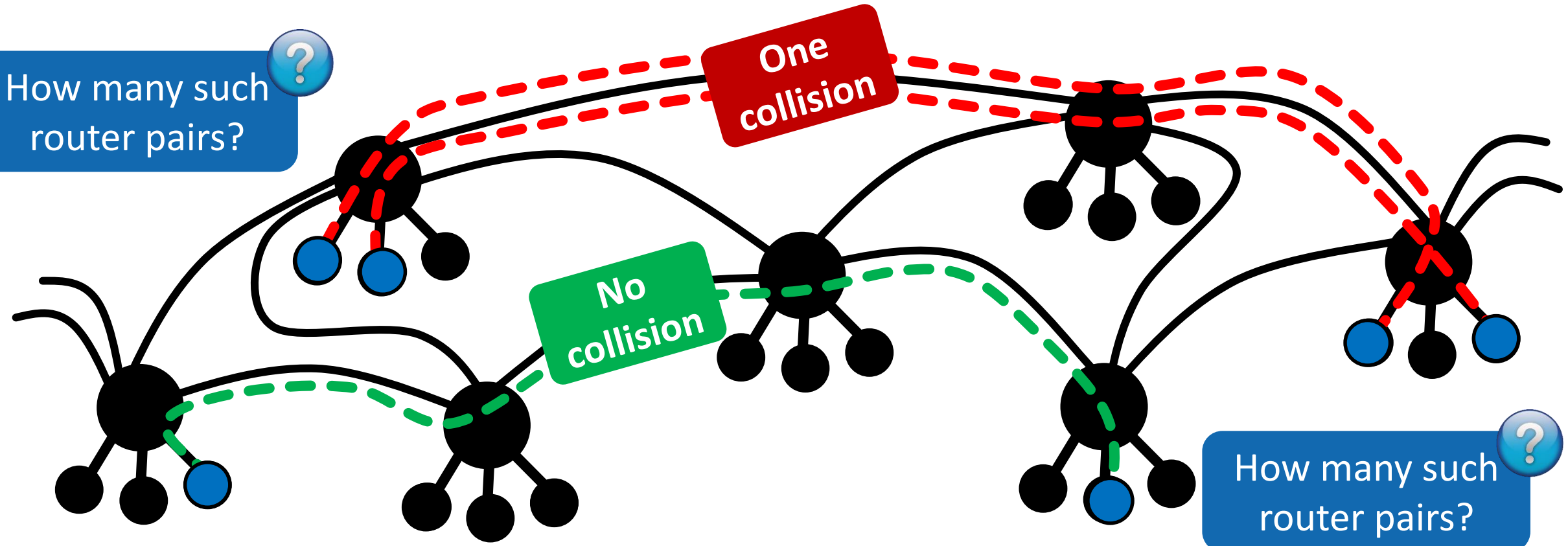
...For different topologies



...For different workloads



How many such router pairs?



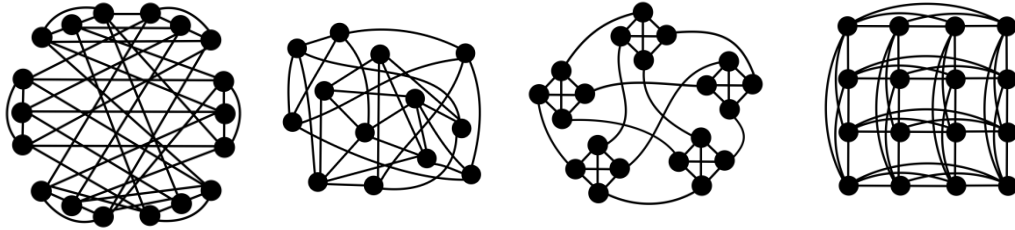
How many such router pairs?



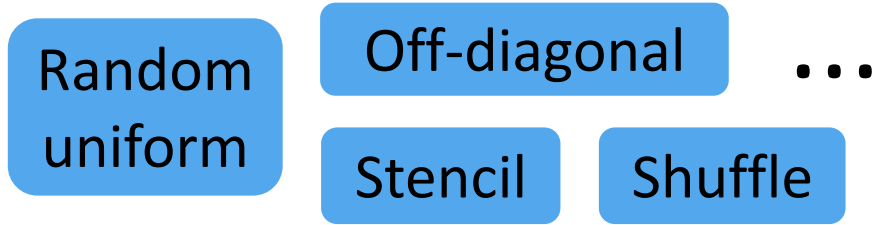
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



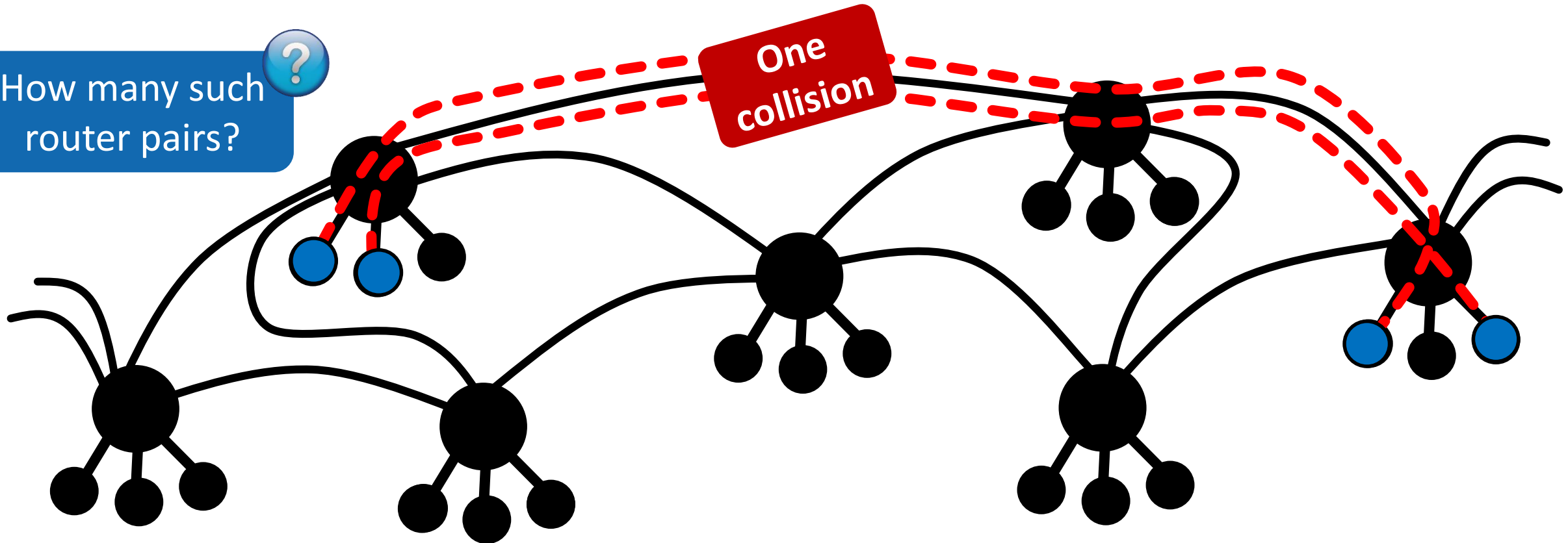
...For different workloads



How many such router pairs?



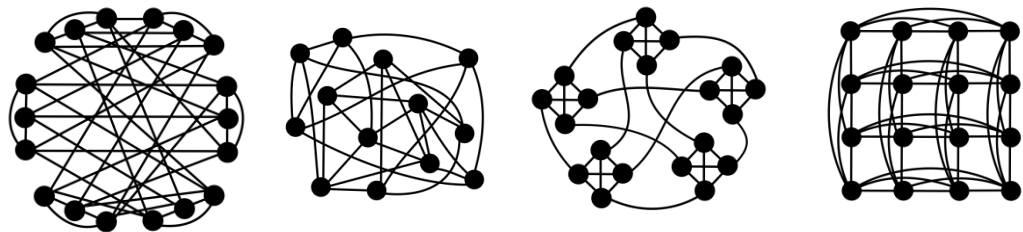
One collision



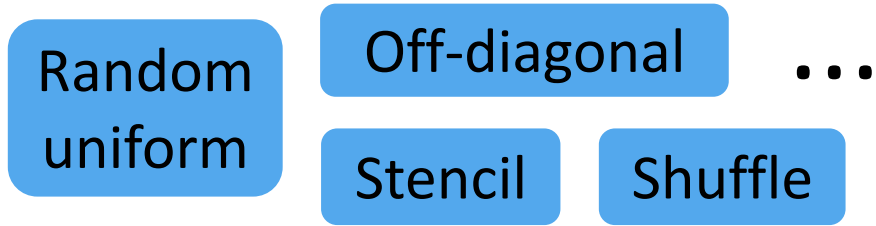
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

...For different topologies



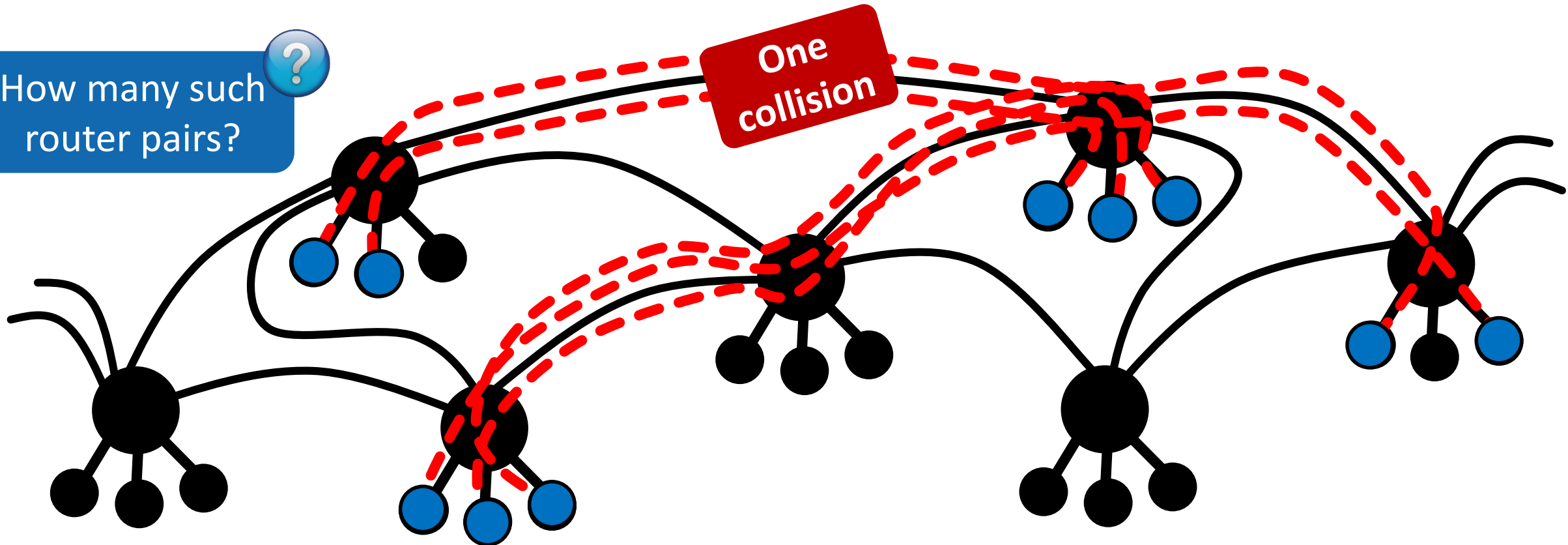
...For different workloads



How many such router pairs?



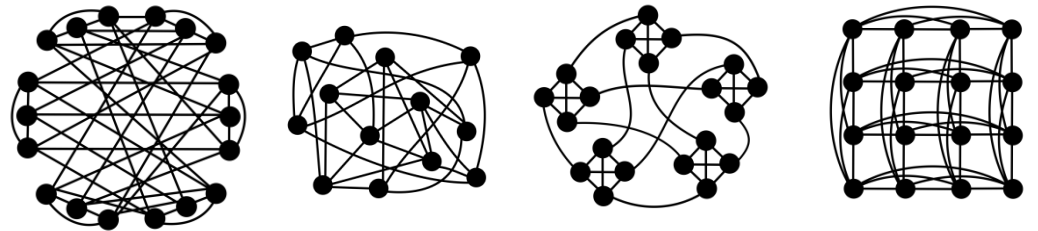
One collision



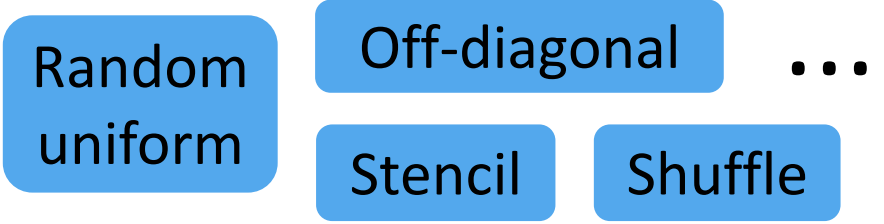
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

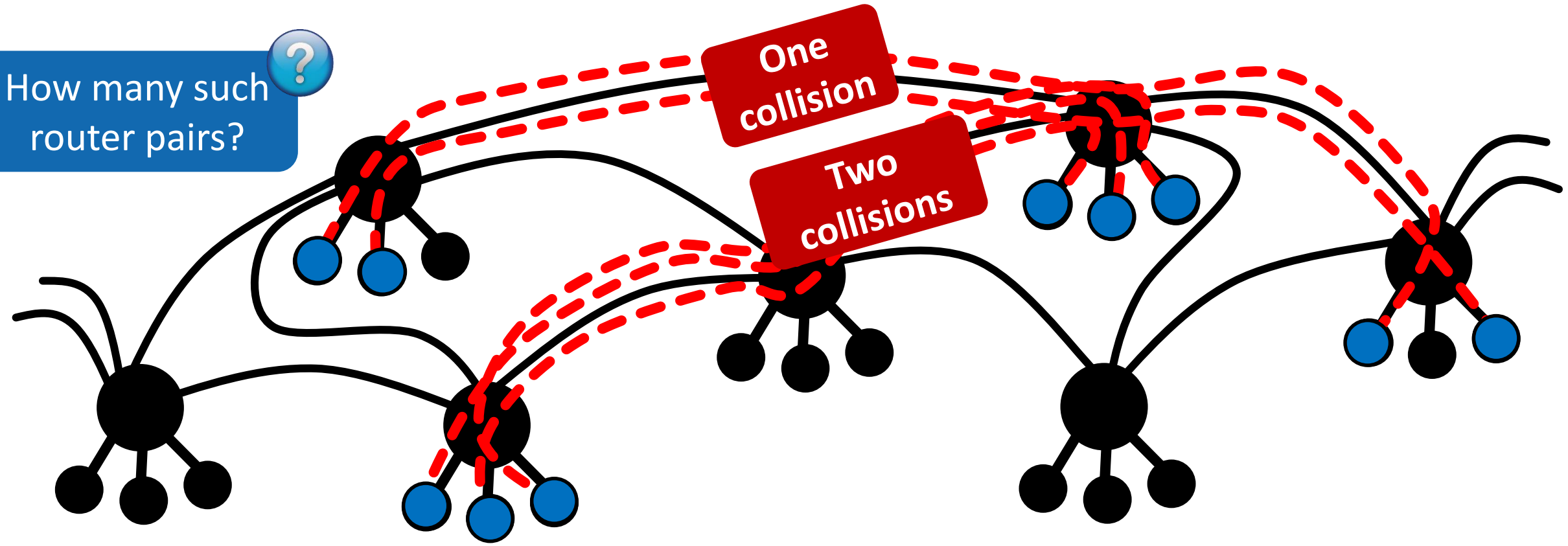
...For different topologies



...For different workloads



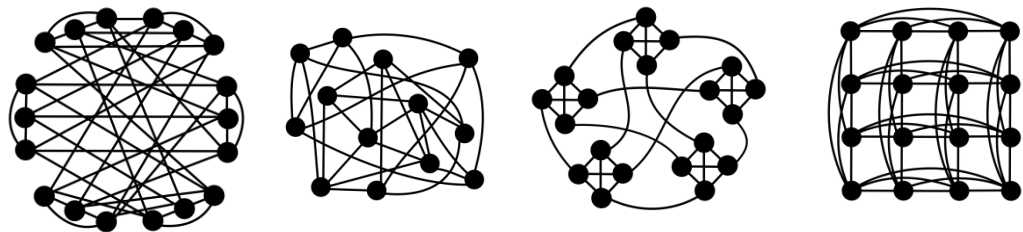
How many such router pairs?



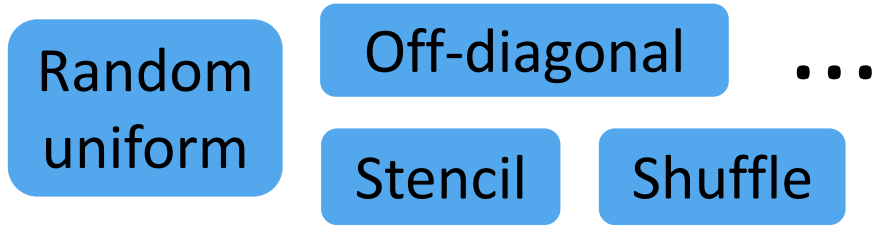
# HOW MANY MULTIPLE PATHS DO WE NEED TO TACKLE FLOW COLLISIONS?

## Part 1

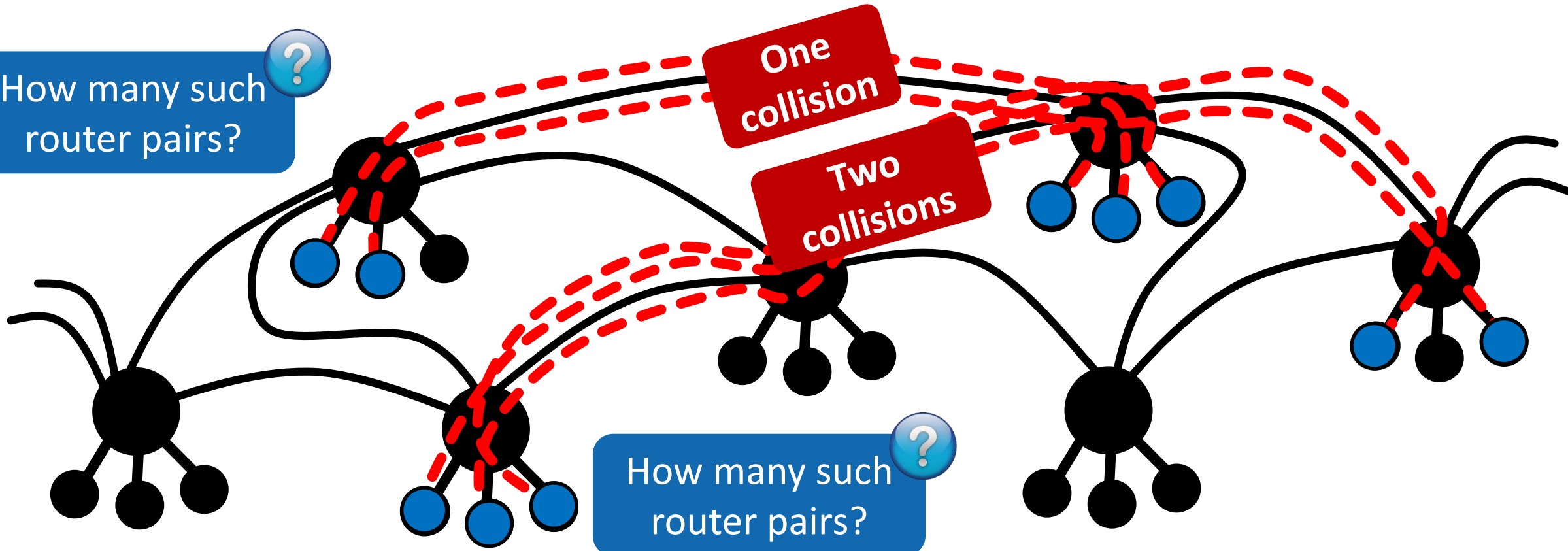
...For different topologies



...For different workloads



How many such router pairs?



How many such router pairs?

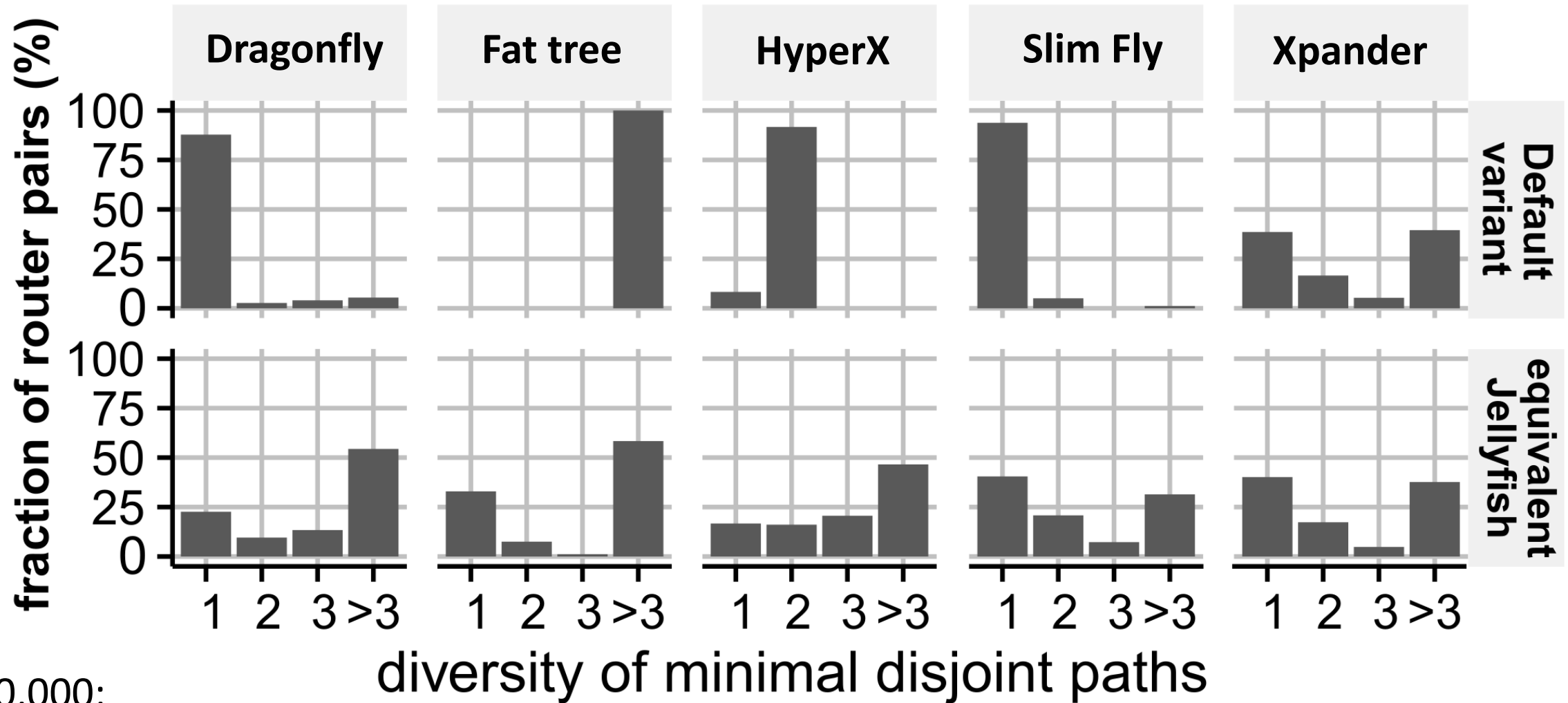
# PATH DIVERSITY ANALYSIS

Do we have enough such shortest paths?



# PATH DIVERSITY ANALYSIS

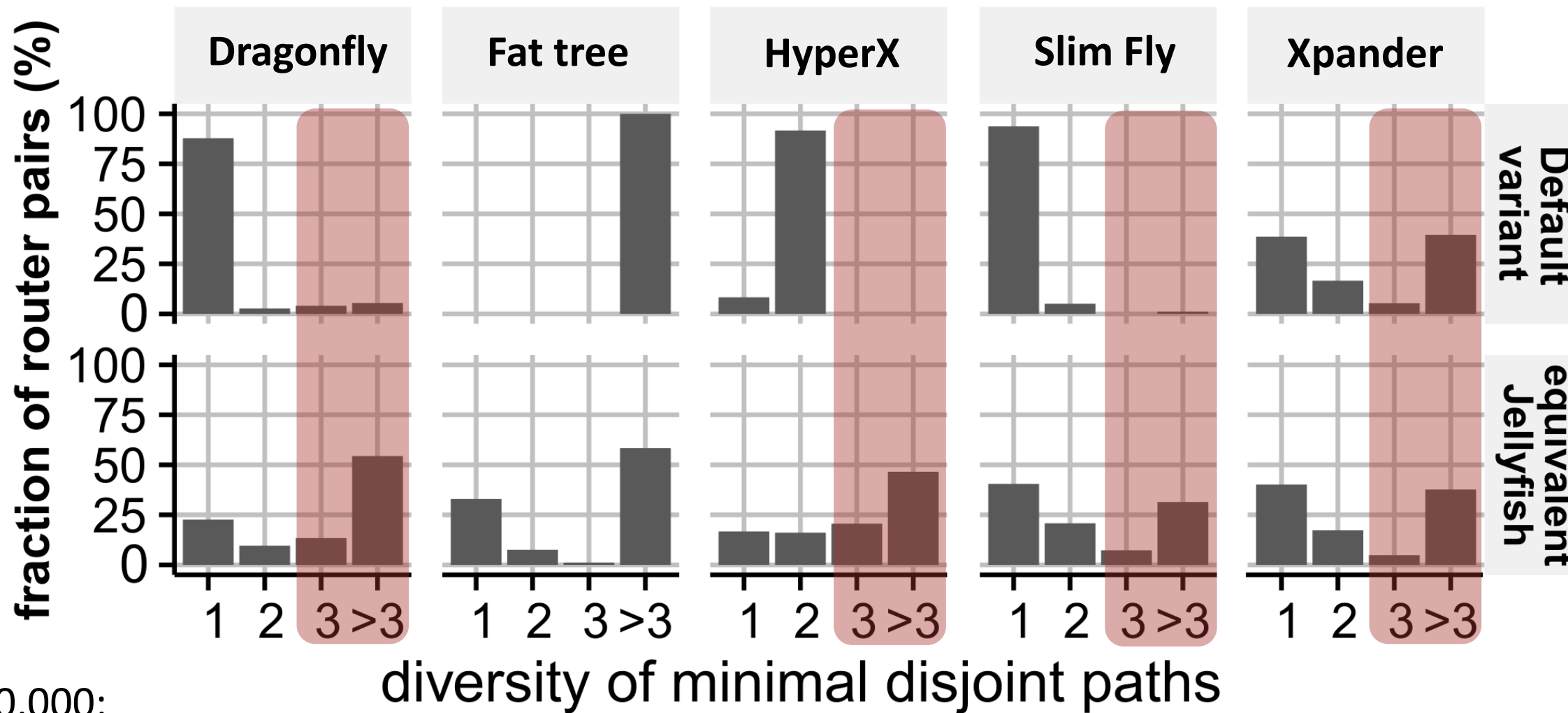
Do we have enough such shortest paths?



N ≈ 10,000;  
 comparable cost

# PATH DIVERSITY ANALYSIS

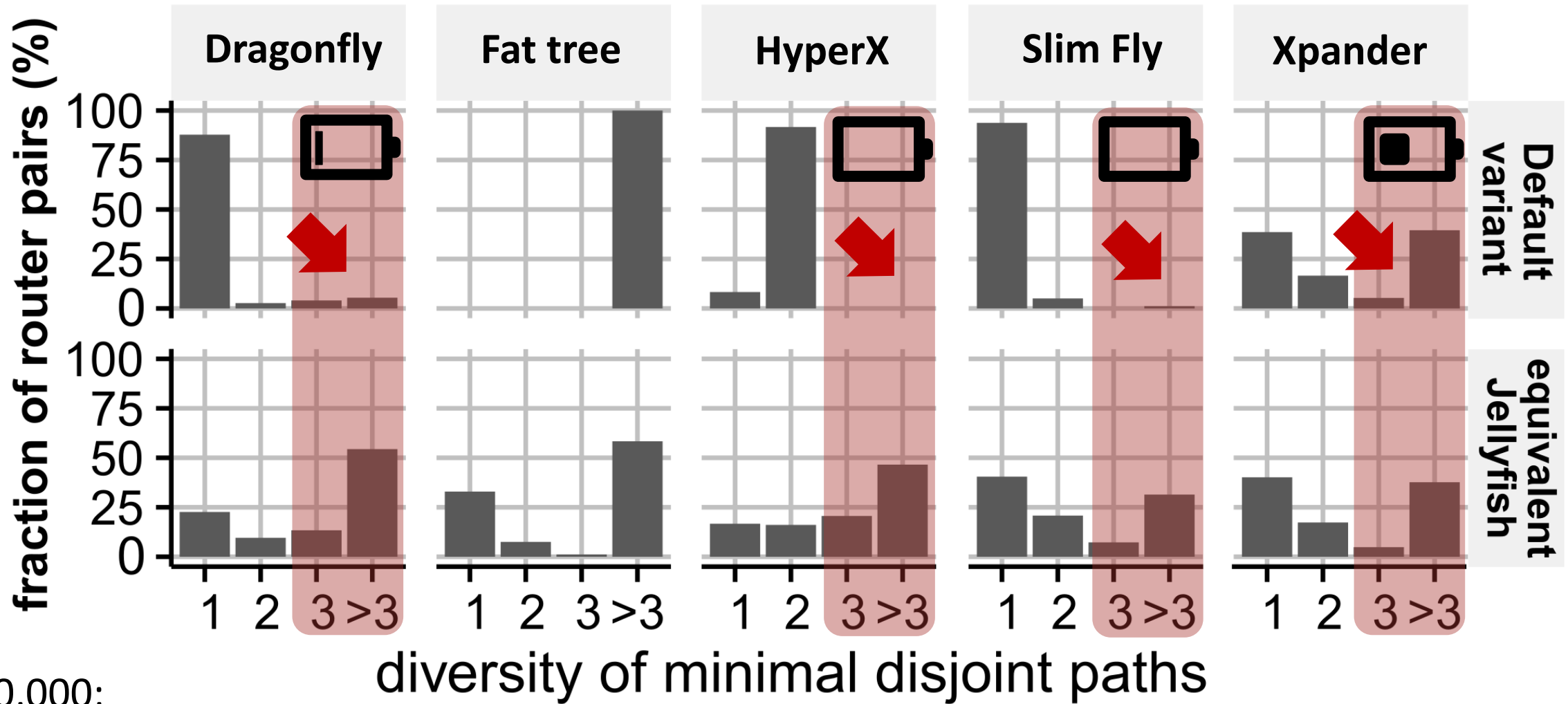
Do we have enough such shortest paths?



N ≈ 10,000;  
 comparable cost

# PATH DIVERSITY ANALYSIS

Do we have enough such shortest paths?

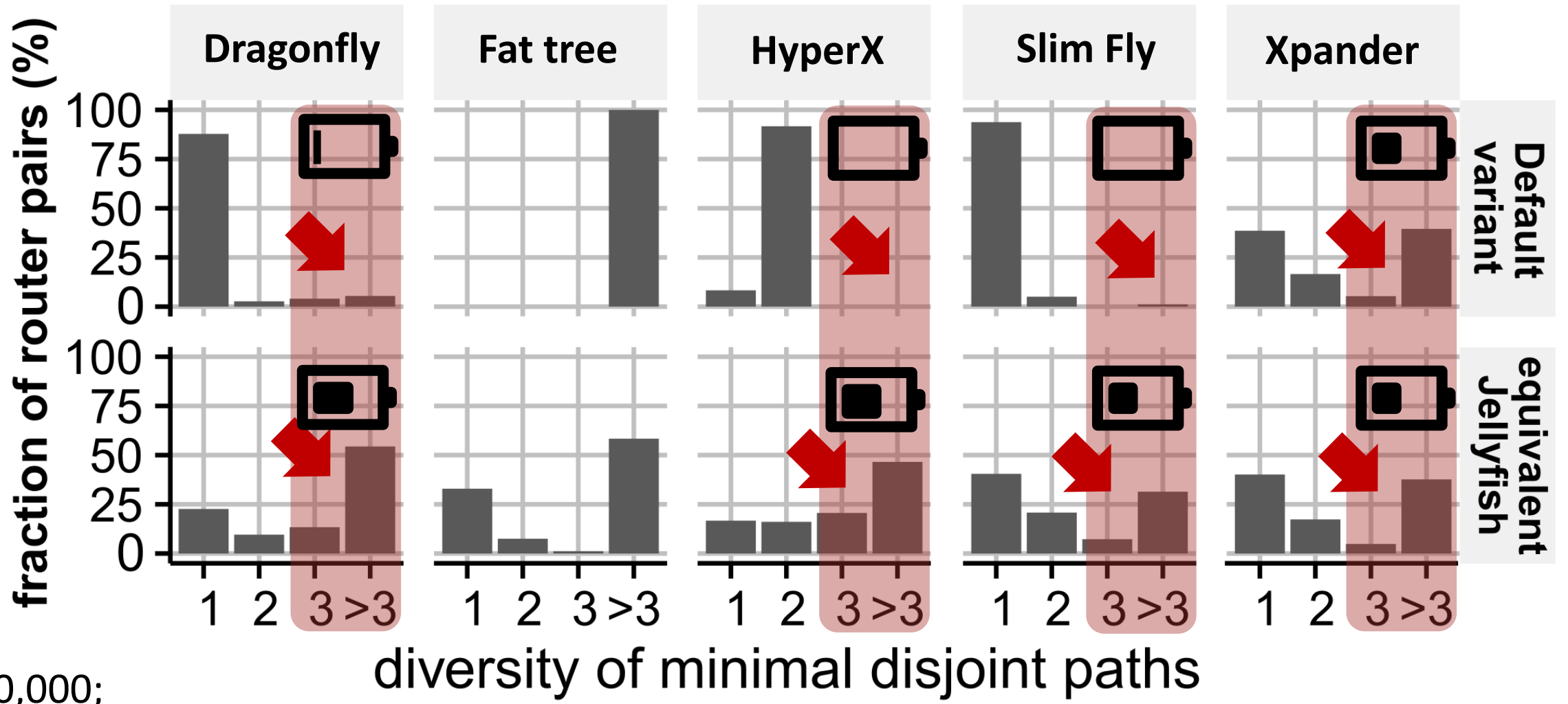


N ≈ 10,000;  
comparable cost



# PATH DIVERSITY ANALYSIS

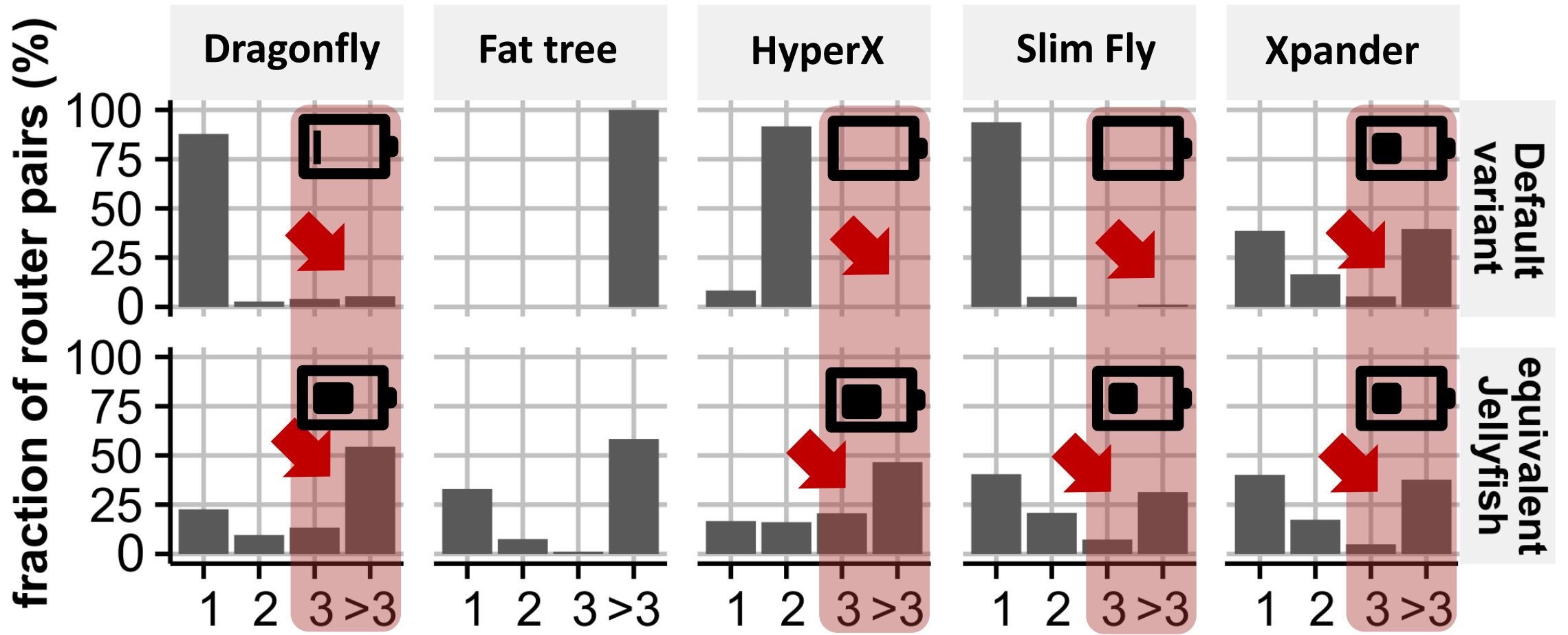
Do we have enough such shortest paths?



N ≈ 10,000;  
 comparable cost

# PATH DIVERSITY ANALYSIS

Do we have enough such shortest paths?

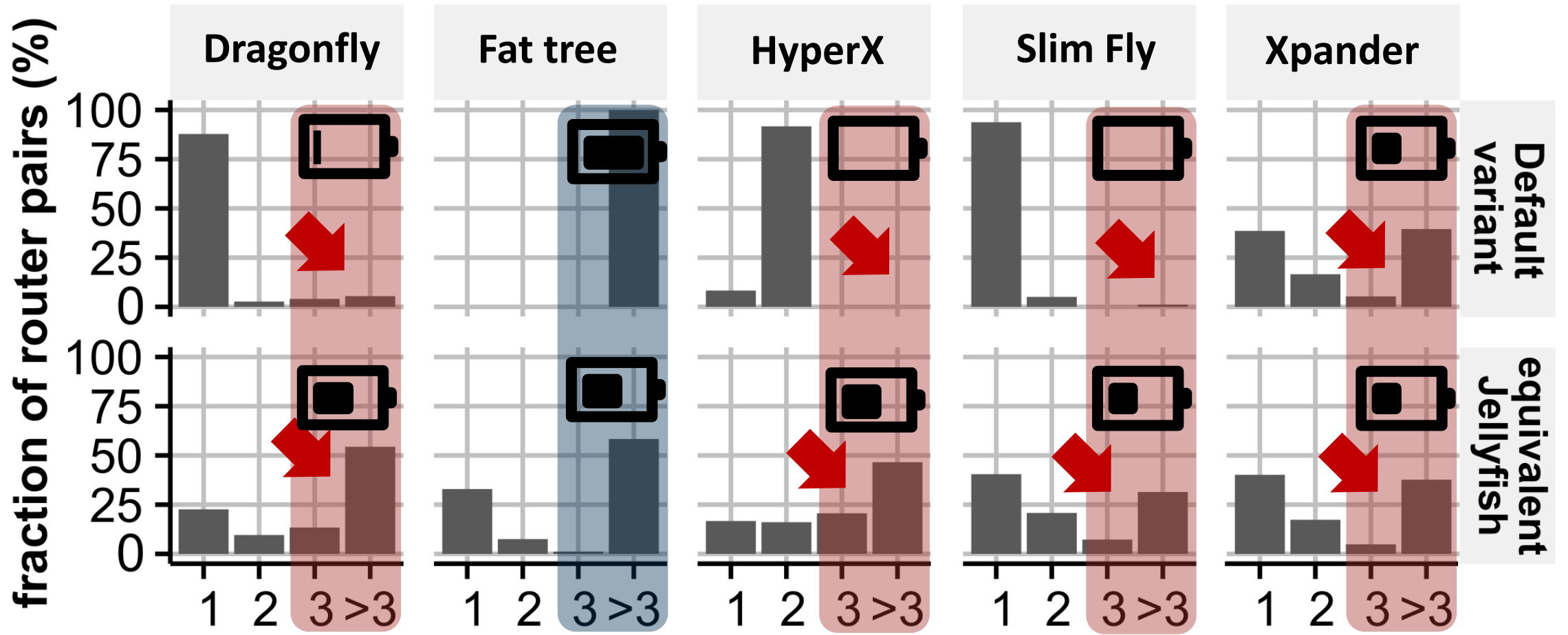


N ≈ 10,000;  
 comparable cost

diversity of minimal paths **Randomness in Jellyfish "smooths out" distributions of minimal path diversities**

# PATH DIVERSITY ANALYSIS

Do we have enough such shortest paths?



N ≈ 10,000;  
 comparable cost

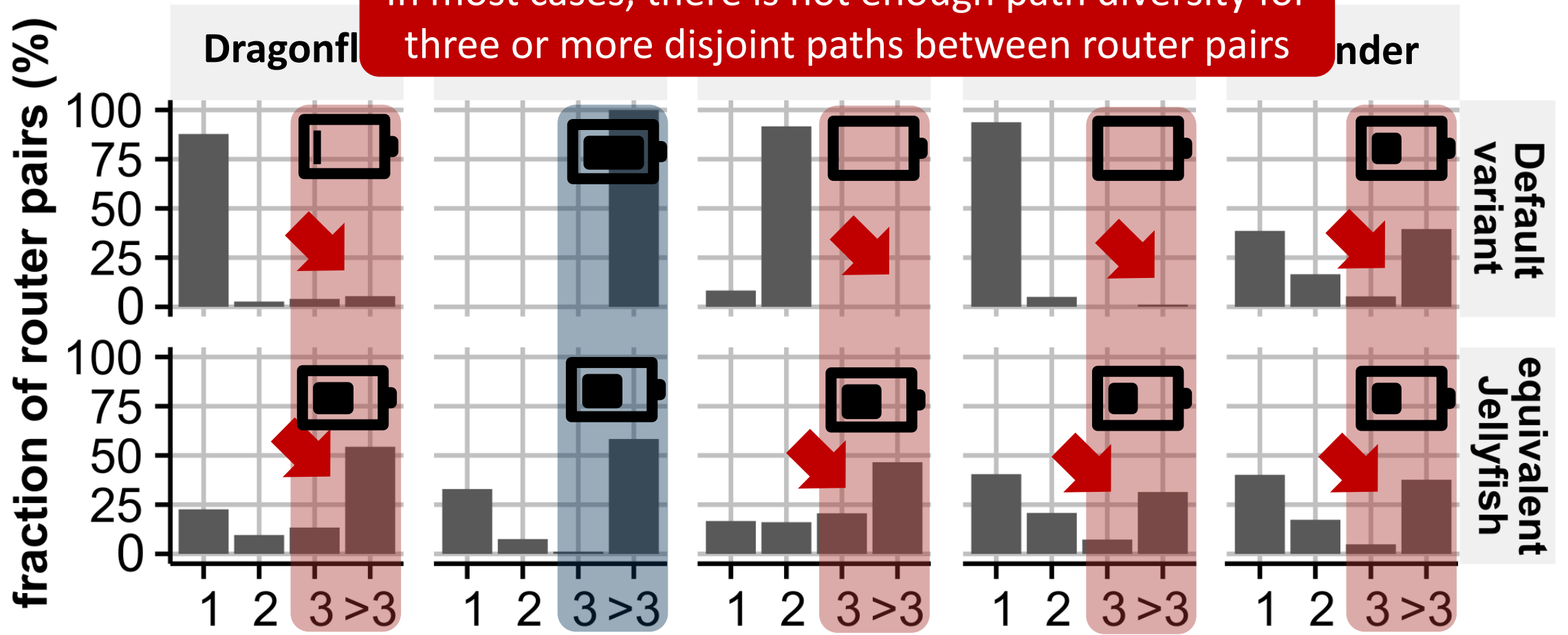
diversity of minimal paths

Randomness in Jellyfish "smooths out" distributions of minimal path diversities

# PATH DIVERSITY ANALYSIS

Do we have enough such shortest paths? ?

In most cases, there is not enough path diversity for three or more disjoint paths between router pairs



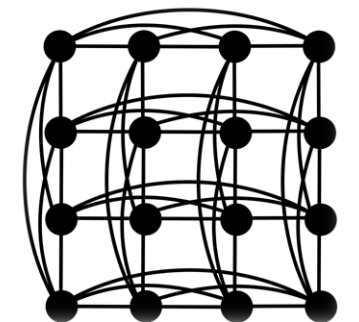
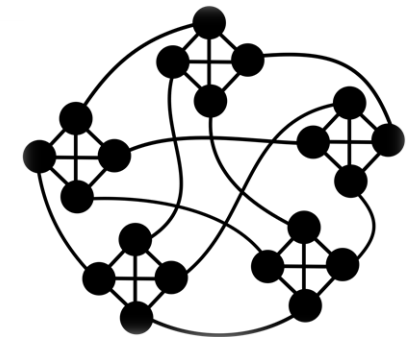
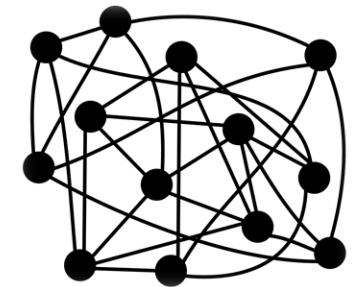
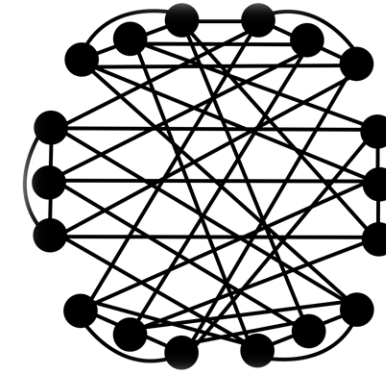
N ≈ 10,000;  
comparable cost

Randomness in Jellyfish "smooths out" distributions of minimal path diversities

# LOW-DIAMETER NETWORK TOPOLOGIES



Can we use multipathing?  
i.e., are there multiple paths between routers?



... shortest

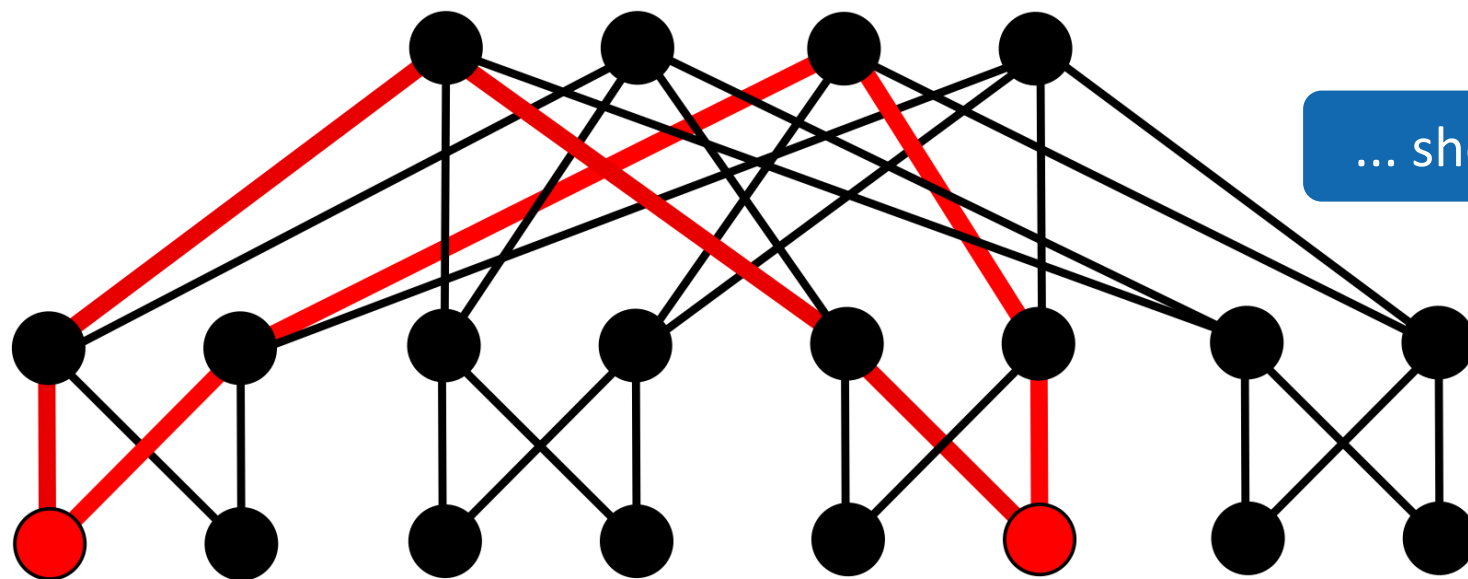
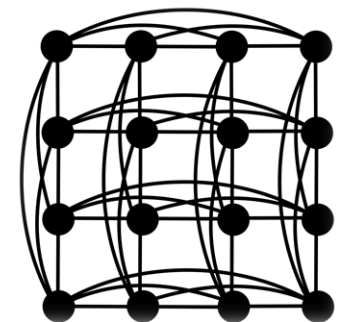
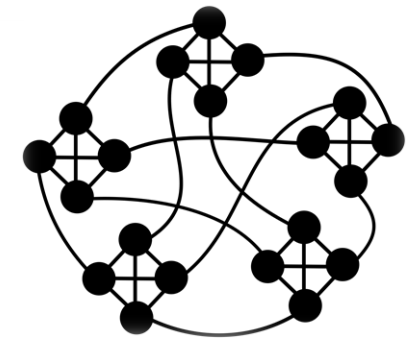
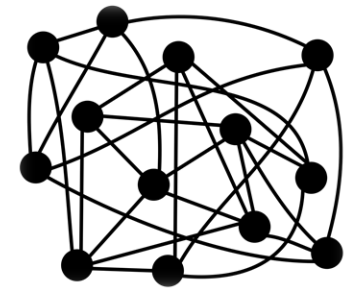
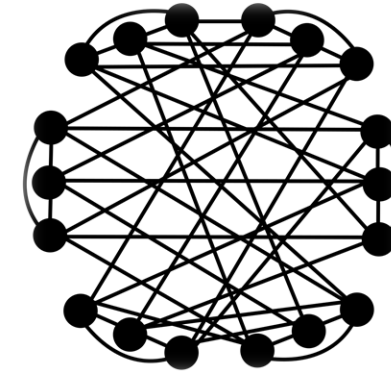
... equal length

... Between all router pairs

# LOW-DIAMETER NETWORK TOPOLOGIES



Can we use multipathing?  
i.e., are the multiple paths between routers?



... shortest

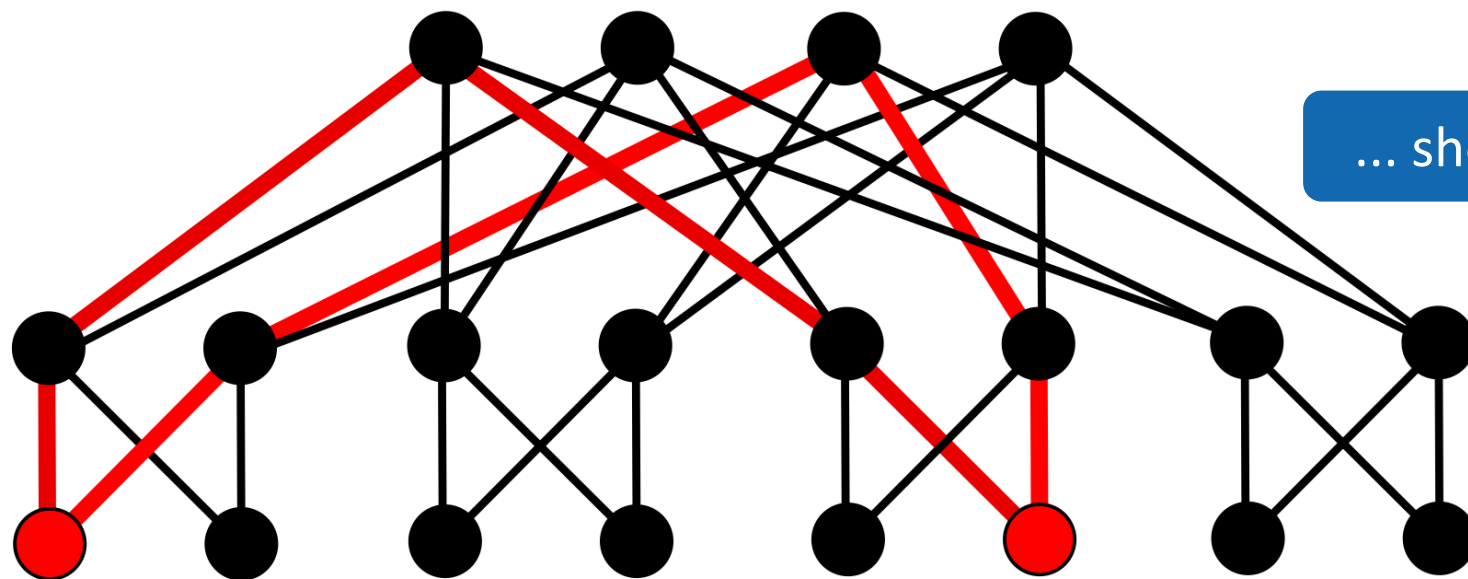
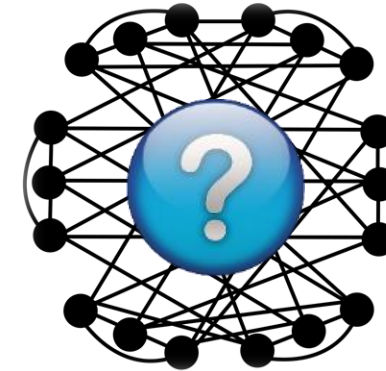
... equal length

... Between all router pairs

# LOW-DIAMETER NETWORK TOPOLOGIES



Can we use multipathing?  
i.e., are there multiple paths between routers?



... shortest

... equal length

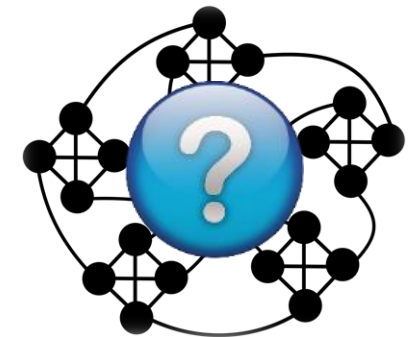
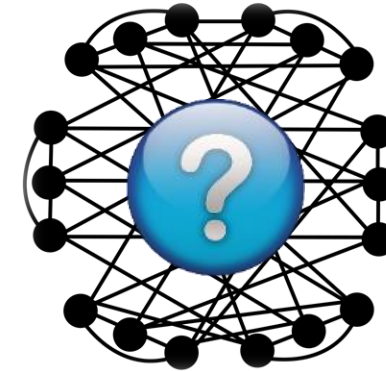
... Between all router pairs

# LOW-DIAMETER NETWORK TOPOLOGIES

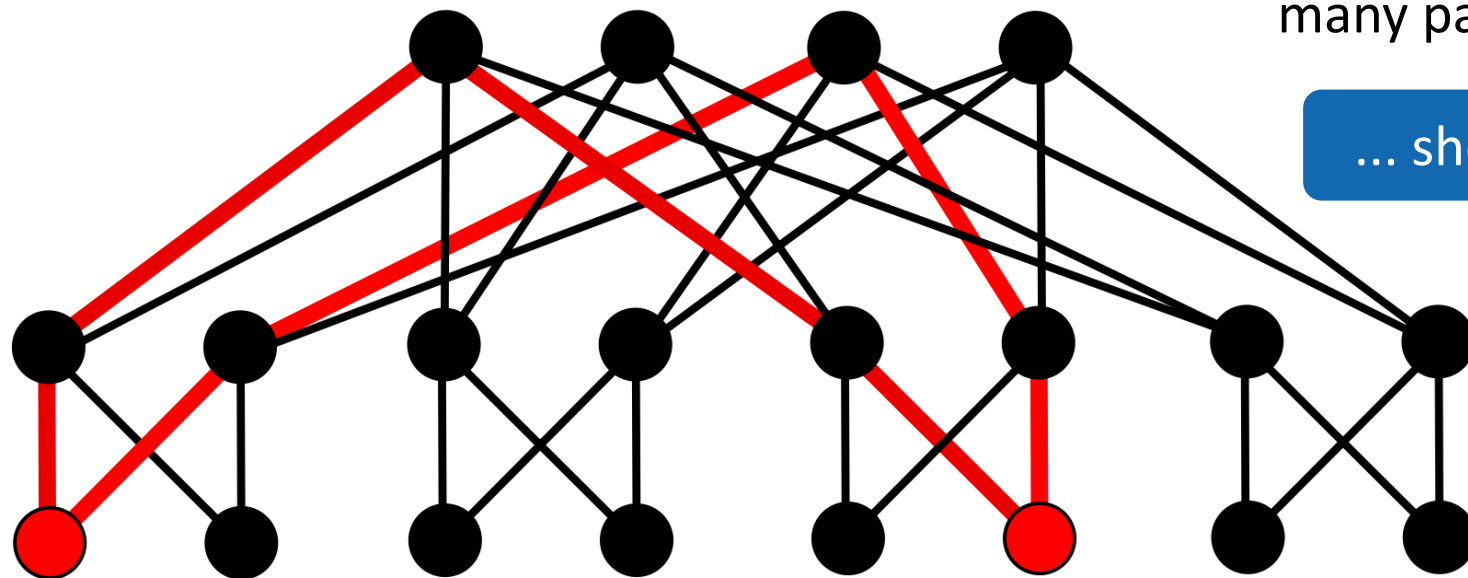


Can we use multipathing?  
i.e., are the multiple paths between routers?

Part 2



In Fat trees, easily, as we have many paths...



... shortest

... equal length

... Between all router pairs



# PATH DIVERSITY ANALYSIS

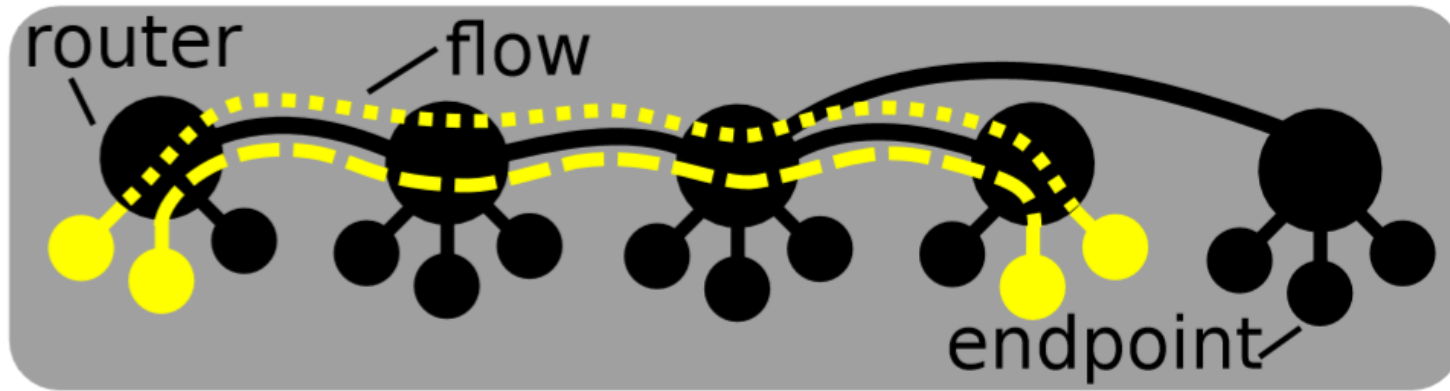
## PATH DIVERSITY ANALYSIS

What are the problems to be tackled with multipathing?



# PATH DIVERSITY ANALYSIS

What are the problems to be tackled with multipathing?

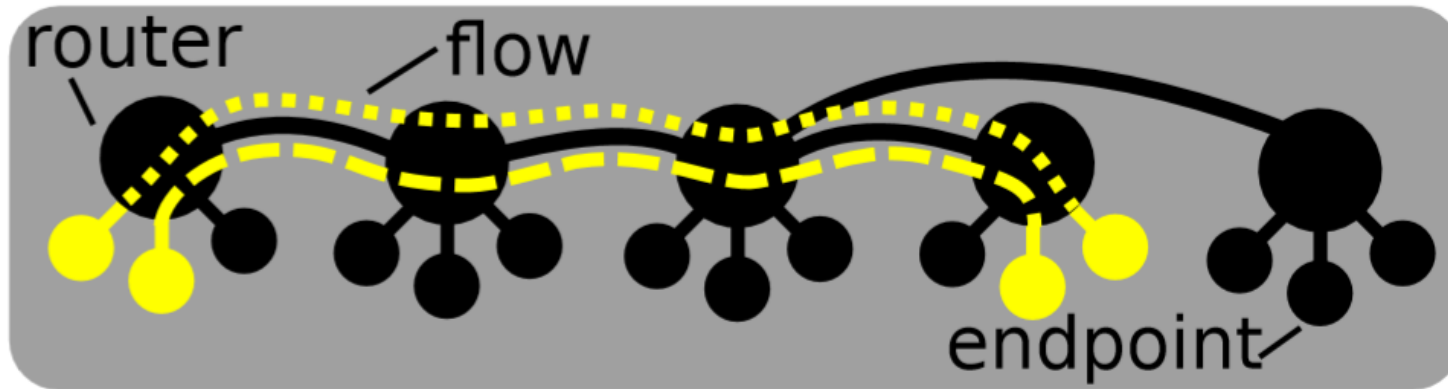


**Flow collisions**

Communicating endpoint pairs mapped to the same routers

## PATH DIVERSITY ANALYSIS

What are the problems to be tackled with multipathing?



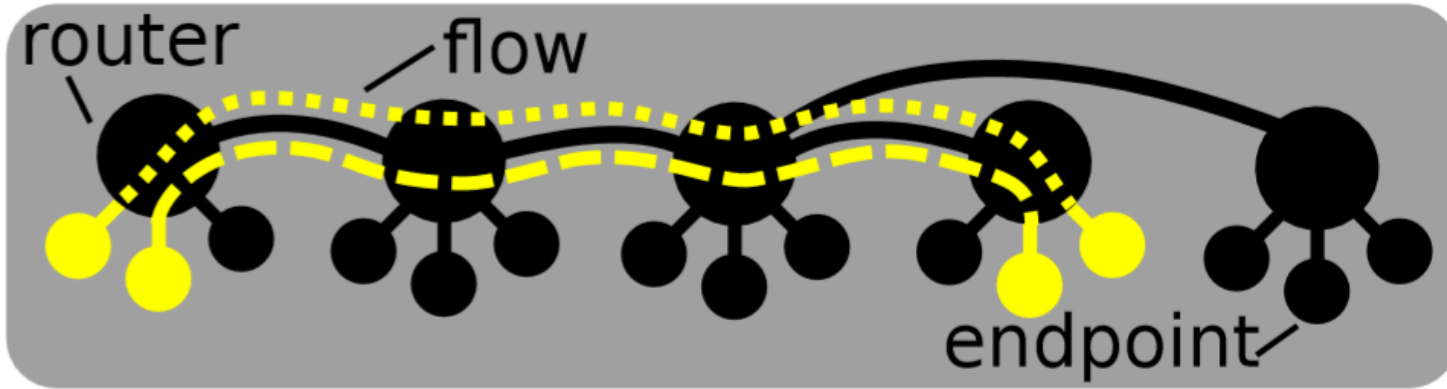
Communicating endpoint pairs mapped to the same routers

**Flow collisions**

Depend on workload mapping  
(assignment of communicating  
endpoints to routers)

# PATH DIVERSITY ANALYSIS

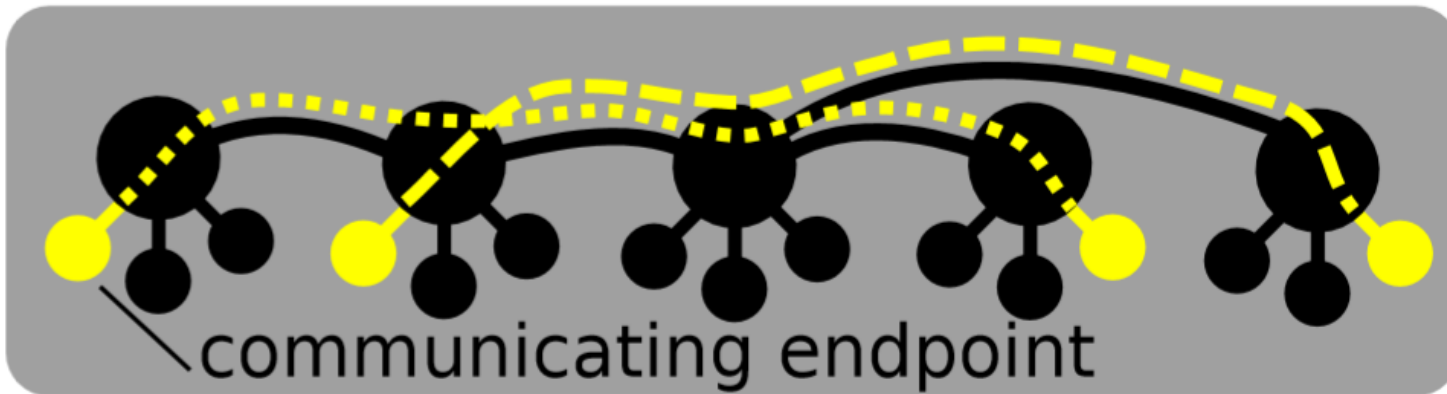
What are the problems to be tackled with multipathing?



Communicating endpoint pairs mapped to the same routers

**Flow collisions**

Depend on workload mapping  
(assignment of communicating endpoints to routers)

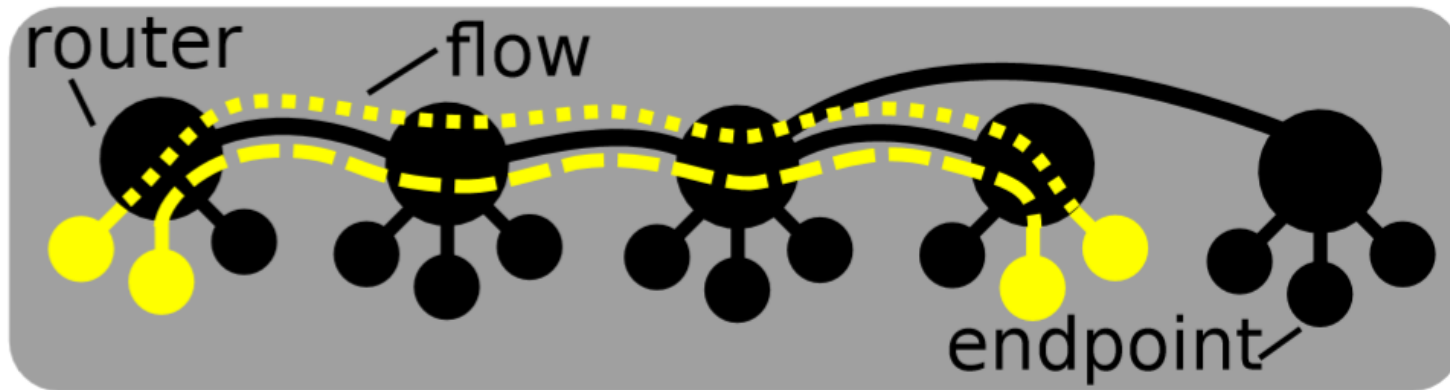


Communicating endpoint pairs mapped to different routers

**Flow overlaps**

# PATH DIVERSITY ANALYSIS

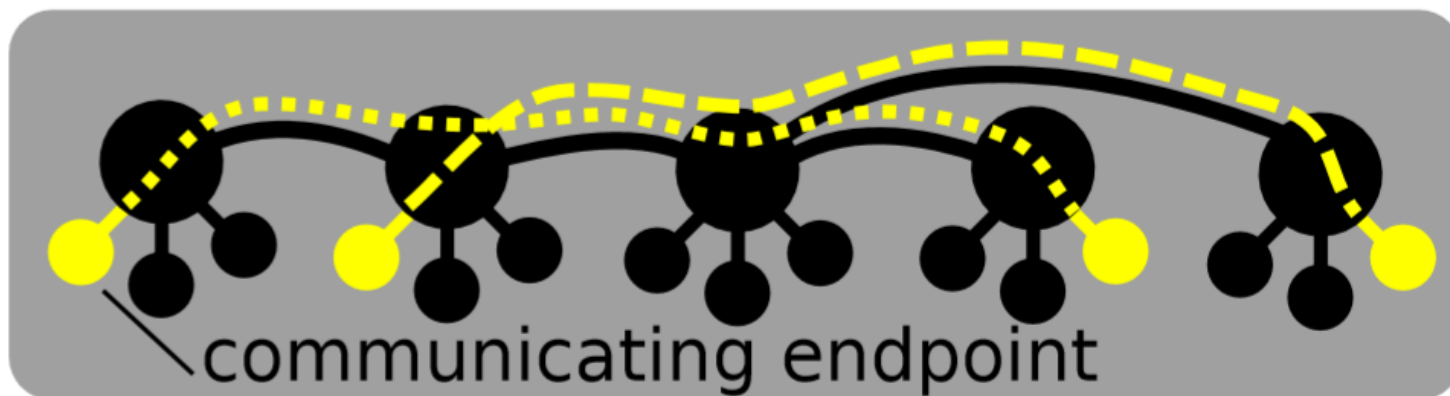
What are the problems to be tackled with multipathing?



Communicating endpoint pairs mapped to the same routers

**Flow collisions**

Depend on workload mapping  
(assignment of communicating endpoints to routers)



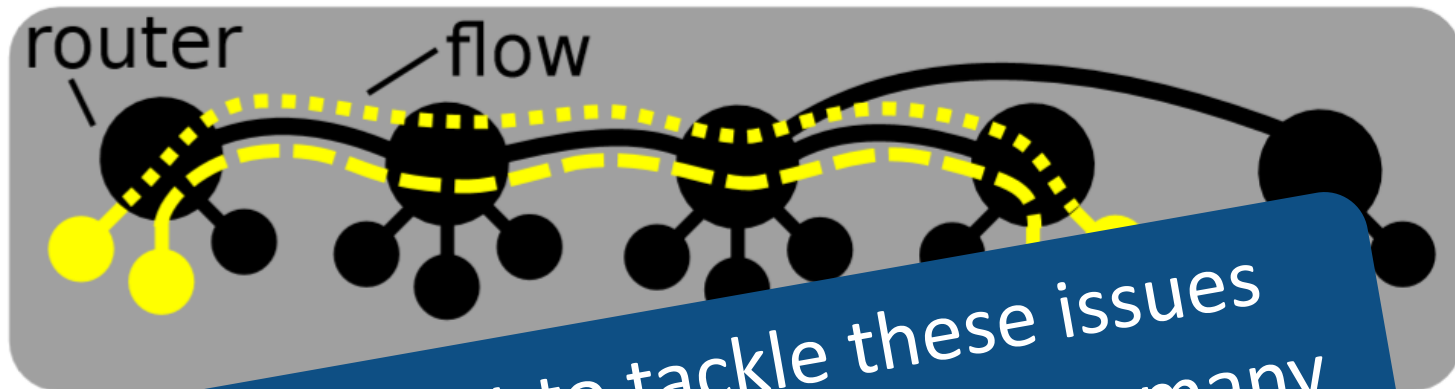
Communicating endpoint pairs mapped to different routers

**Flow overlaps**

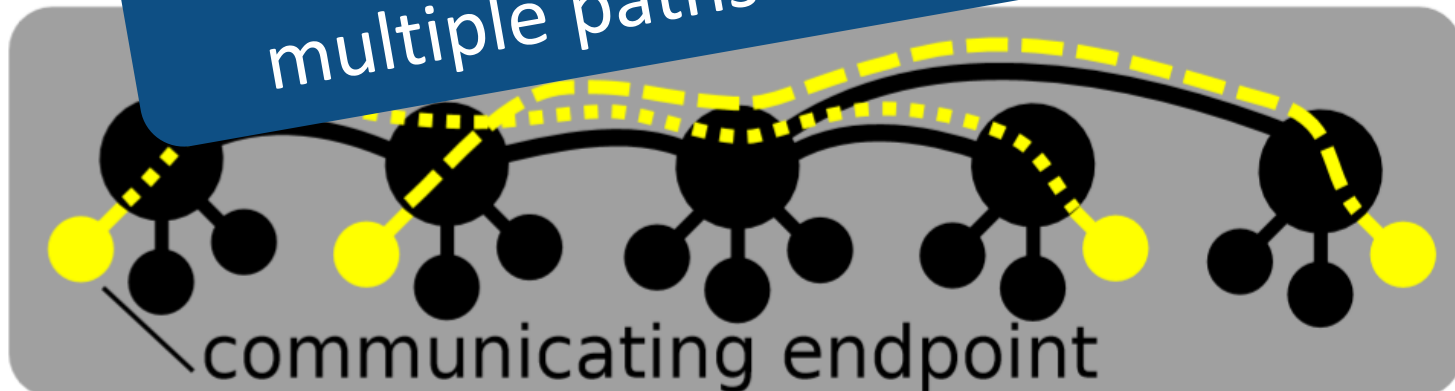
Depend on topology details  
(router-router connections)

# PATH DIVERSITY ANALYSIS

What are the problems to be tackled with multipathing?



Communicating endpoint pairs mapped to different routers



Communicating endpoint pairs mapped to different routers

**Flow collisions**

Depend on workload mapping (assignment of communicating endpoints to routers)

**Flow overlaps**

Depend on topology details (router-router connections)

We want to tackle these issues with multipathing. But how many multiple paths do we need...?

# WORKLOAD MAPPING



# WORKLOAD MAPPING

How to further enhance performance?



## WORKLOAD MAPPING

Random workload  
remapping

How to further enhance performance?



## WORKLOAD MAPPING

How to further enhance performance?



Random workload  
remapping

**Observation:** In  
low-diameter  
topologies, locality  
is less important

# WORKLOAD MAPPING

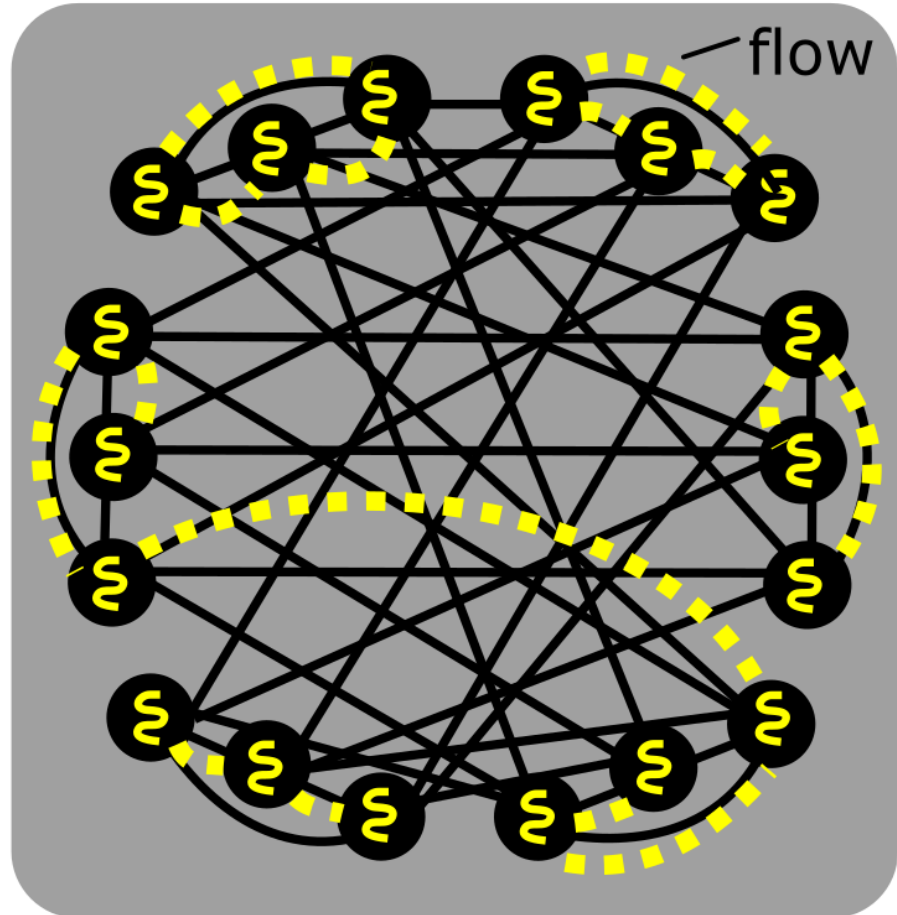
How to further enhance performance? 

Random workload remapping

Observation: In low-diameter topologies, locality is less important

A simple stencil running on a Slim Fly

Before remapping



# WORKLOAD MAPPING

How to further enhance performance?

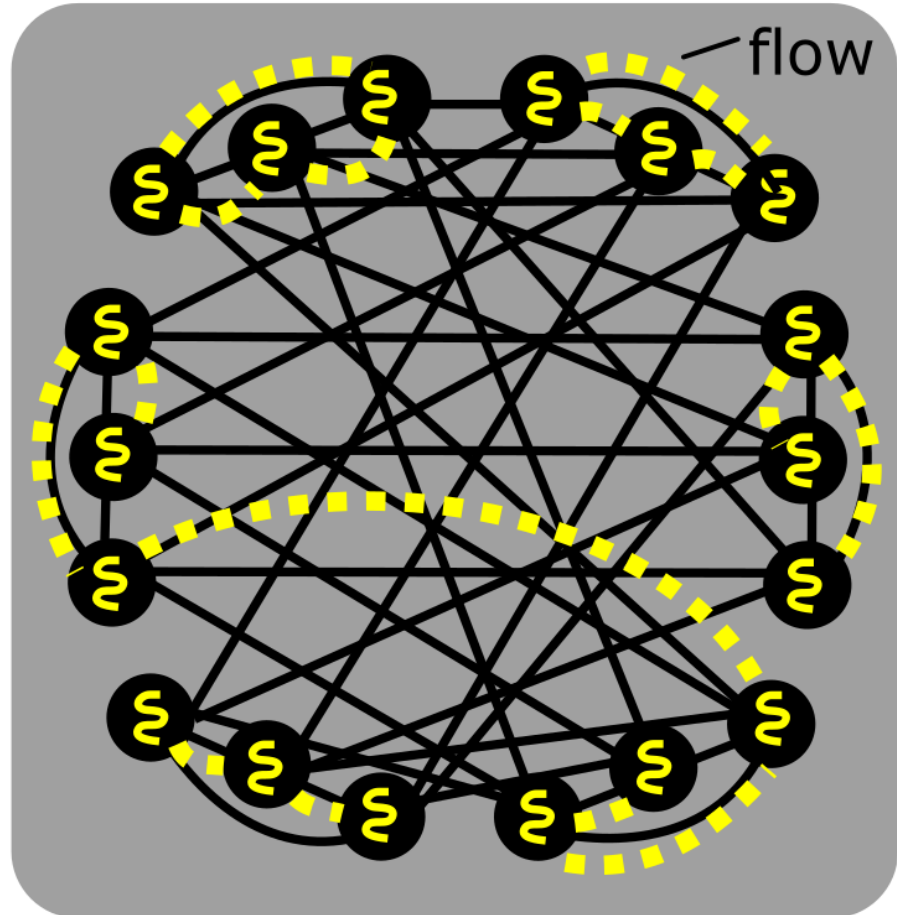
Random workload remapping

Observation: In low-diameter topologies, locality is less important

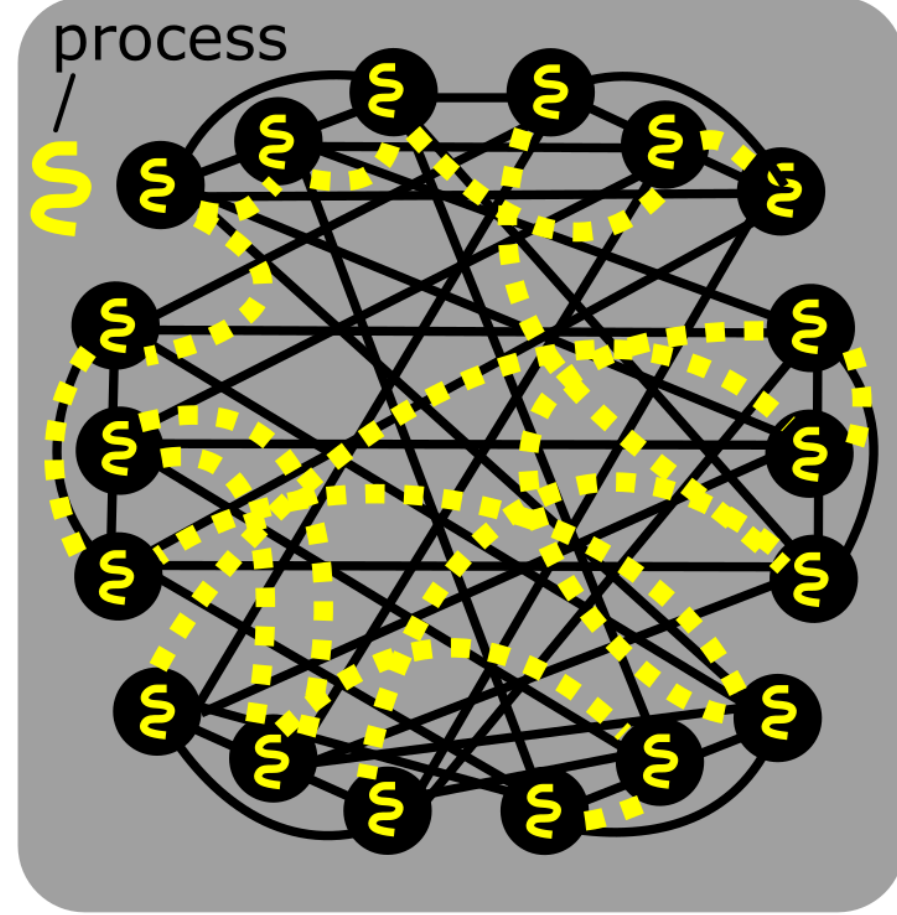
Random mapping uses rich diversity of inter-group paths

A simple stencil running on a Slim Fly

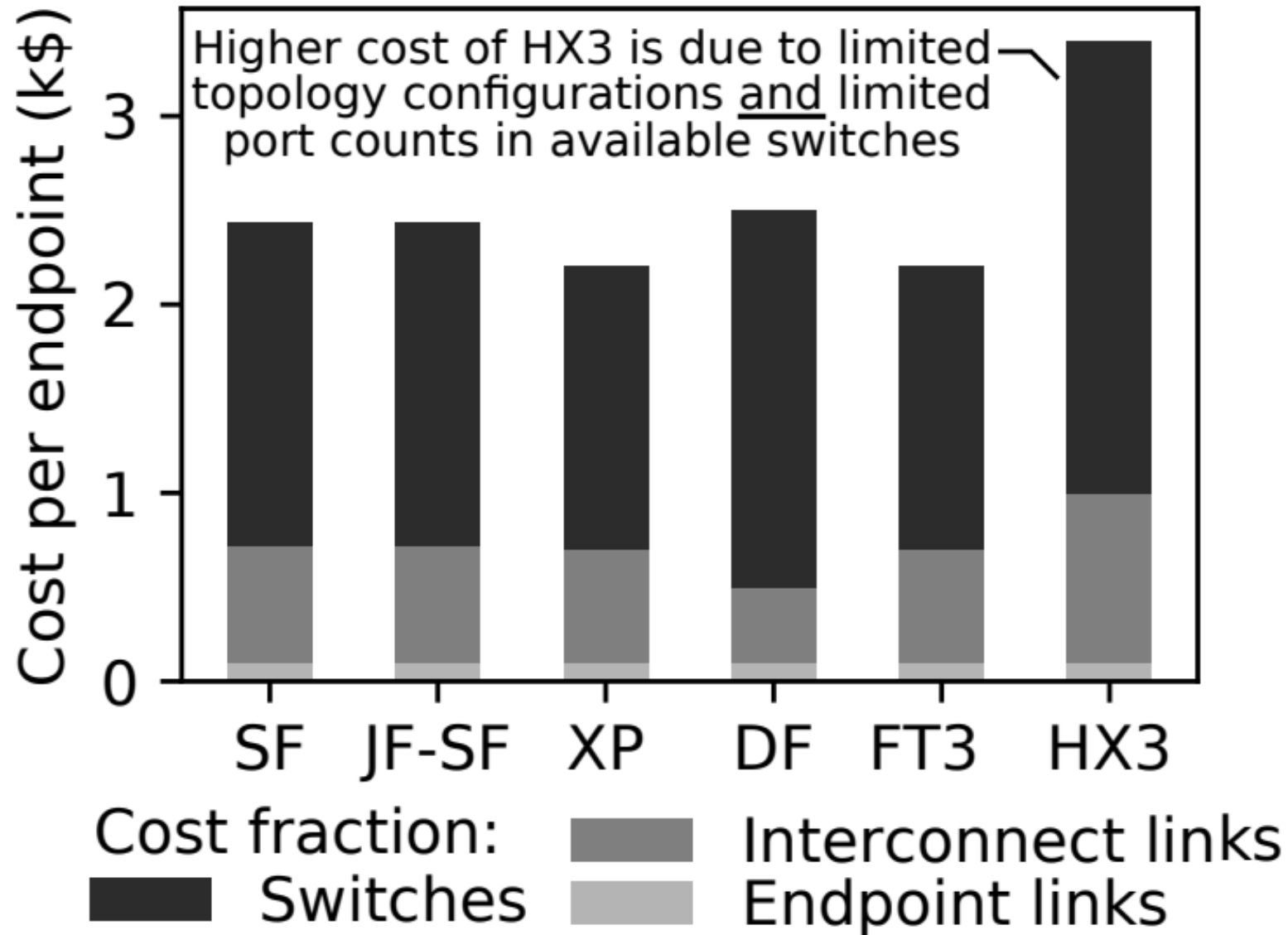
Before remapping



After remapping



# COST MODEL



# EVALUATION

## EVALUATION

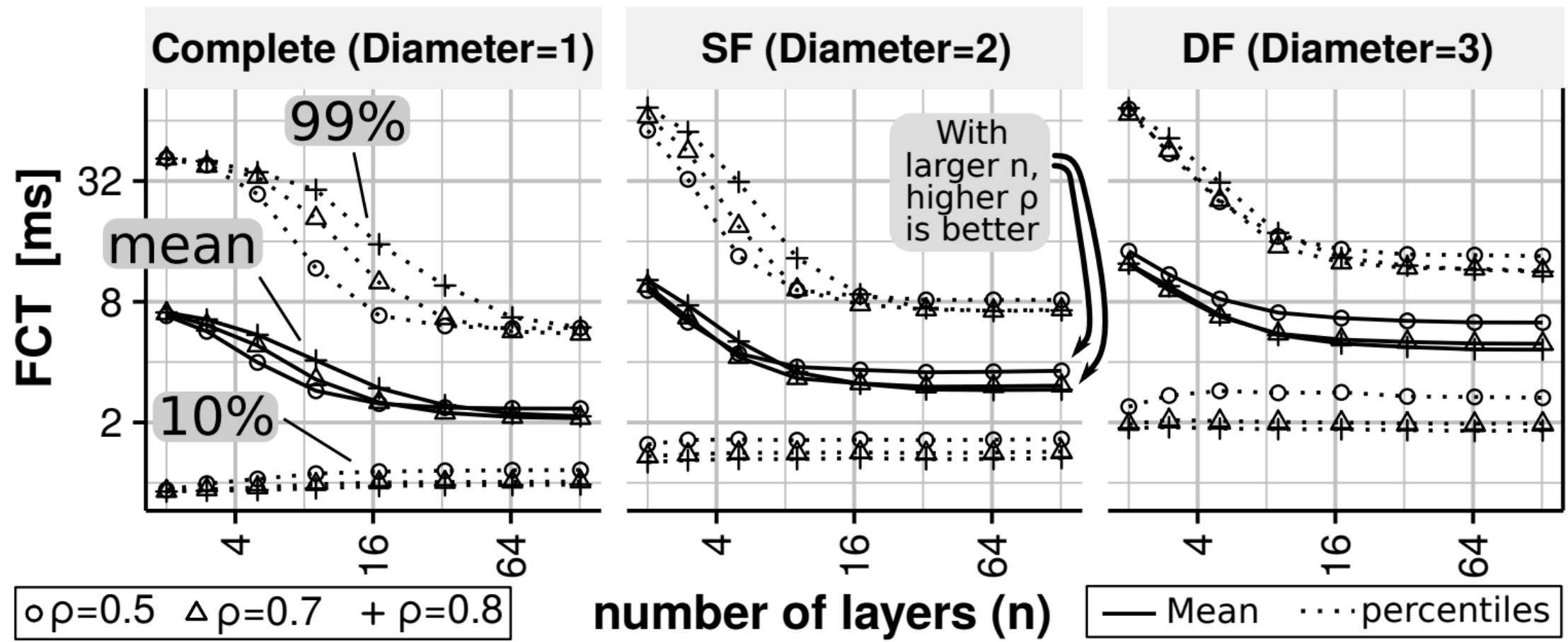
What layer setup fares best? 



# EVALUATION

What layer setup fares best?

$N \approx 10,000$ ; comparable cost; random uniform traffic; „bare Ethernet”

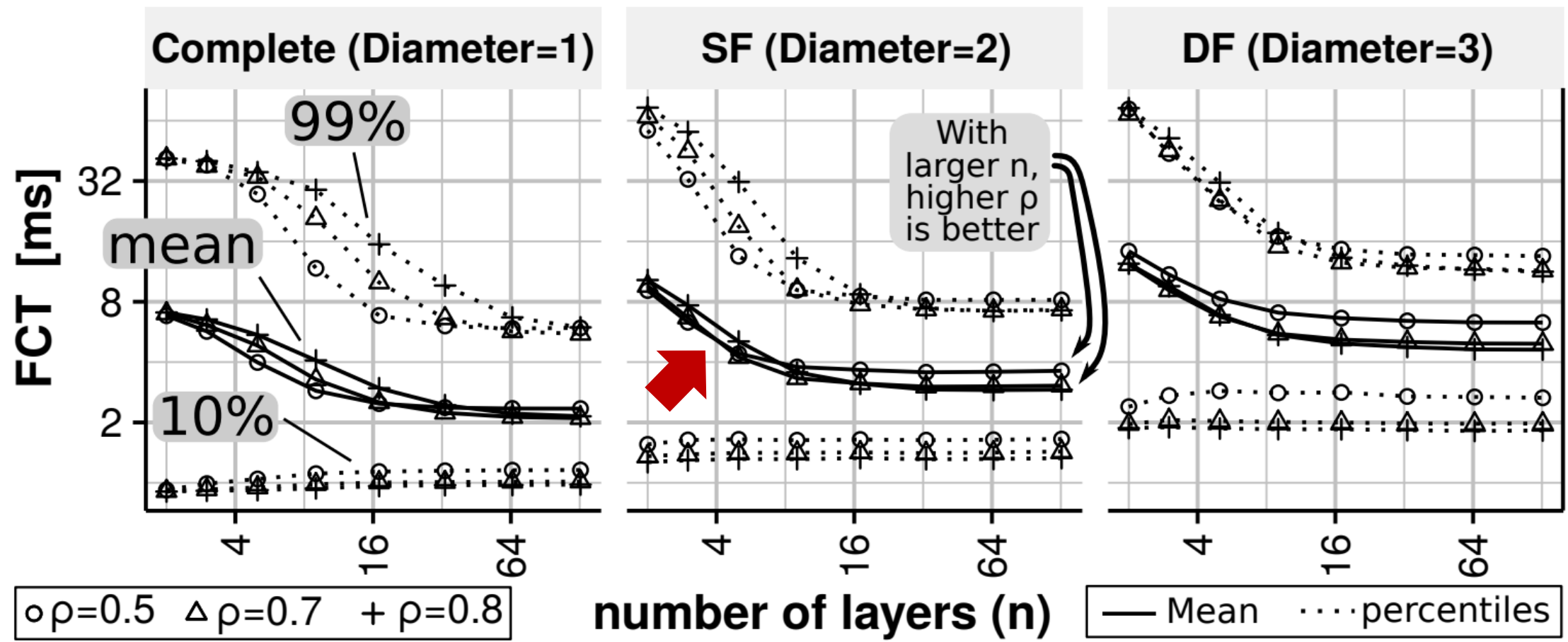


rho: fraction of links kept in a layer

# EVALUATION

What layer setup fares best?

$N \approx 10,000$ ; comparable cost; random uniform traffic; „bare Ethernet”



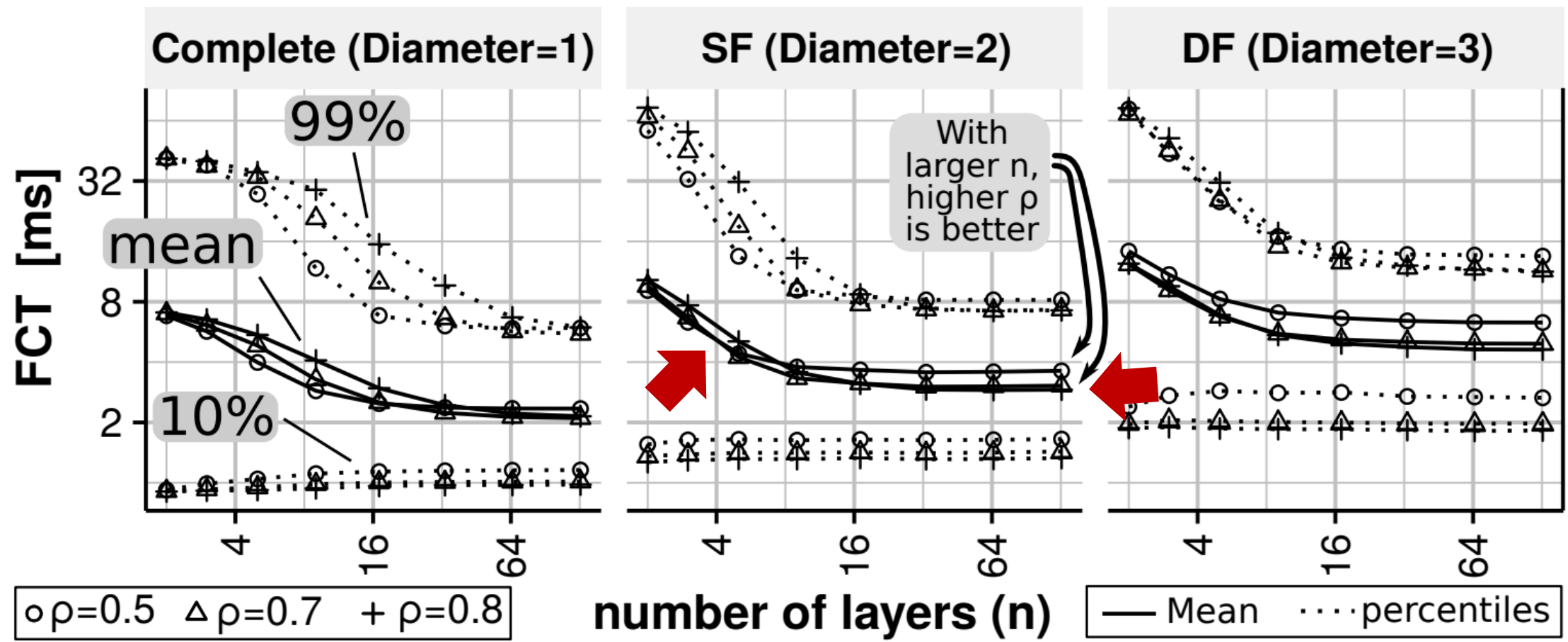
rho: fraction of links kept in a layer

With fewer layers, lower rho is better (sparser layers)

# EVALUATION

What layer setup fares best?

$N \approx 10,000$ ; comparable cost; random uniform traffic; „bare Ethernet”



**rho**: fraction of links kept in a layer

With fewer layers, lower rho is better (sparser layers)

With more layers, higher rho is better (denser layers)

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

Part 2

Non- shortest  
paths

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

Part 2

Non-shortest  
paths

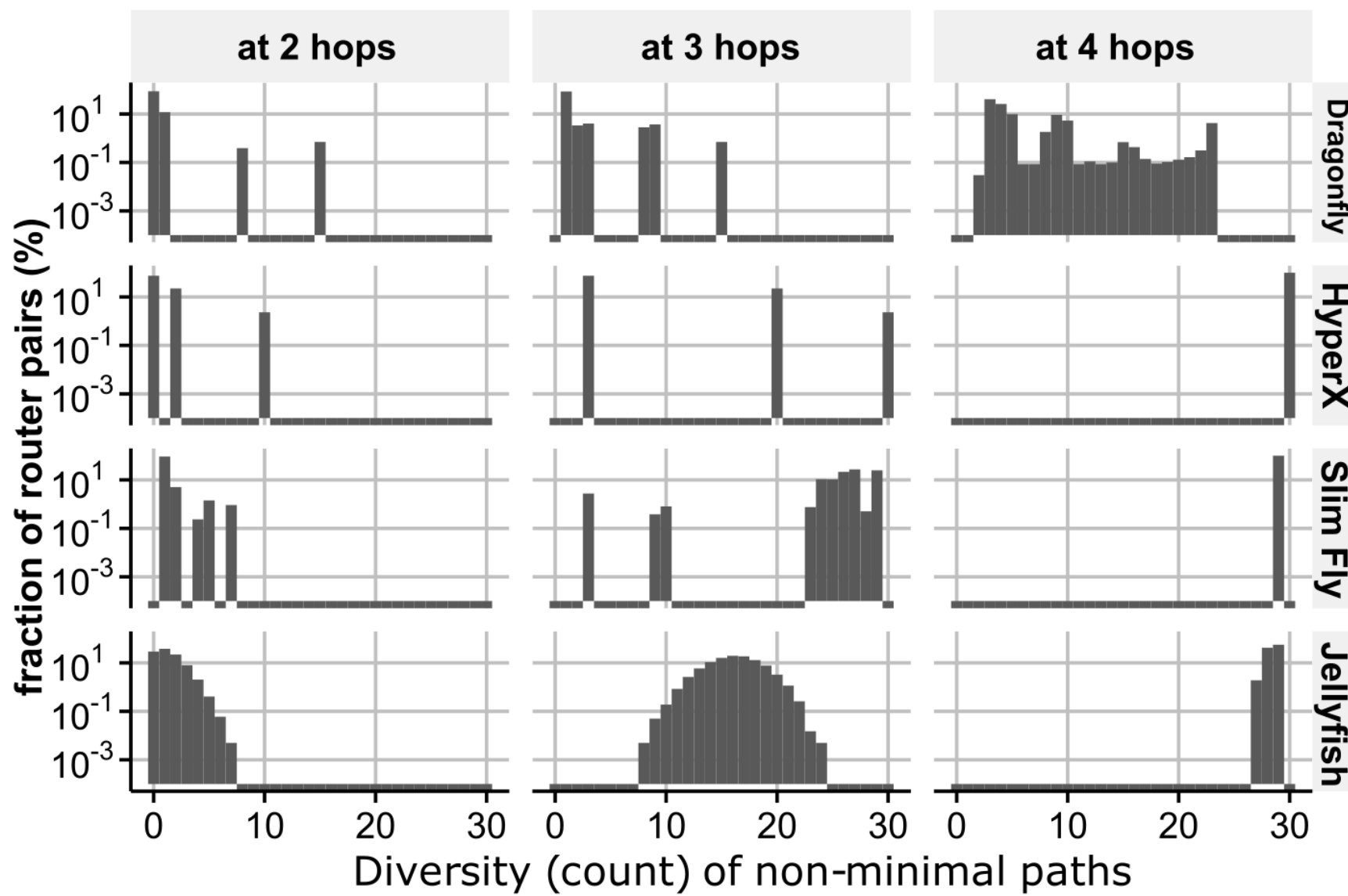
How about  
non-shortest  
paths?

# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**

How about  non-shortest paths?

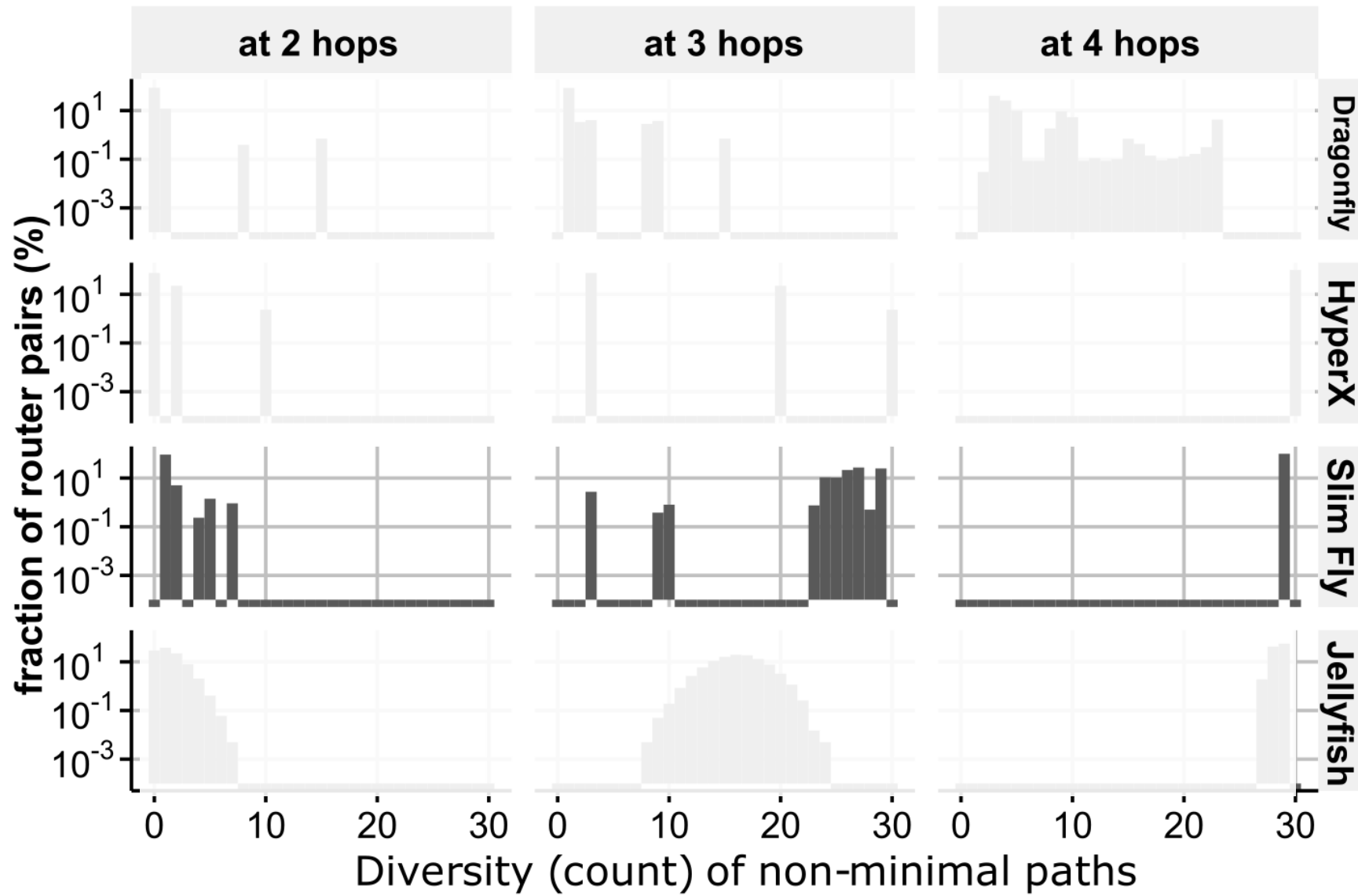


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non- shortest paths**

How about  non-shortest paths?

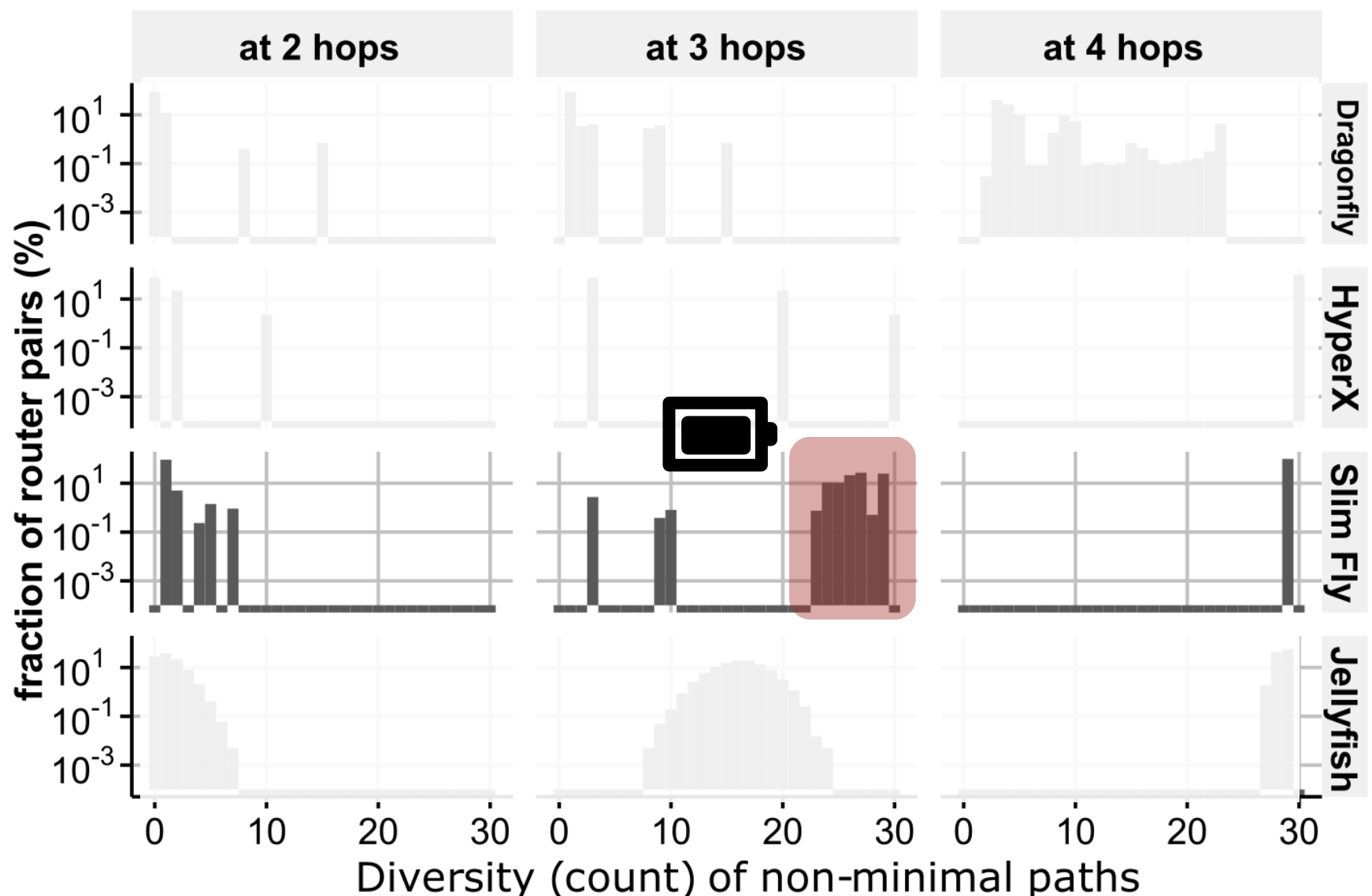


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non- shortest paths**

How about  non-shortest paths?



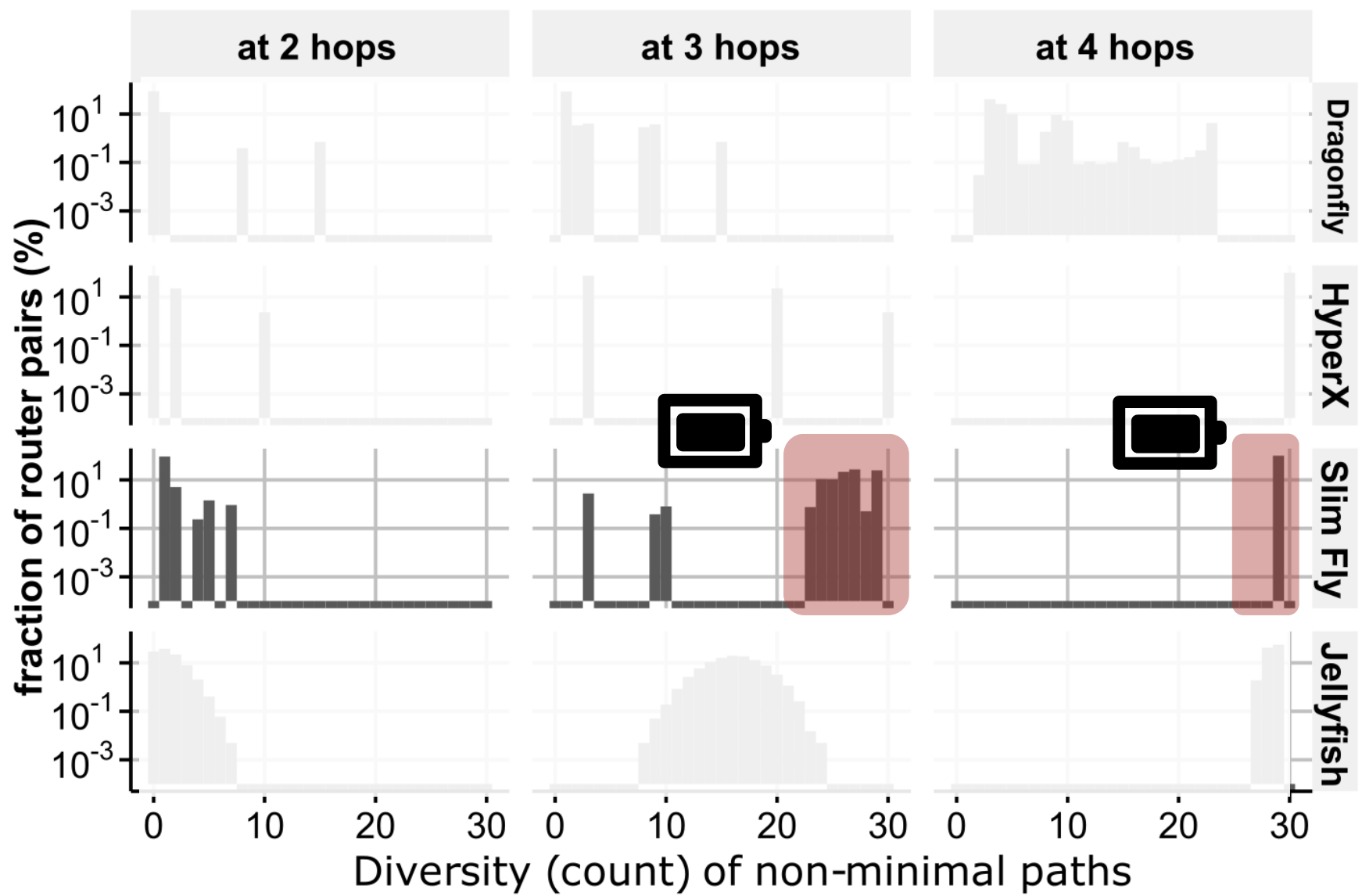


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

## Part 2

### Non-shortest paths

How about  non-shortest paths?

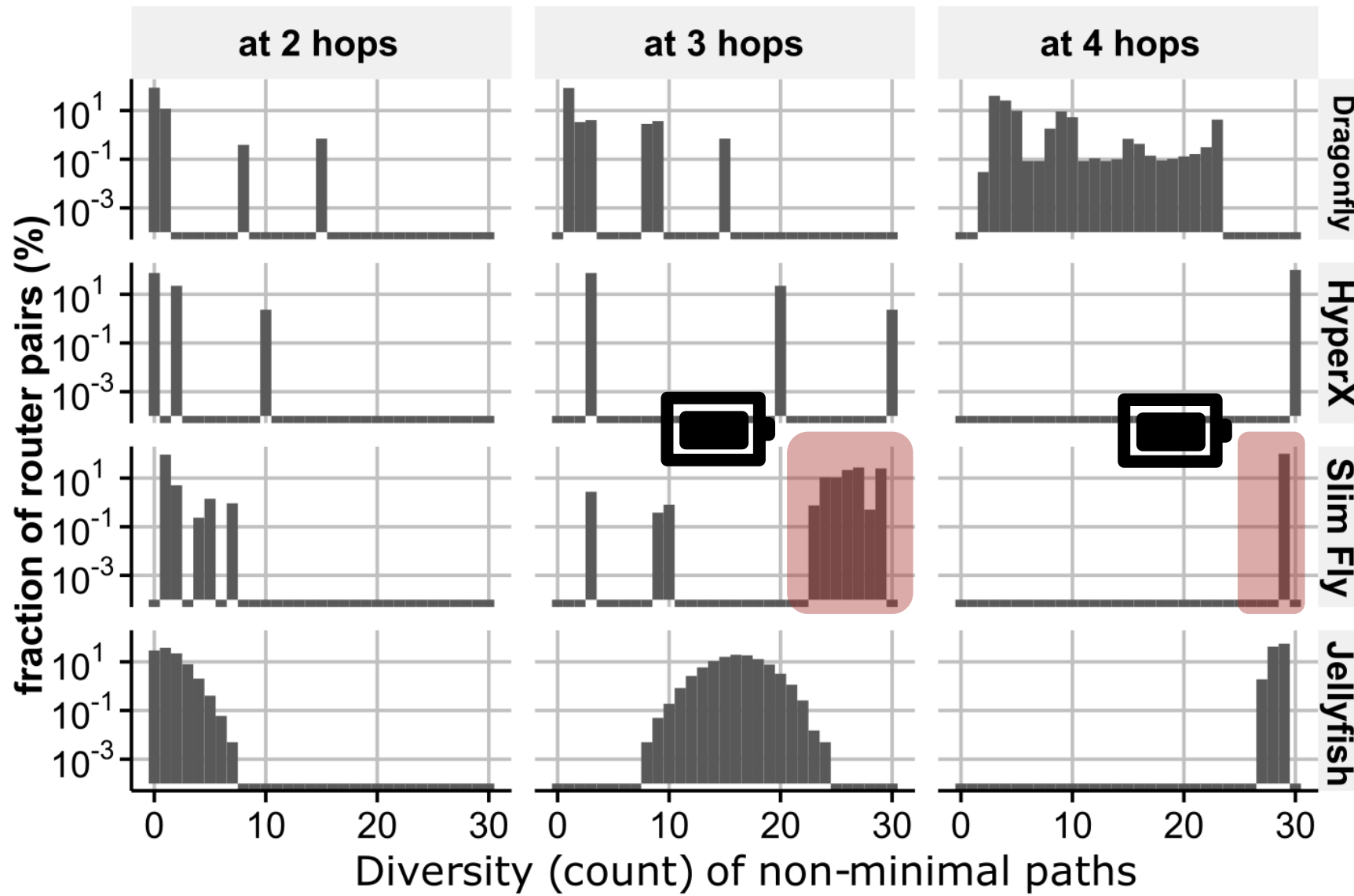


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non-shortest paths**

How about  non-shortest paths?

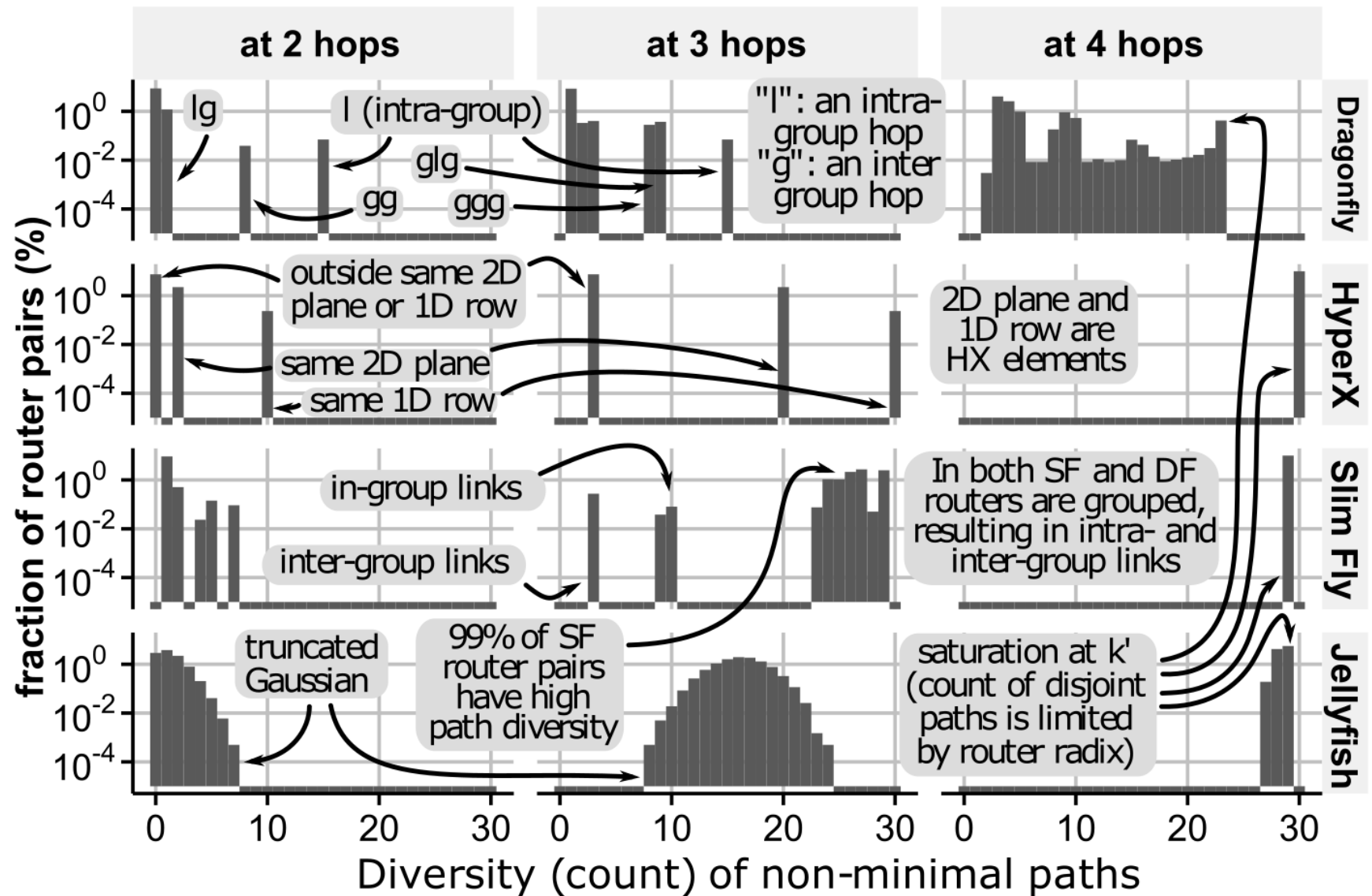


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

## Part 2

### Non-shortest paths

How about non-shortest paths?

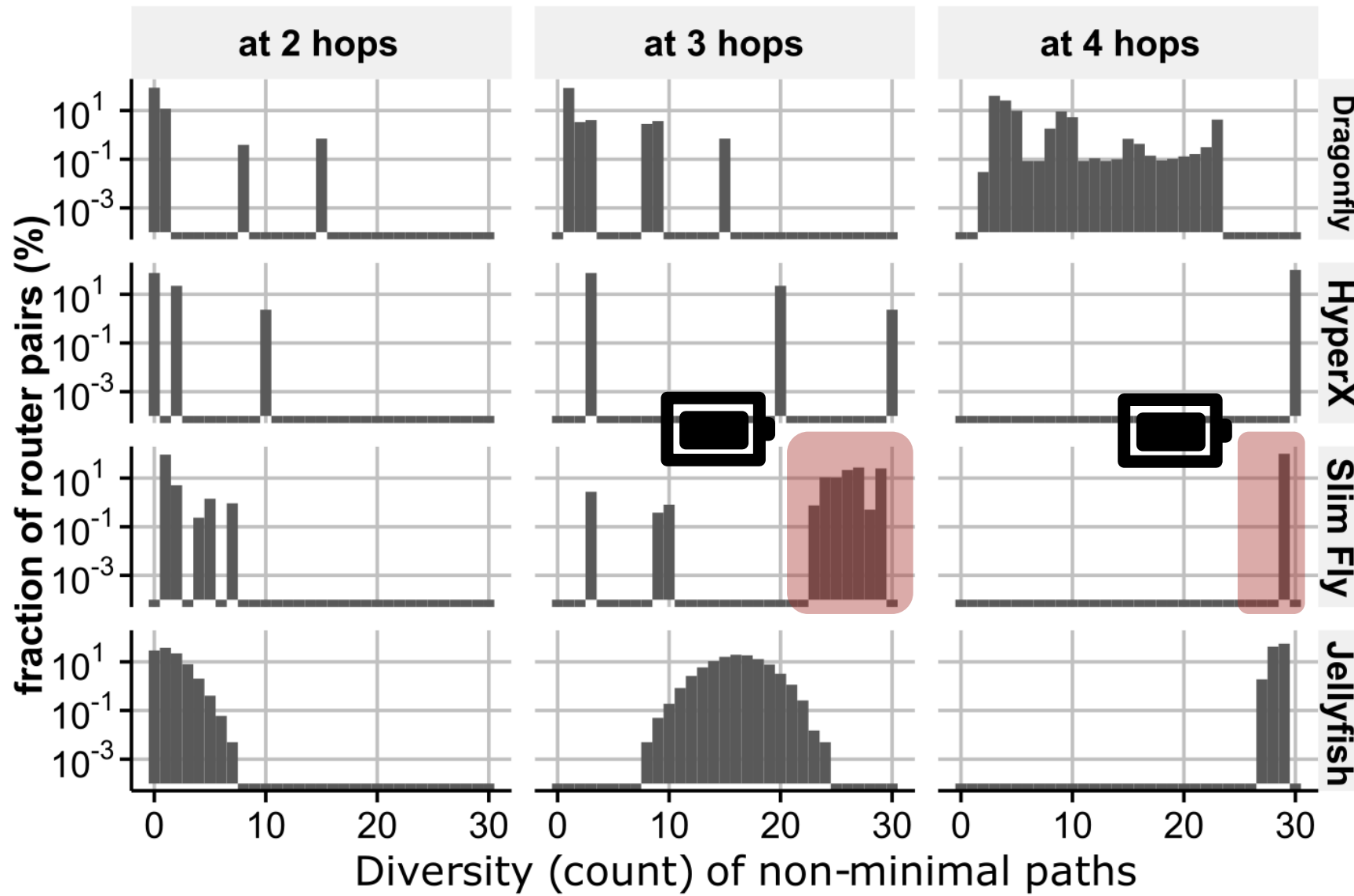


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

**Part 2**

**Non- shortest paths**

How about  non-shortest paths?

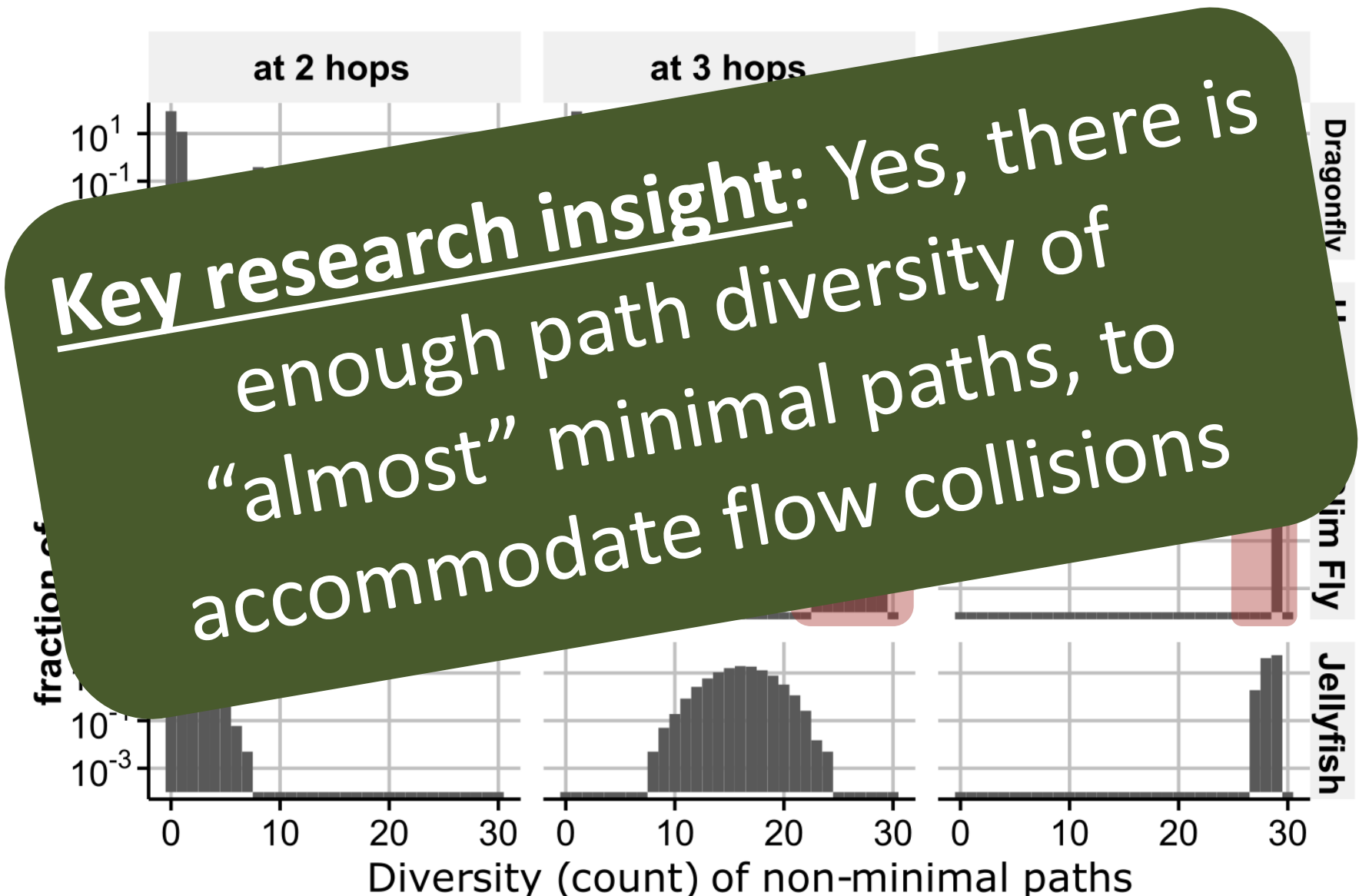


# ARE THERE ENOUGH MULTIPLE PATHS IN LOW-DIAMETER TOPOLOGIES?

## Part 2

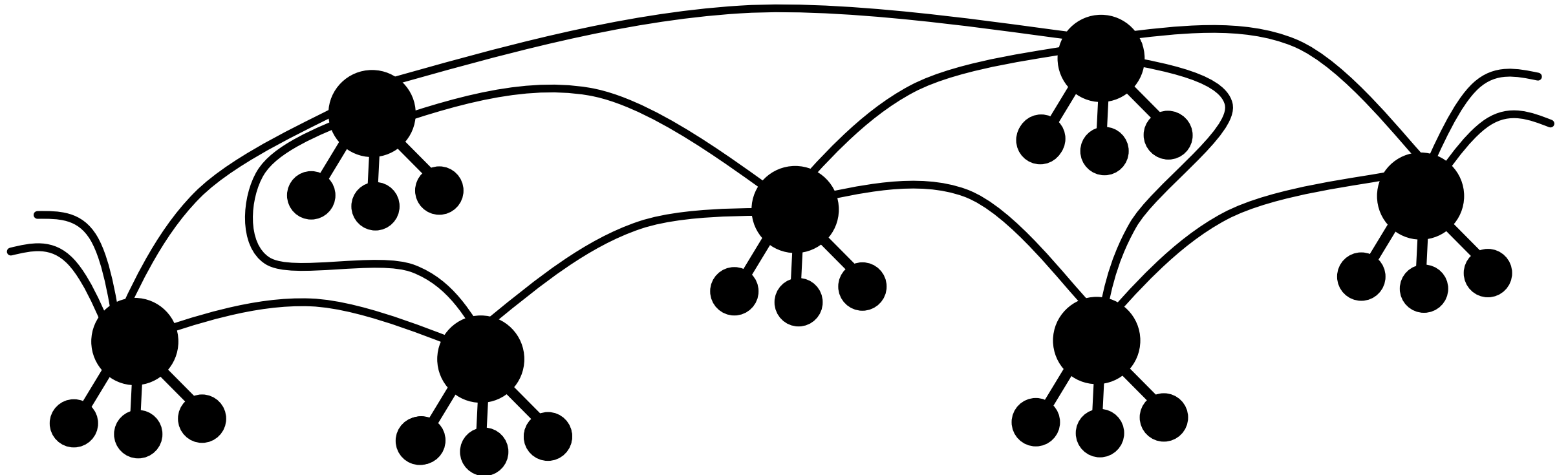
### Non-shortest paths

How about  non-shortest paths?

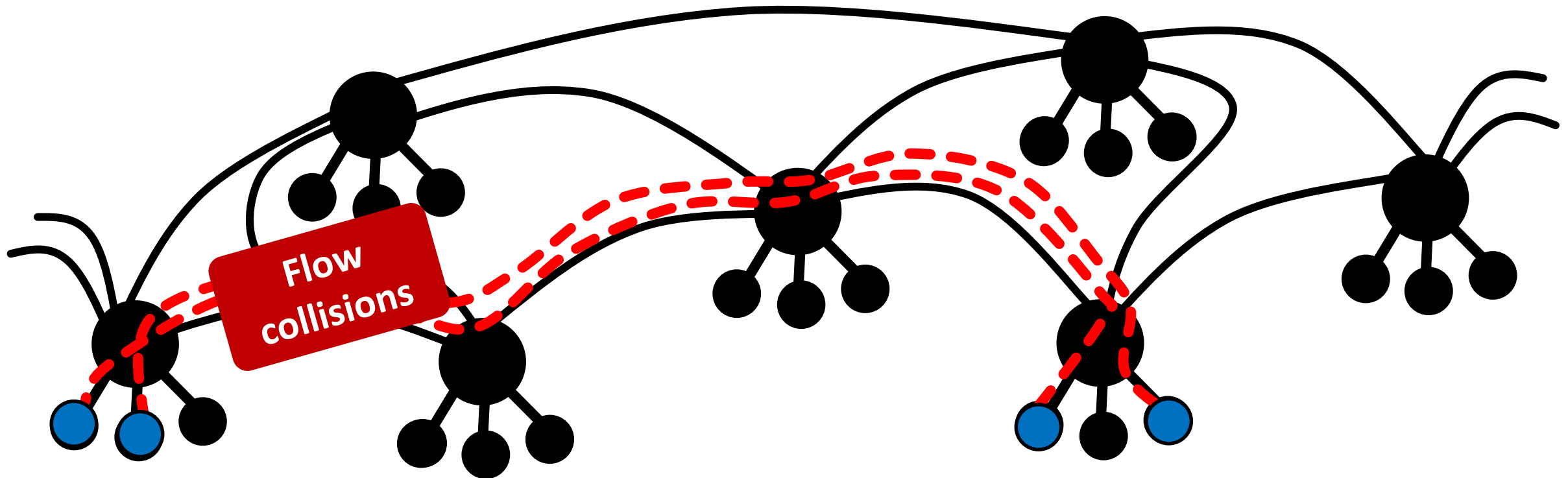


**Key research insight:** Yes, there is enough path diversity of “almost” minimal paths, to accommodate flow collisions

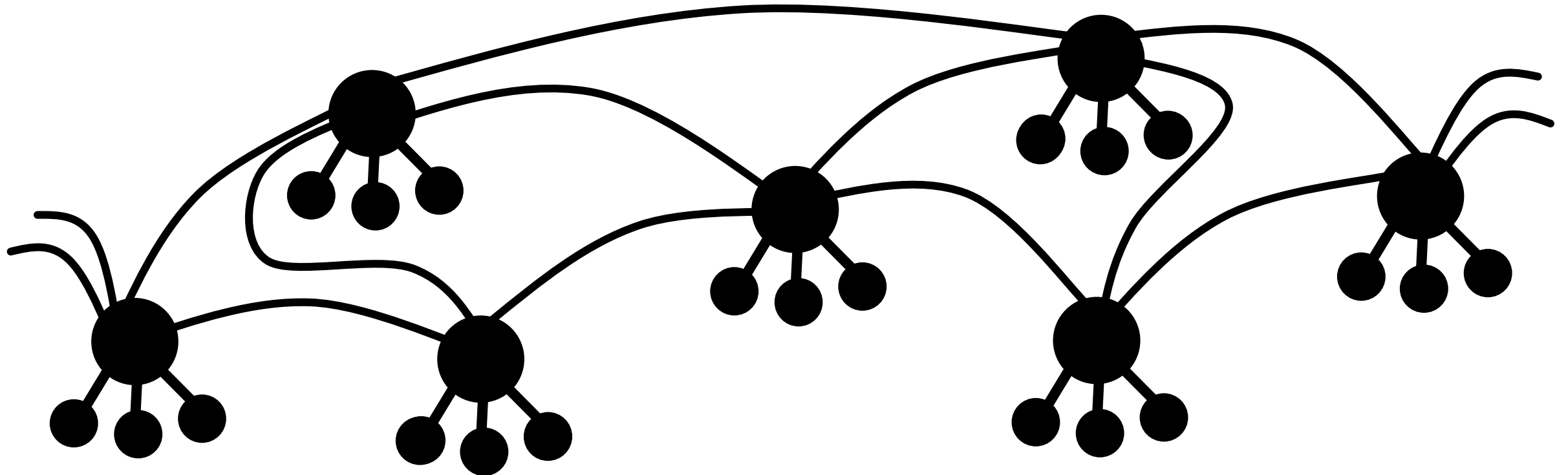
THERE IS MUCH MORE...



THERE IS MUCH MORE...

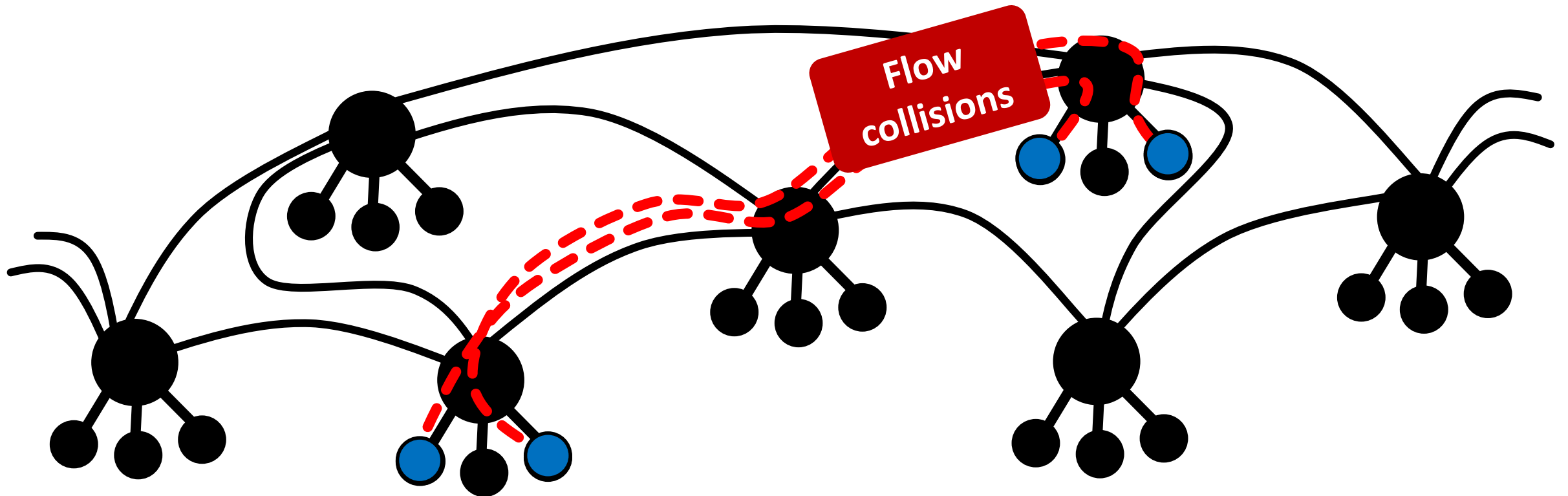


THERE IS MUCH MORE...

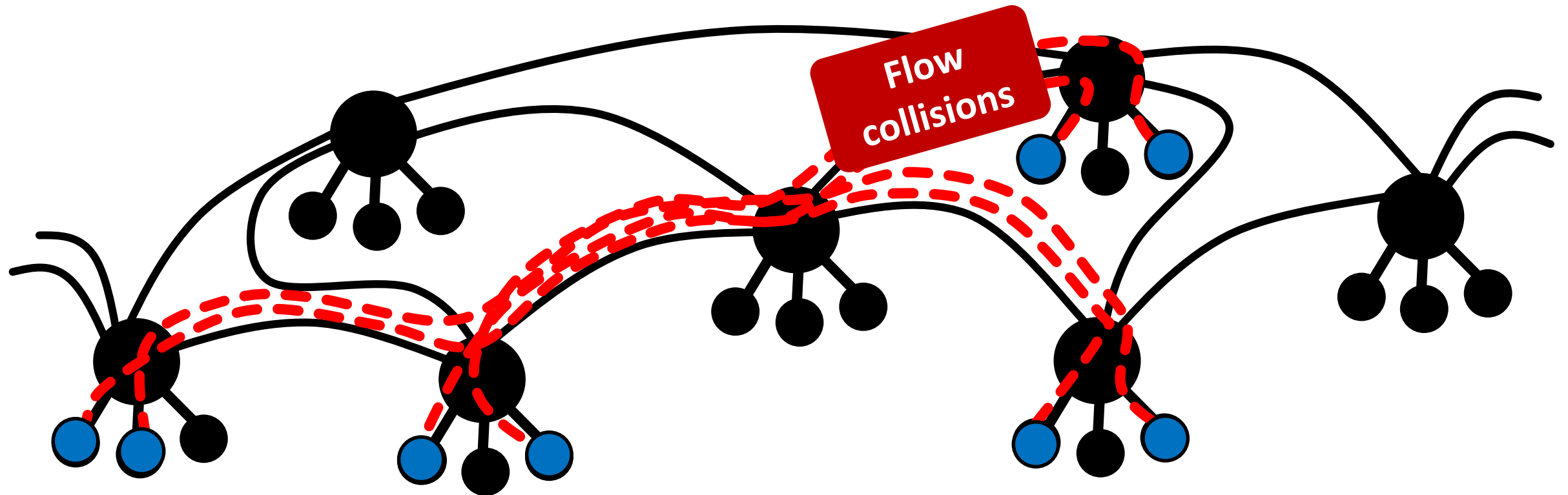




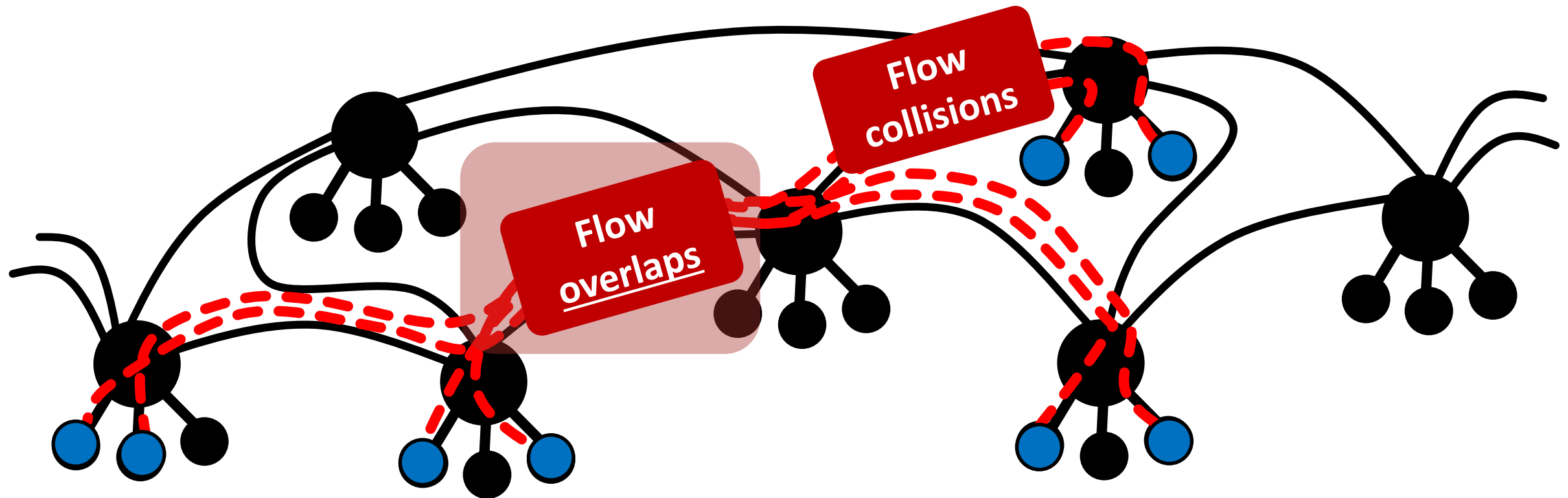
THERE IS MUCH MORE...



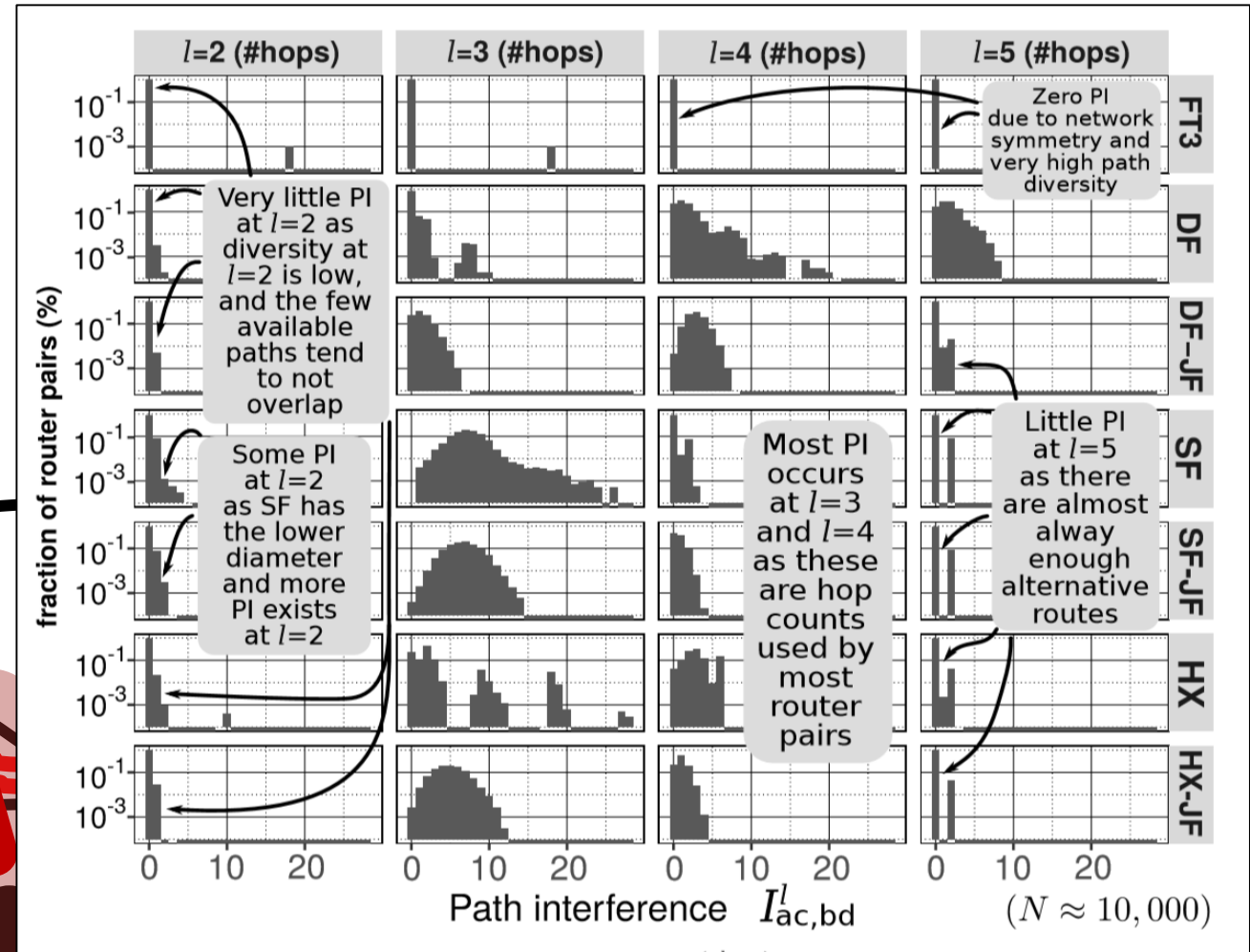
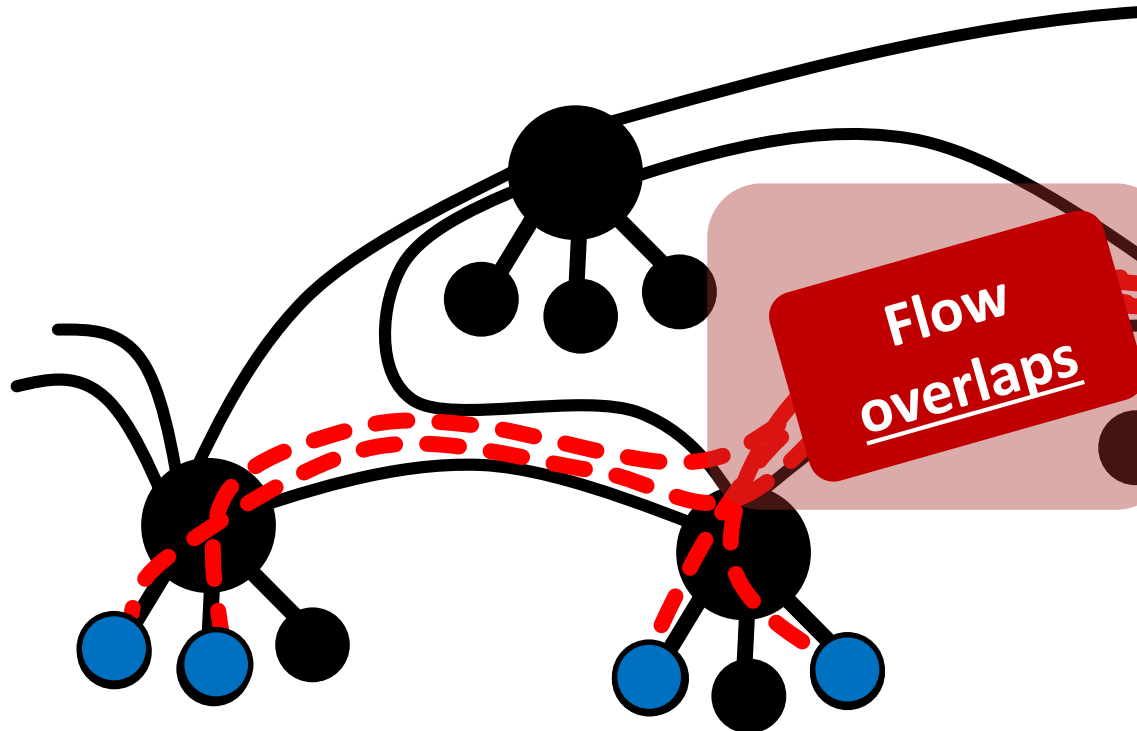
THERE IS MUCH MORE...



THERE IS MUCH MORE...



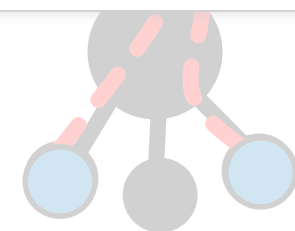
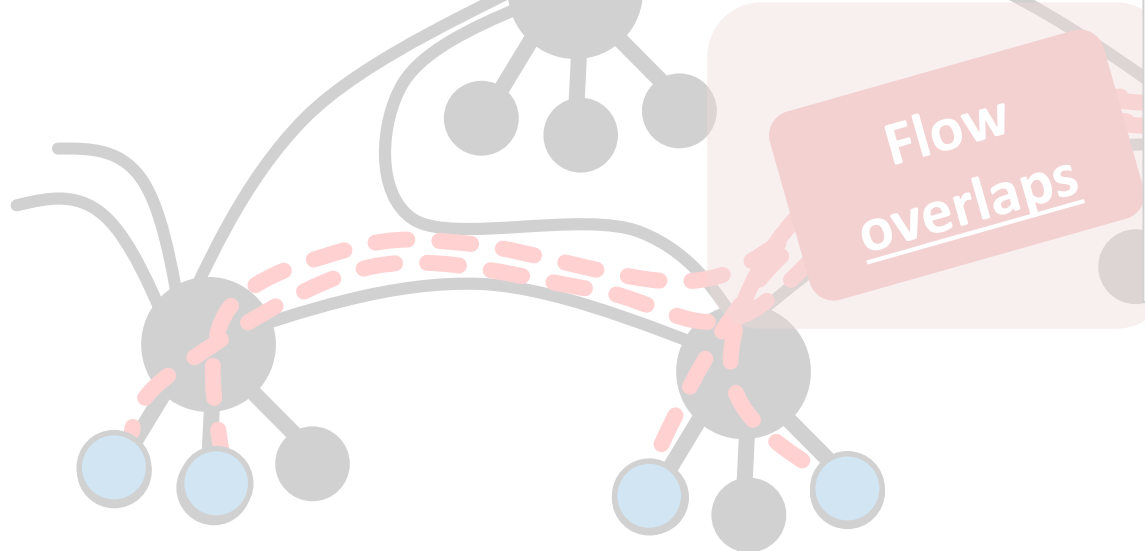
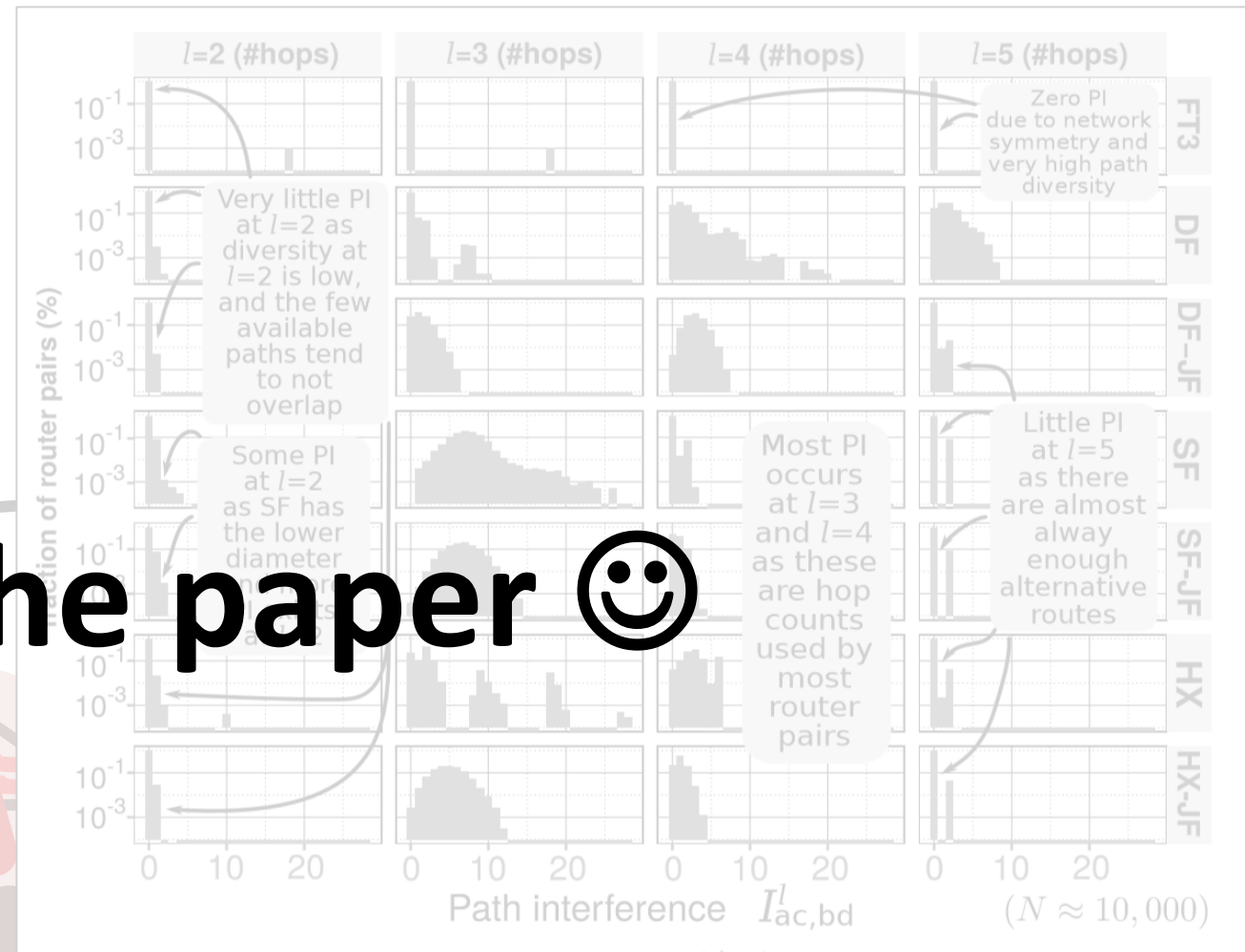
# THERE IS MUCH MORE...



THERE IS MUCH MORE...

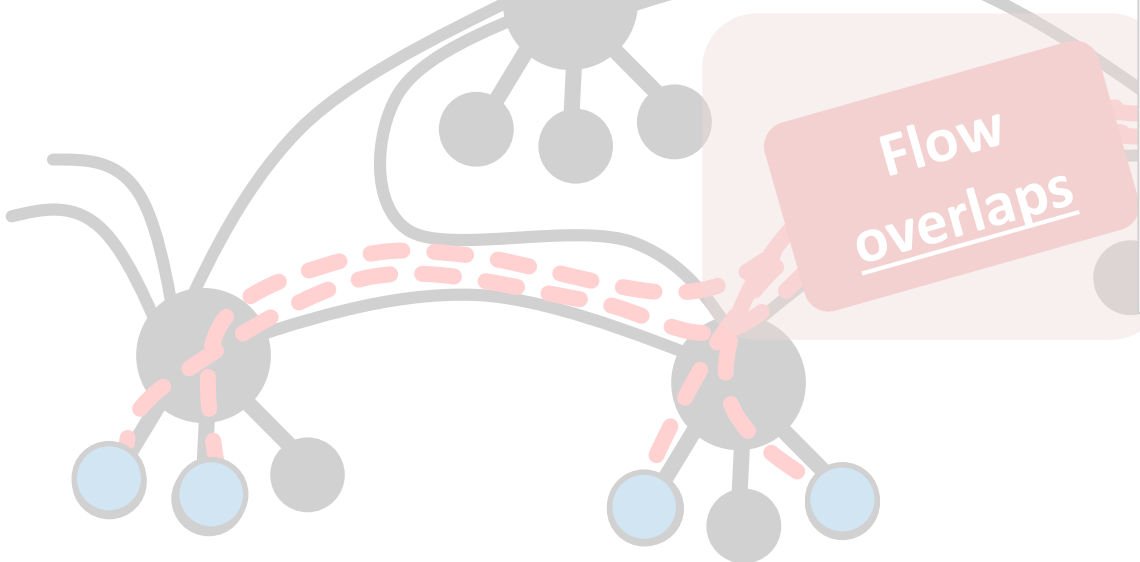
Details in the paper 😊

Flow overlaps



THERE IS MUCH MORE...

Details in the paper 😊



Let's sum up...

# FATPATHS ARCHITECTURE: WHICH TECHNOLOGY?

## FATPATHS ARCHITECTURE: WHICH TECHNOLOGY?



InfiniBand



# FATPATHS ARCHITECTURE: WHICH TECHNOLOGY?



InfiniBand



Myrinet

# FATPATHS ARCHITECTURE: WHICH TECHNOLOGY?



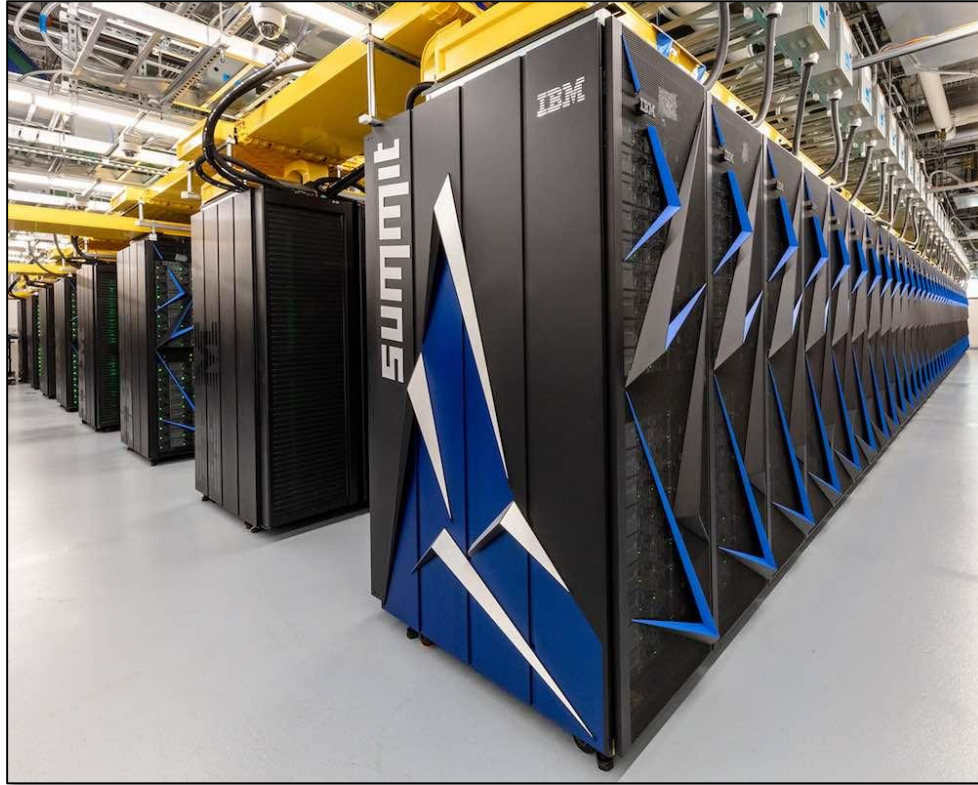
InfiniBand

Ethernet



Myrinet

# FATPATHS ARCHITECTURE: WHICH TECHNOLOGY?



InfiniBand



Others

Ethernet

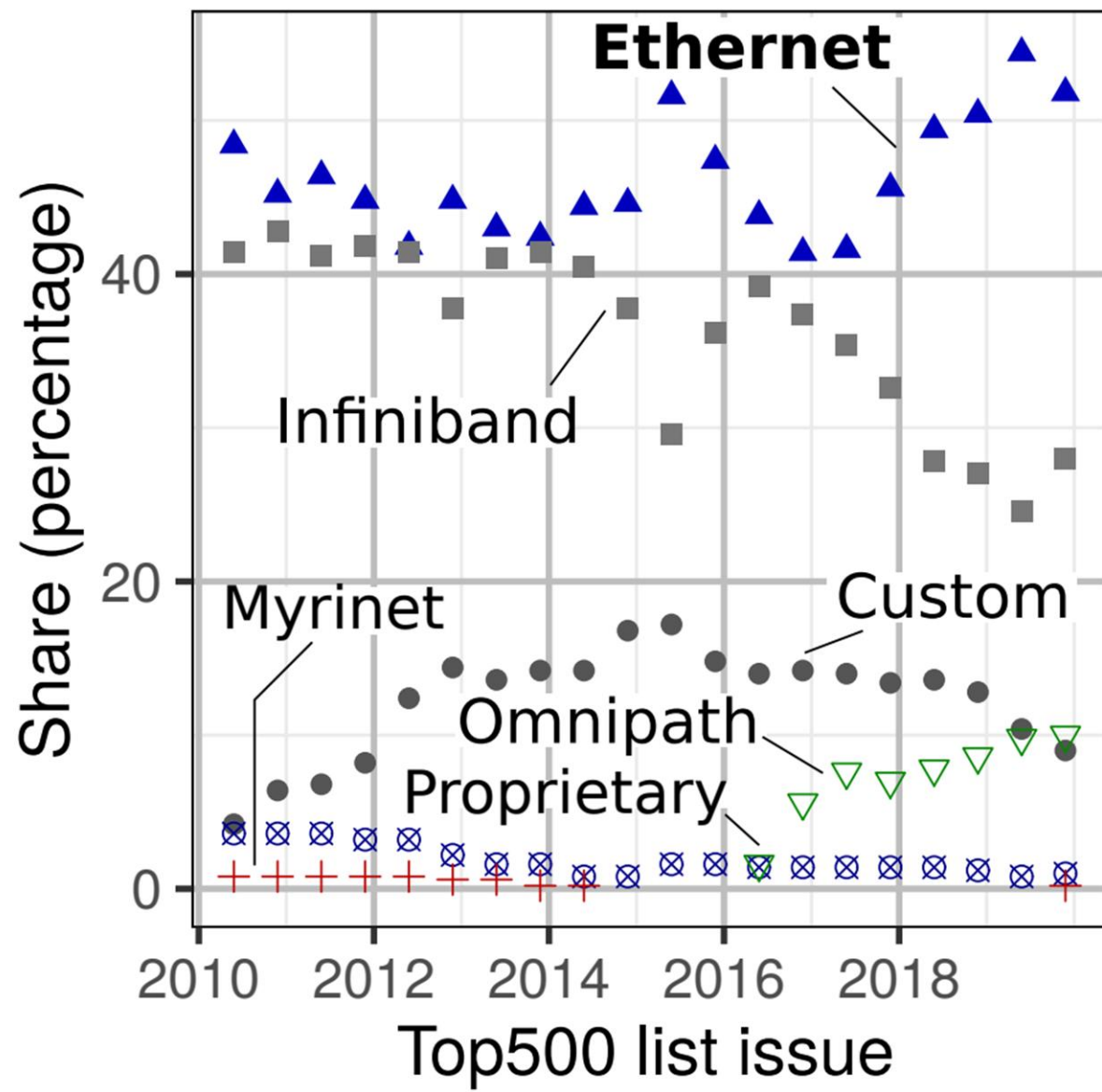


Myrinet

# ETHERNET & HPC

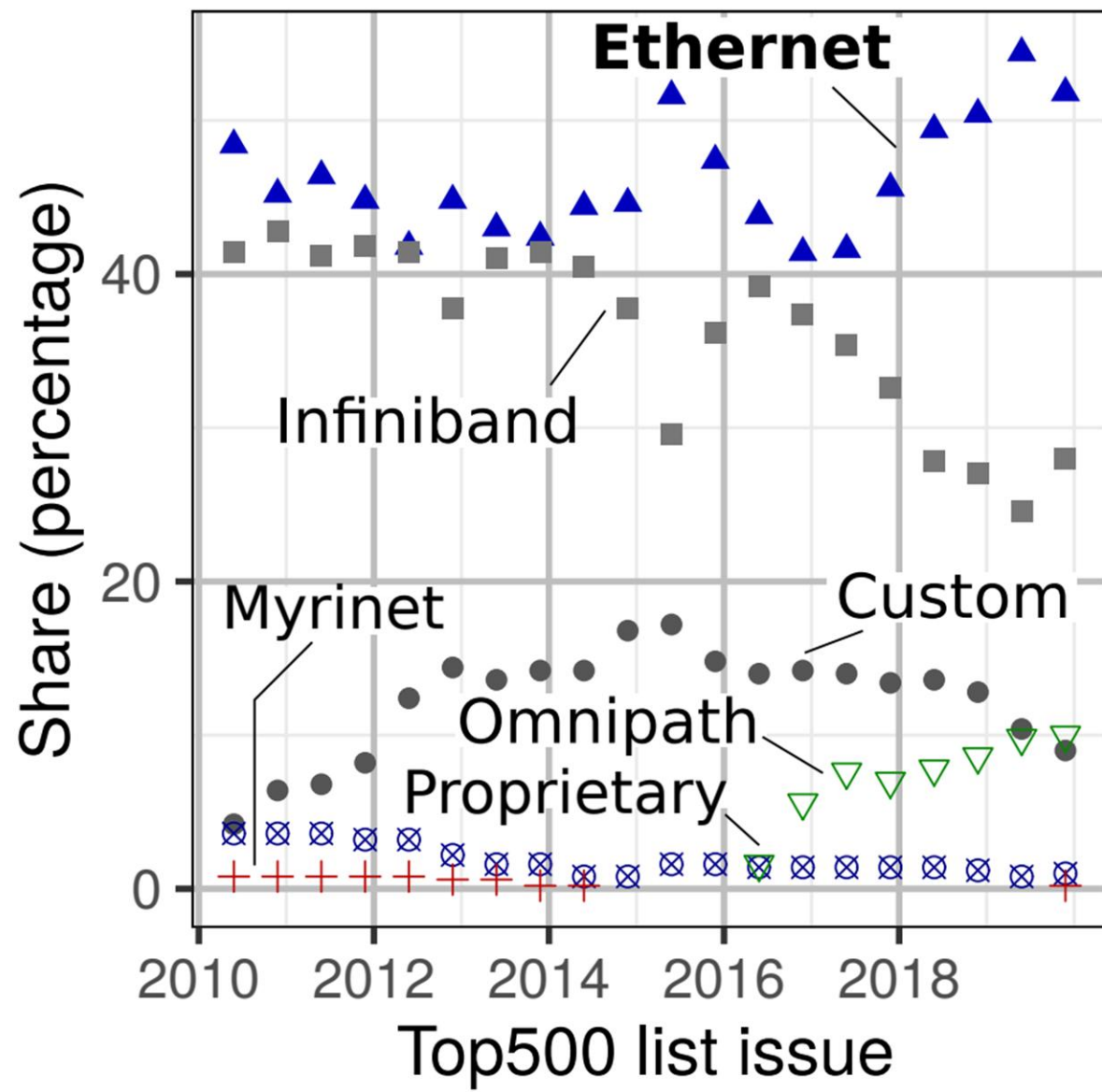


# ETHERNET & HPC



# ETHERNET & HPC

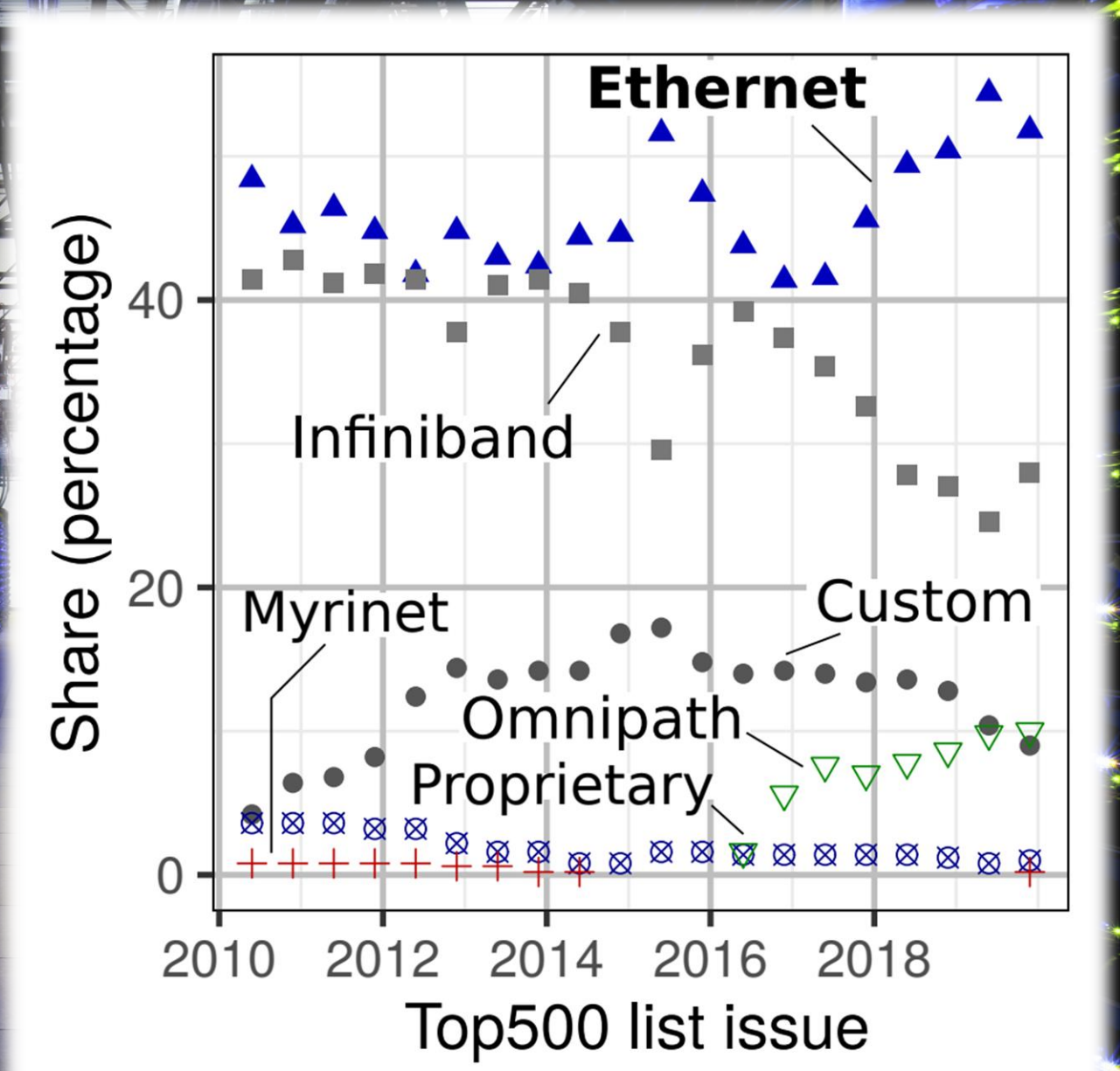
More than 50% of systems in Top500 use Ethernet



# ETHERNET & HPC

More than 50% of systems in Top500 use Ethernet

However, systems based on Ethernet are not as efficient as others:  $\approx$  (details in the paper)

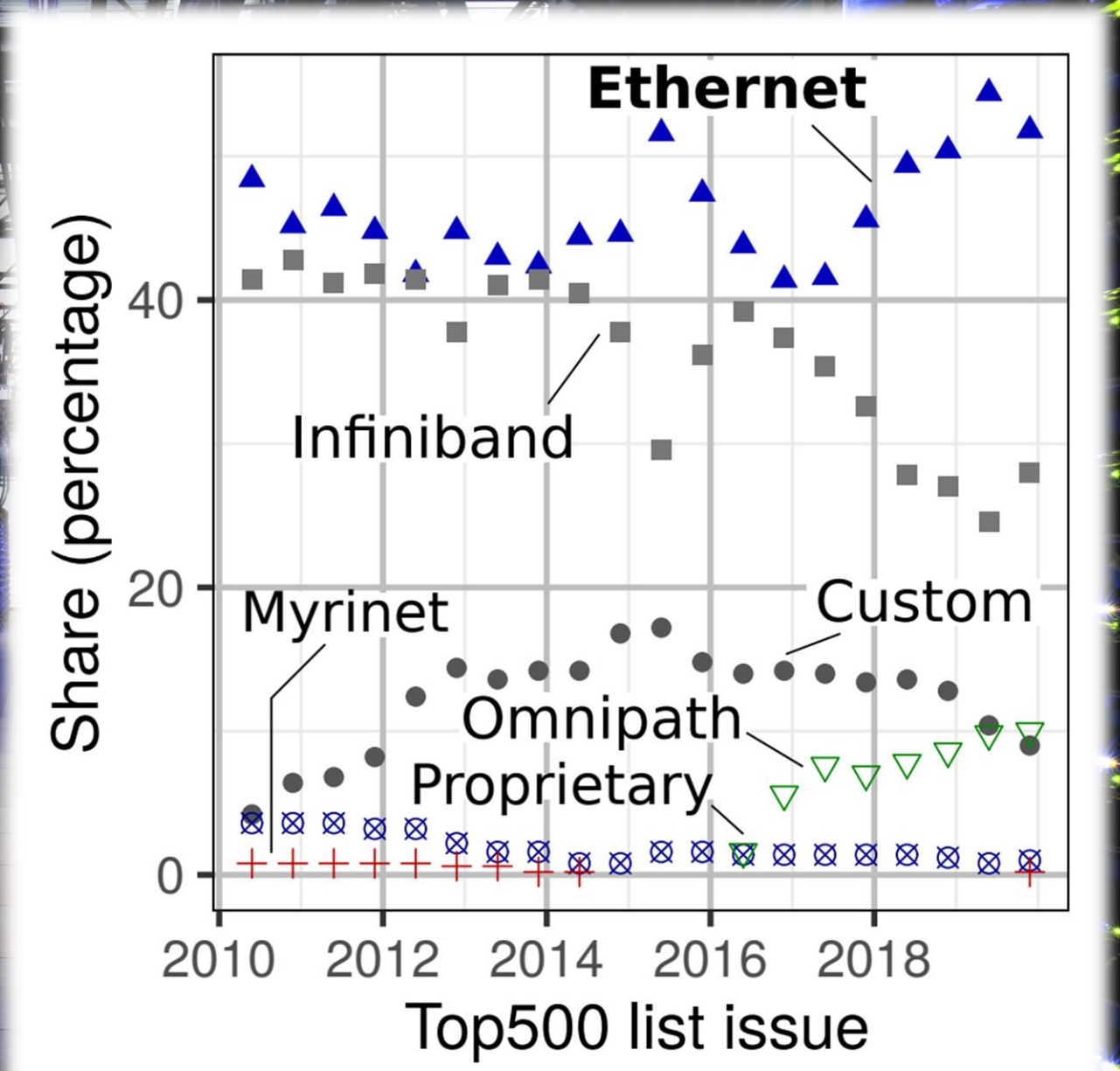


# ETHERNET & HPC

More than 50% of systems in Top500 use Ethernet

However, systems based on Ethernet are not as efficient as others:  $\approx$  (details in the paper)

LINPACK Efficiency:  
 $\approx$ 51% for 100G Ethernet  
 $\approx$ 65% for 100G IB HDR  
 (details in the paper)



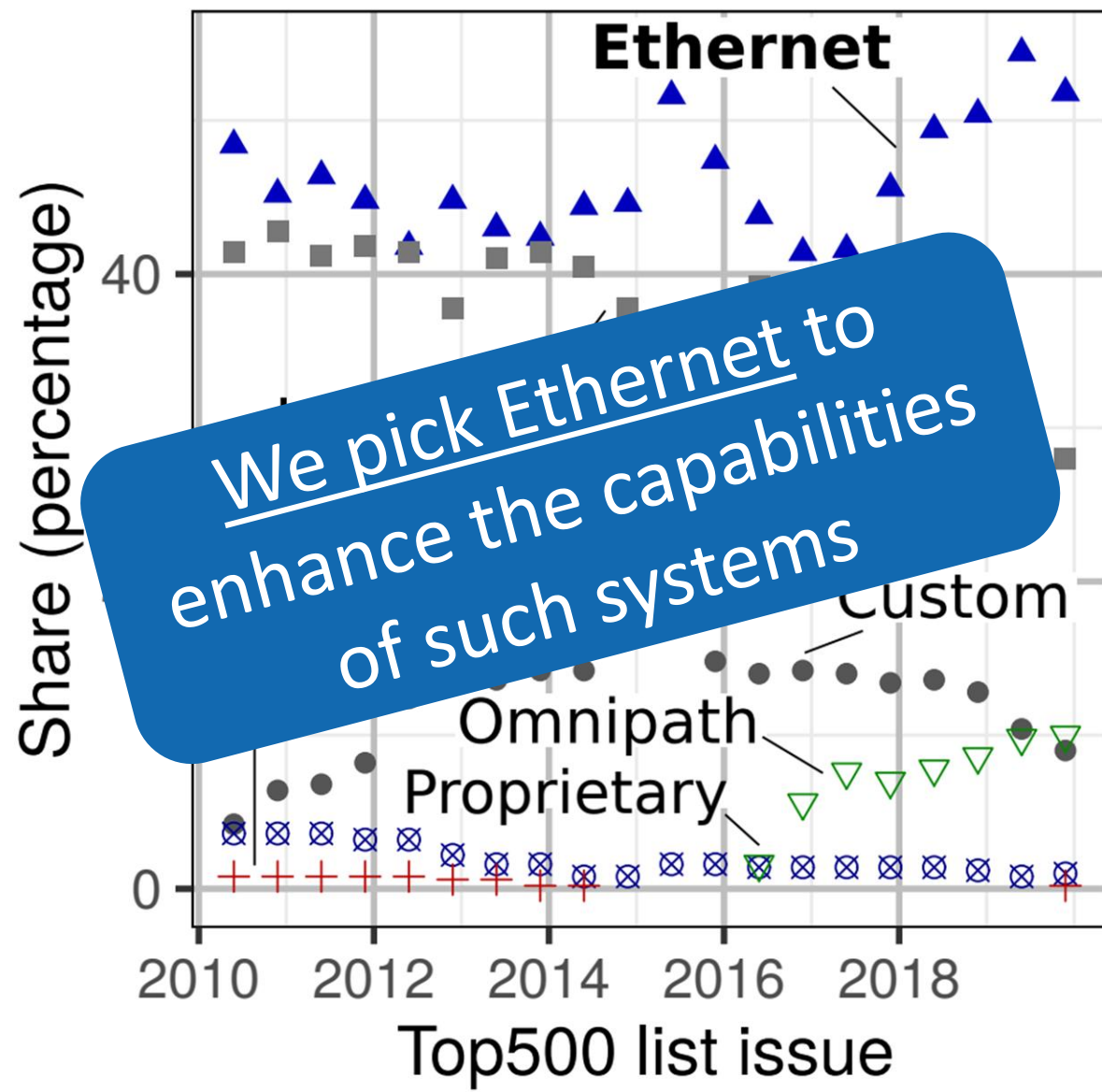


# ETHERNET & HPC

More than 50% of systems in Top500 use Ethernet

However, systems based on Ethernet are not as efficient as others:  $\approx$  (details in the paper)

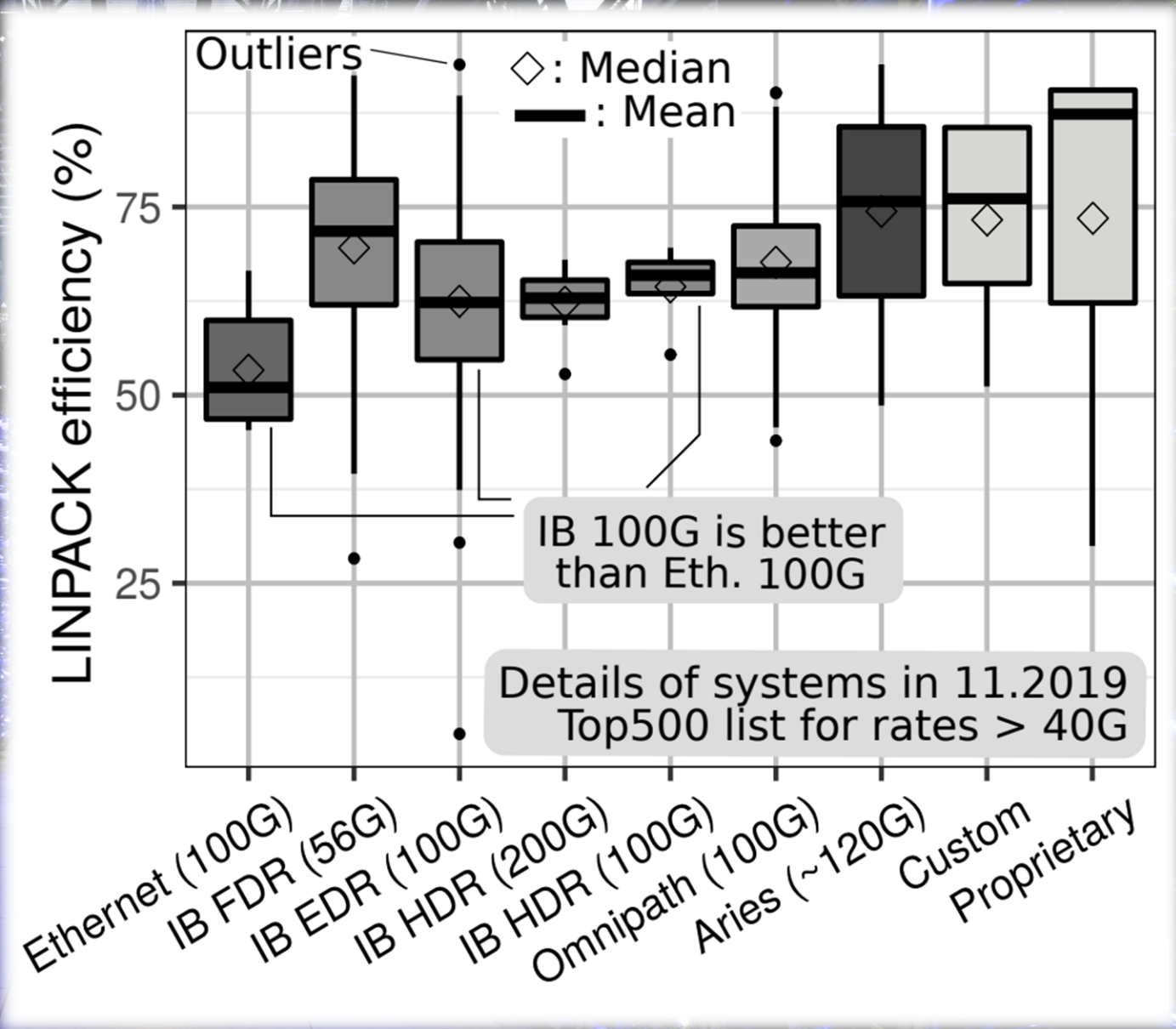
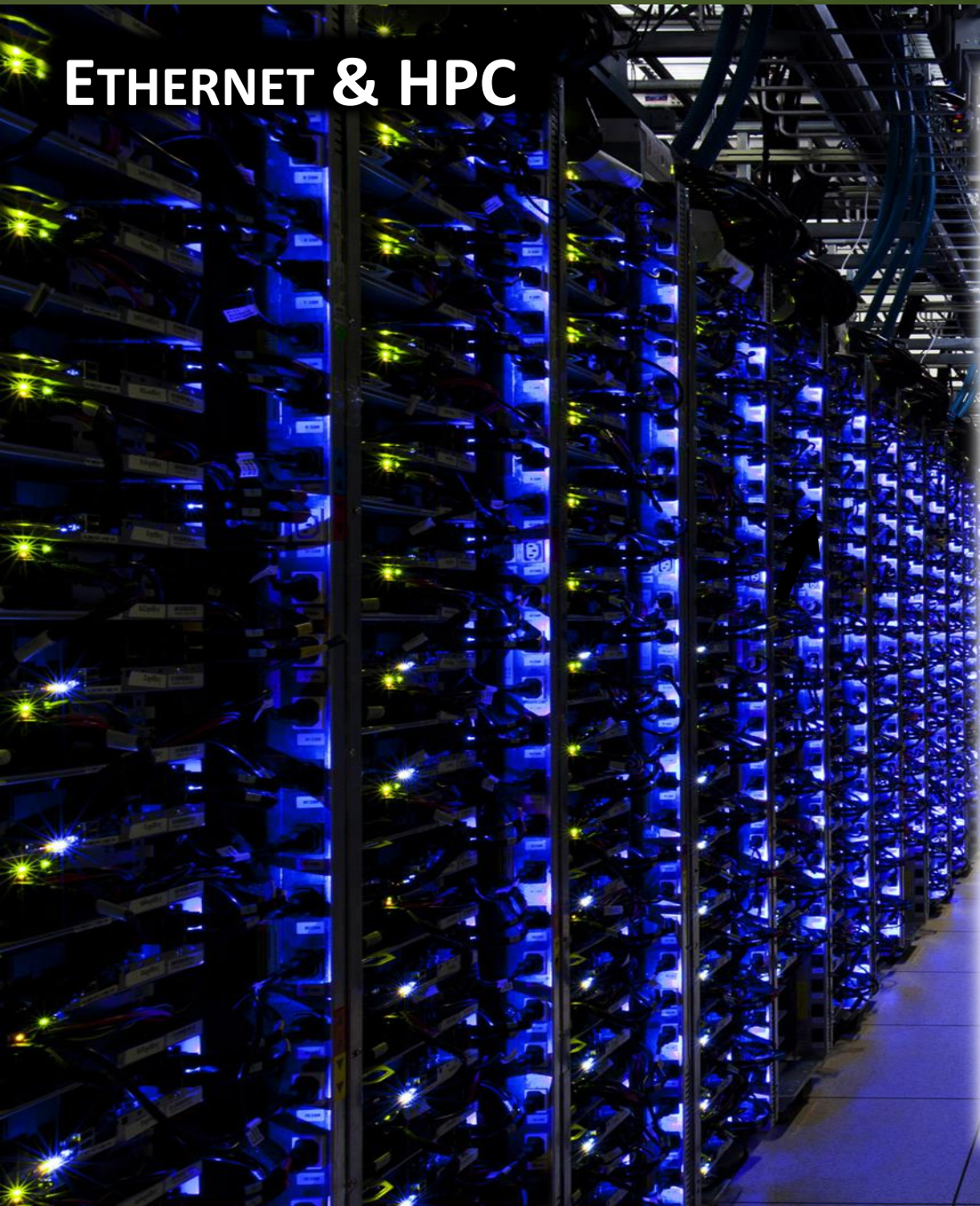
LINPACK Efficiency:  
 $\approx$ 51% for 100G Ethernet  
 $\approx$ 65% for 100G IB HDR  
 (details in the paper)



# ETHERNET & HPC

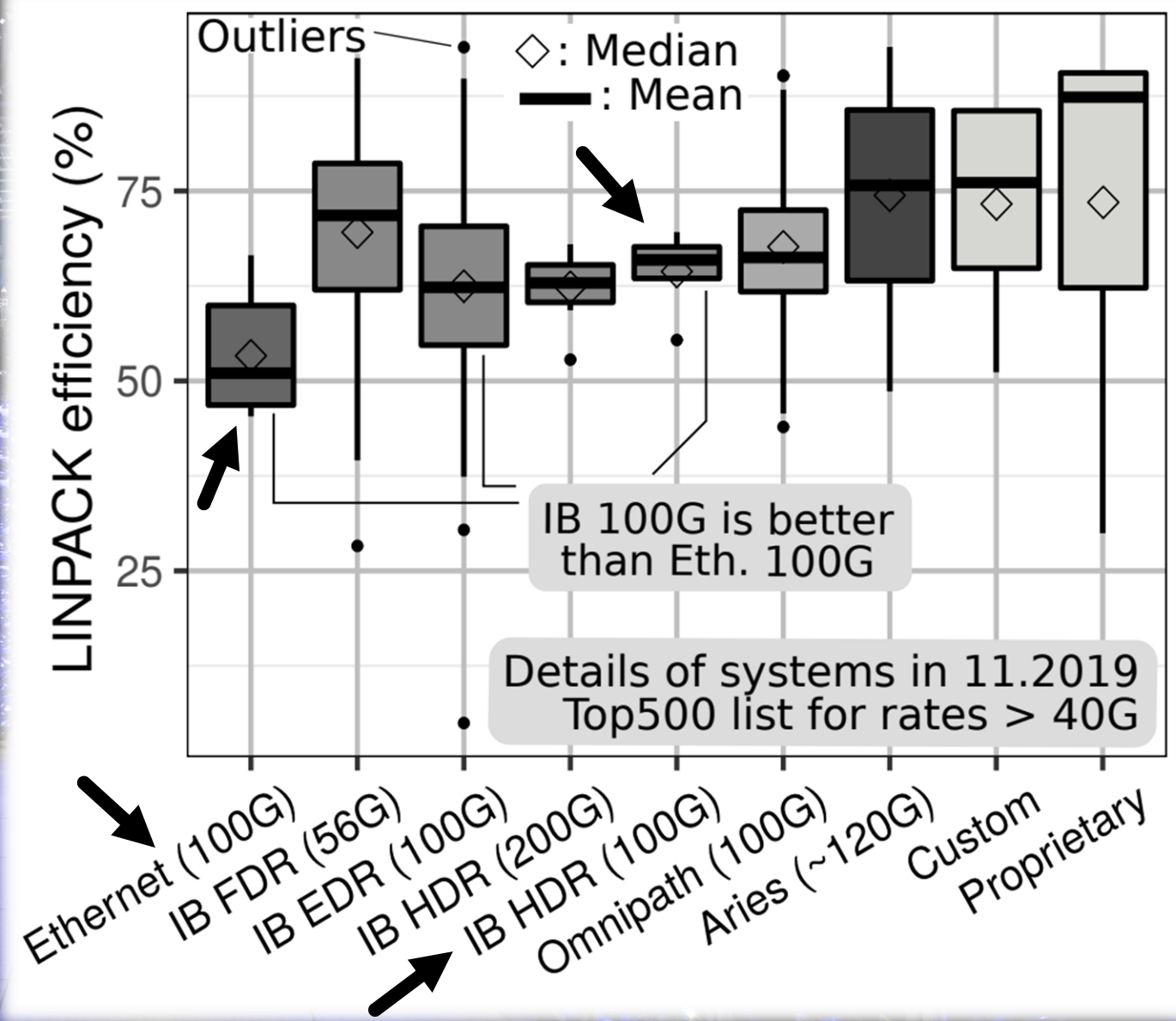


# ETHERNET & HPC



# ETHERNET & HPC

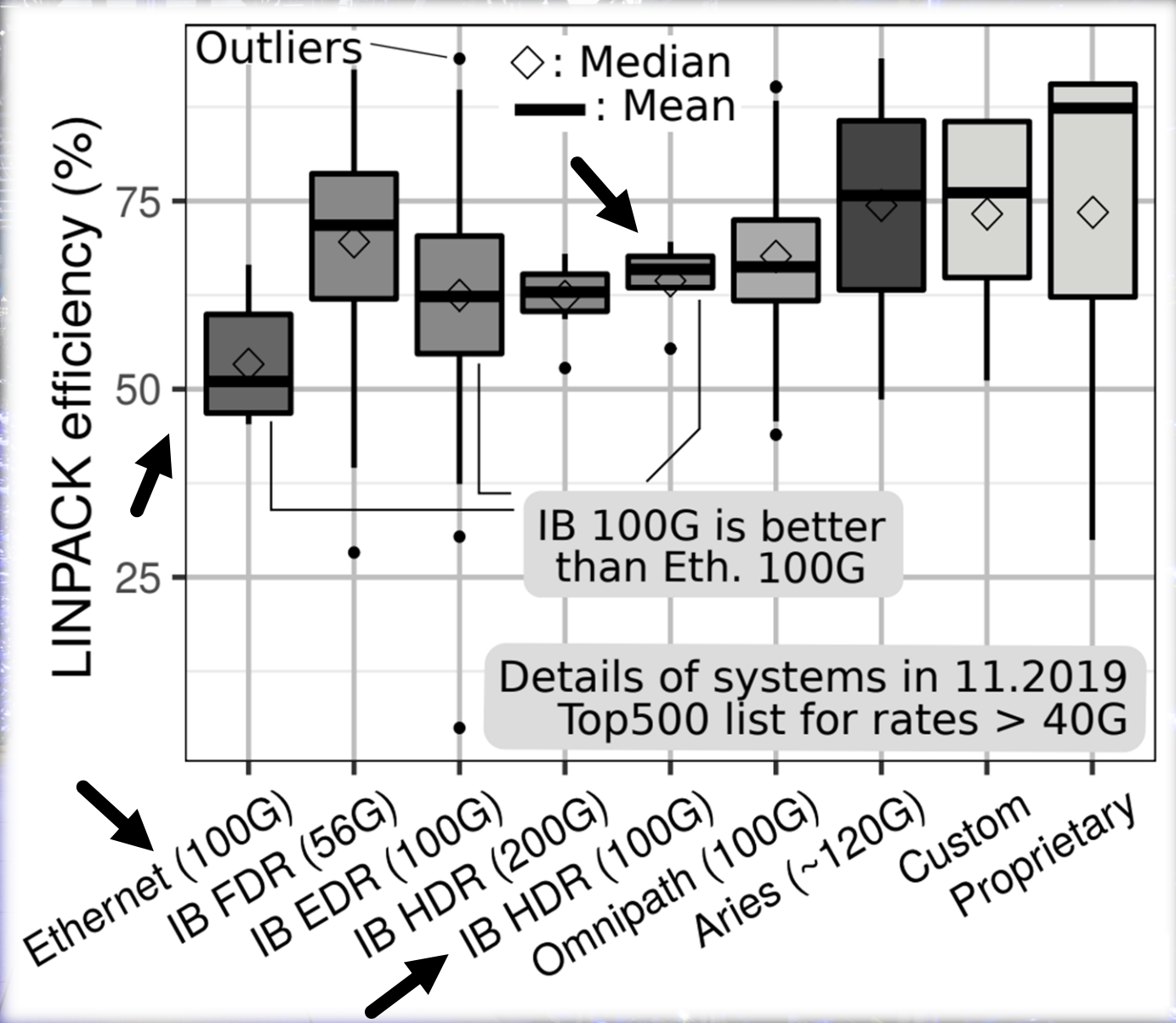
However, systems based on Ethernet are not as efficient as others



# ETHERNET & HPC

However, systems based on Ethernet are not as efficient as others

We pick Ethernet to enhance the capabilities of such systems



# TRANSPORT DESIGN

# TRANSPORT DESIGN

How to maximize performance of the transport layer?



## TRANSPORT DESIGN

How to maximize performance of the transport layer?



Adapt recent flow  
control for fat trees [1]

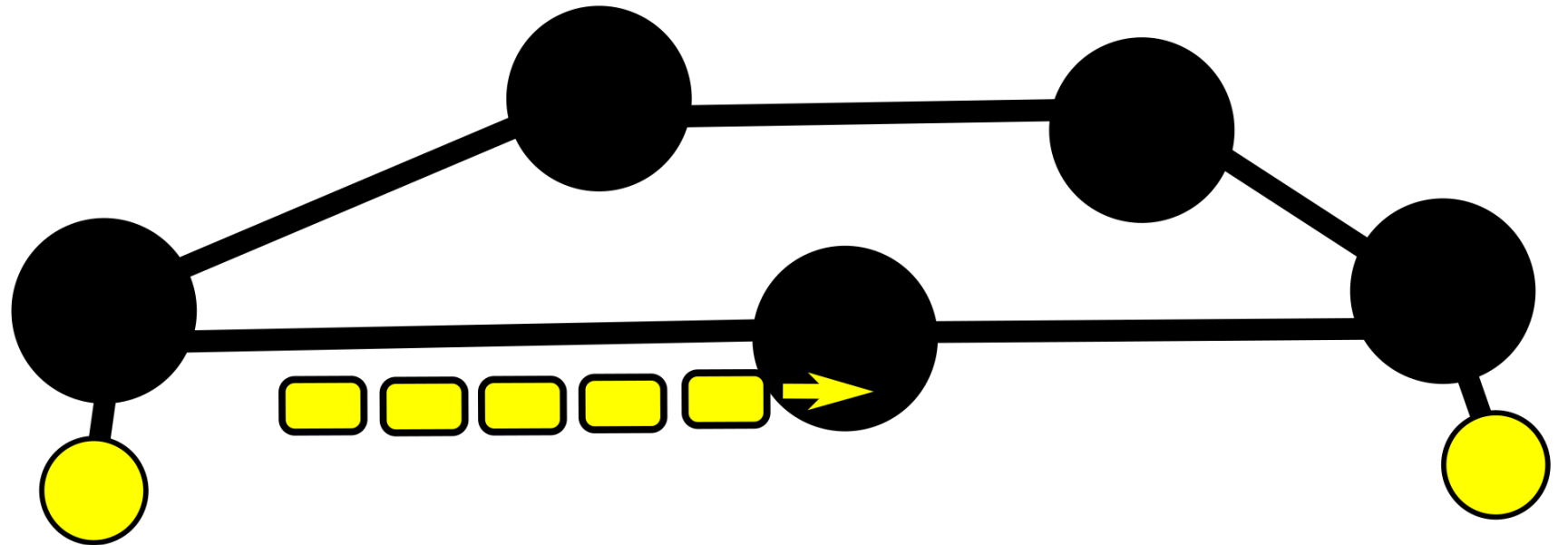


## TRANSPORT DESIGN

How to maximize performance of the transport layer?



Adapt recent flow control for fat trees [1]



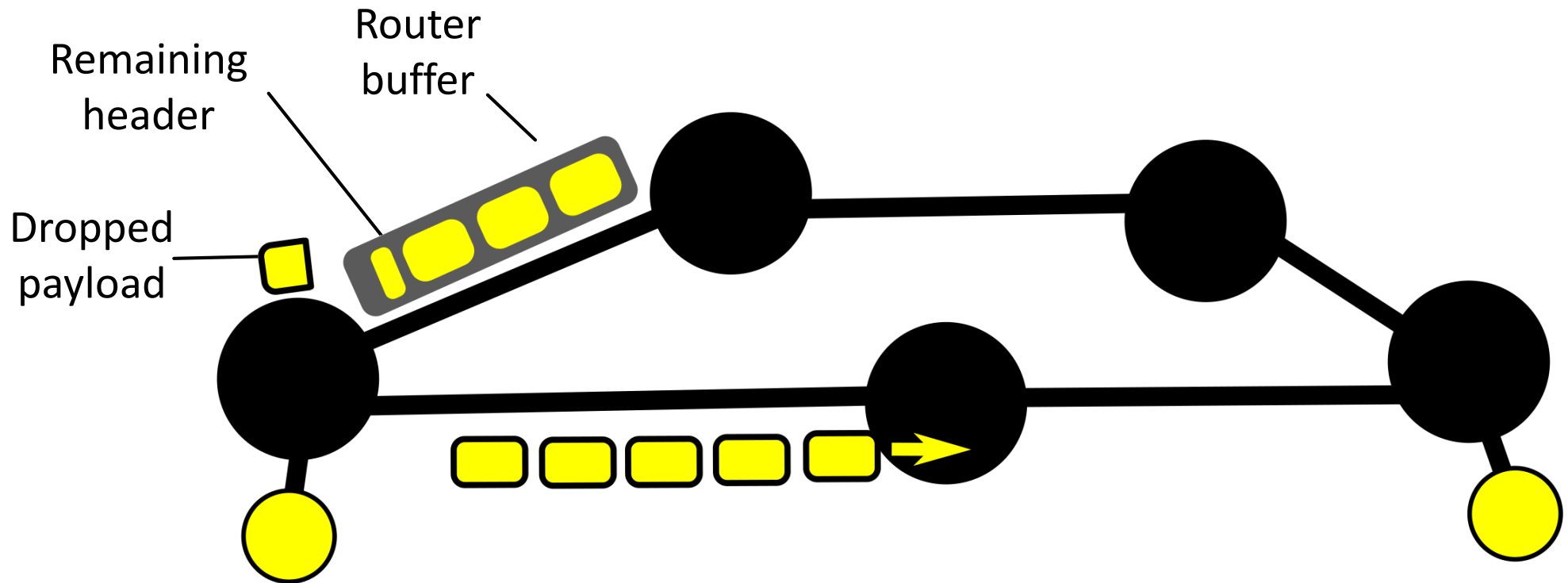
# TRANSPORT DESIGN

How to maximize performance of the transport layer?



Adapt recent flow control for fat trees [1]

Key design choice: Drop only payload if router buffers fill up



[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

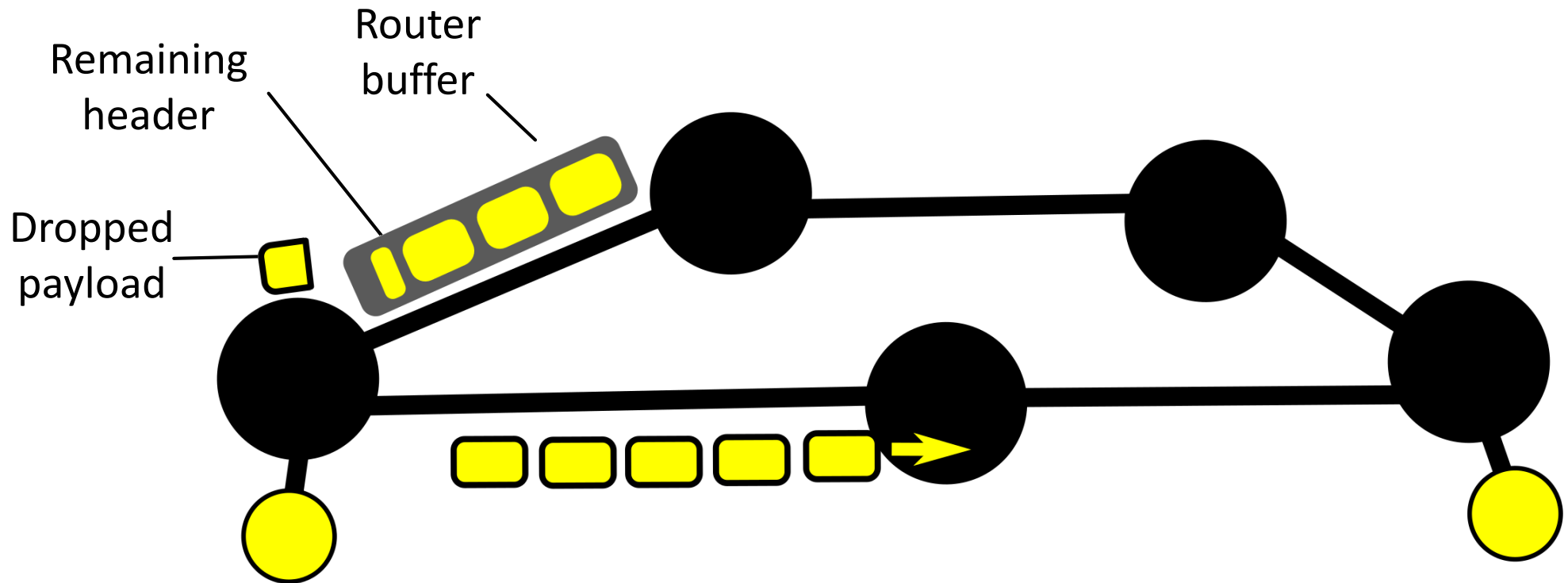
# TRANSPORT DESIGN

How to maximize performance of the transport layer?

Adapt recent flow control for fat trees [1]

Prioritize packets with dropped payload and retransmitted packets

Key design choice: Drop only payload if router buffers fill up



[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

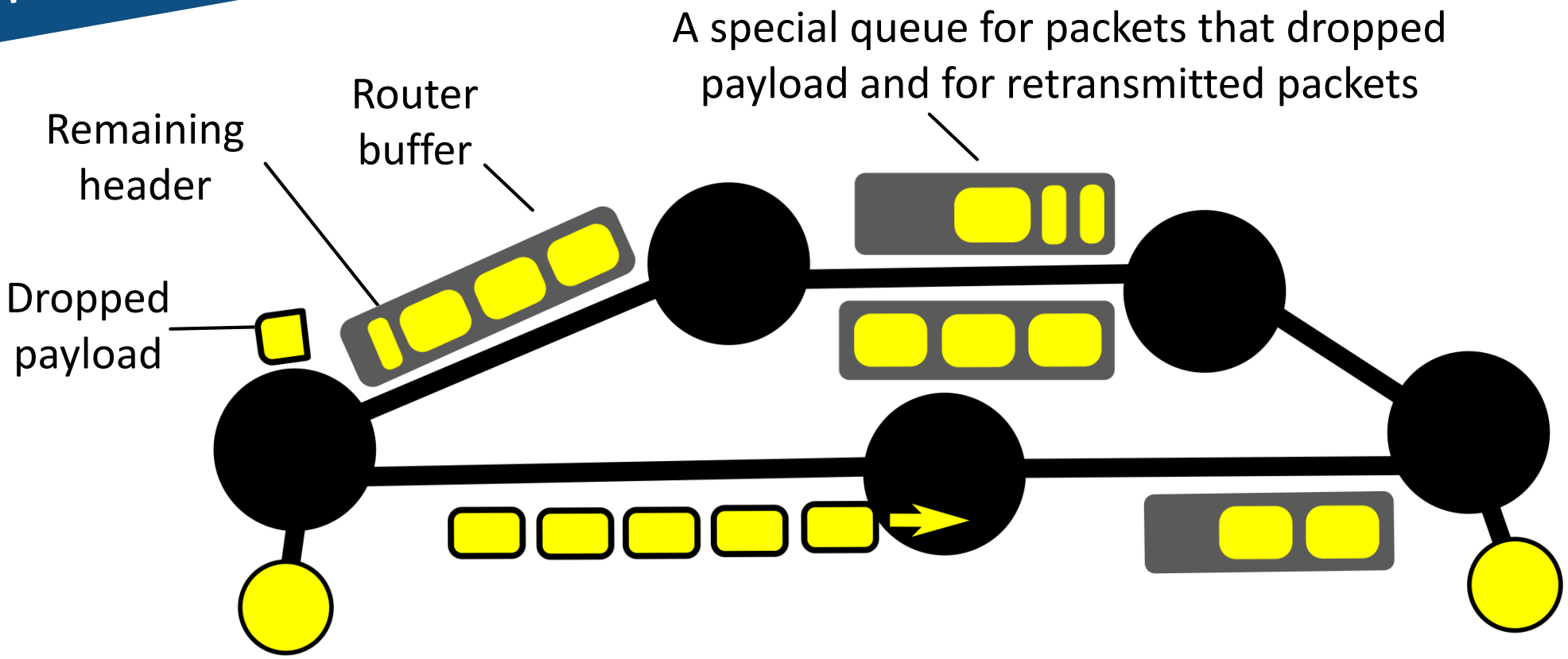
# TRANSPORT DESIGN

How to maximize performance of the transport layer?

Adapt recent flow control for fat trees [1]

Prioritize packets with dropped payload and retransmitted packets

Key design choice: Drop only payload if router buffers fill up



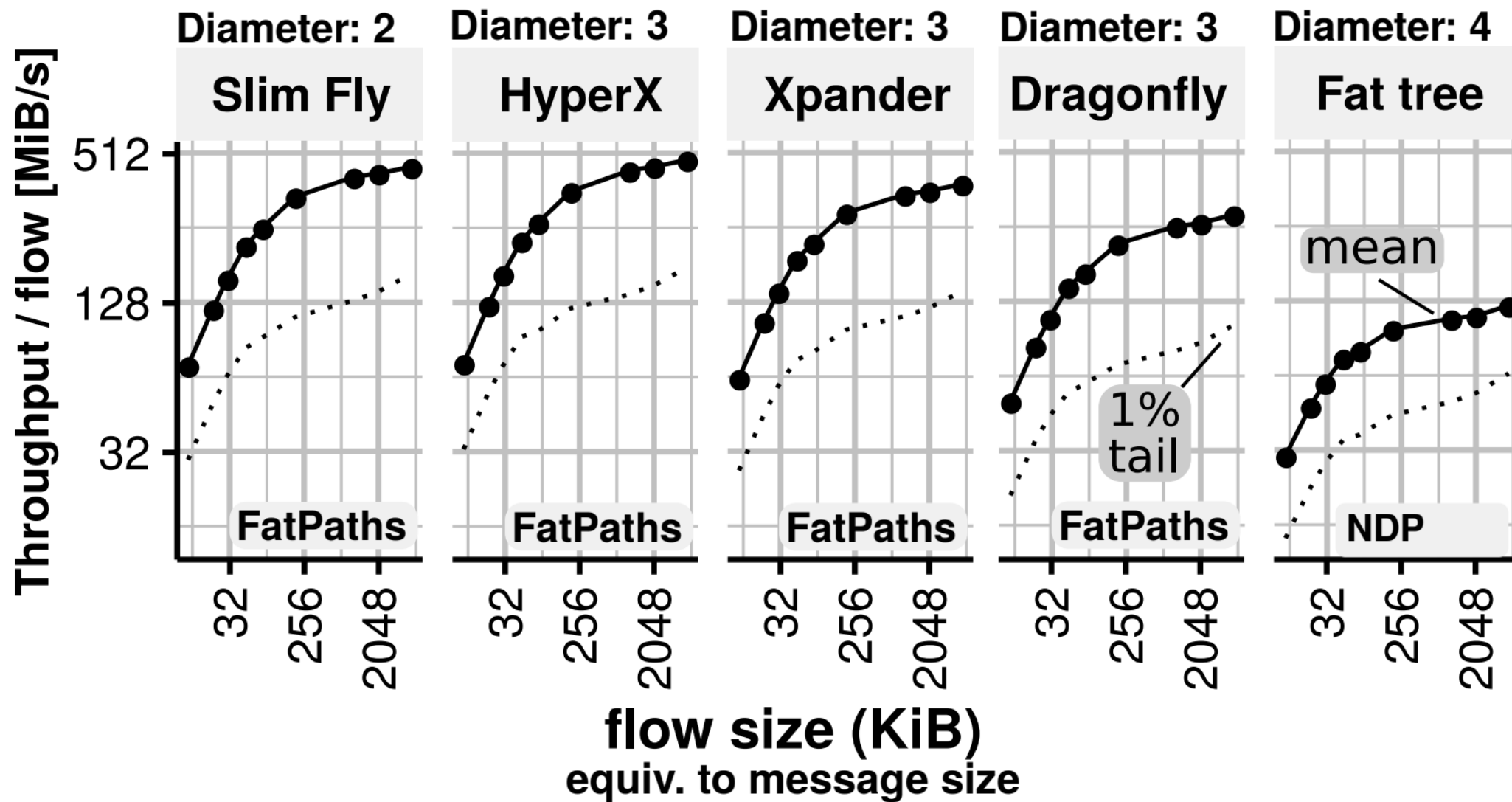
[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

## EVALUATION

$N \approx 10,000$ ; comparable cost; random uniform traffic; “bare Ethernet” setting

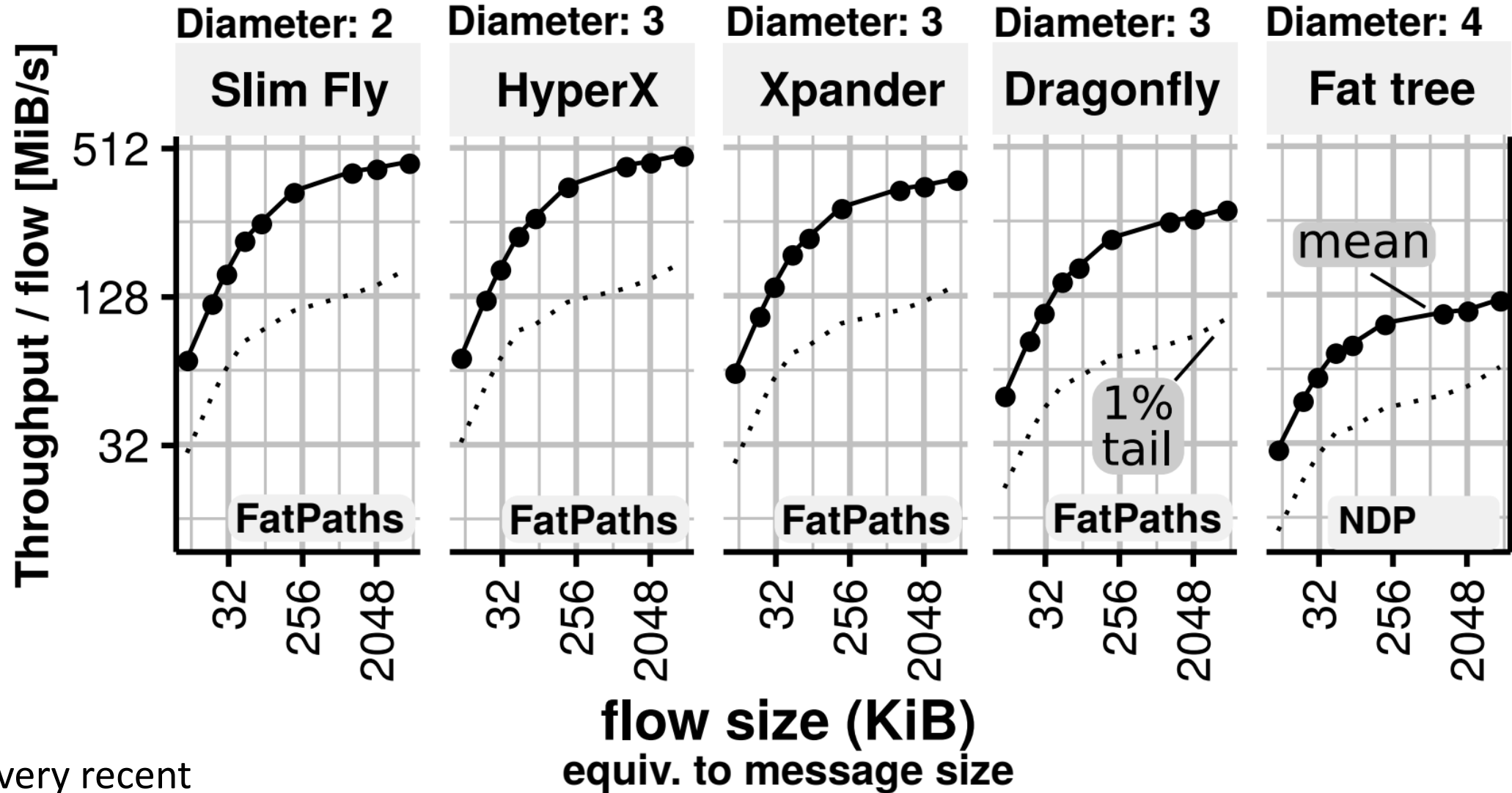
# EVALUATION

$N \approx 10,000$ ; comparable cost; random uniform traffic; “bare Ethernet” setting



# EVALUATION

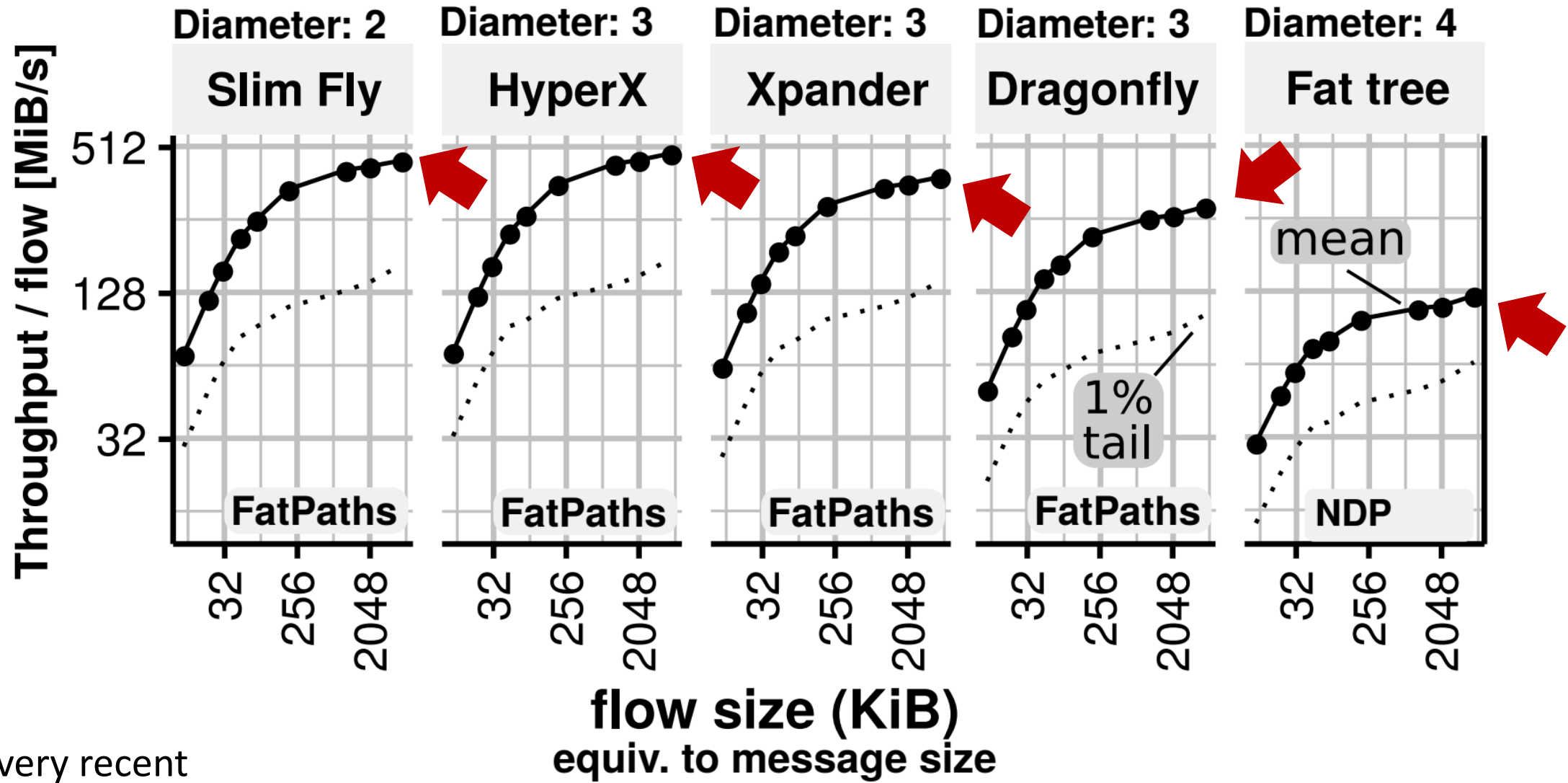
N ≈ 10,000; comparable cost; random uniform traffic; “bare Ethernet” setting



NDP: a very recent baseline for fat trees

# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “bare Ethernet” setting

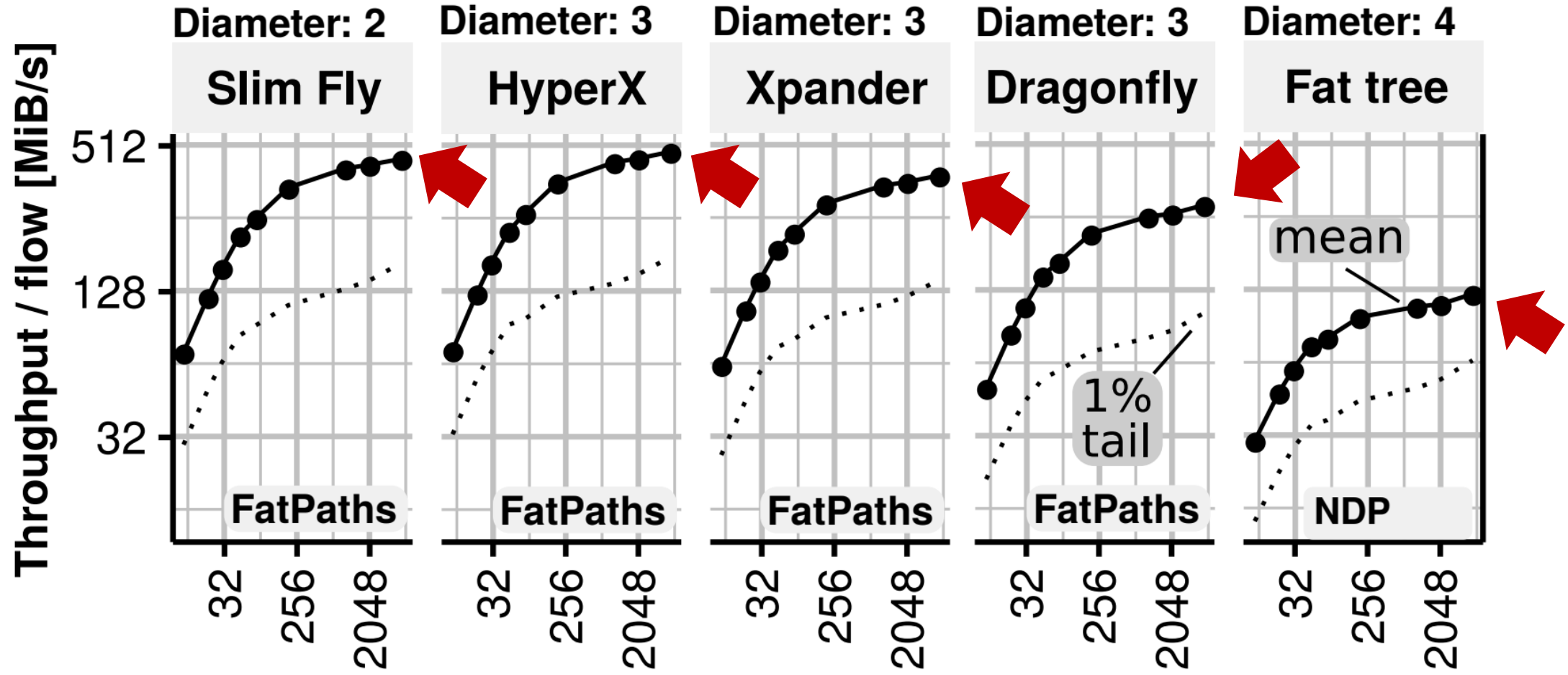


NDP: a very recent baseline for fat trees



# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “bare Ethernet” setting



flow size (KiB)  
equiv. to message size

Low-diameter topologies with FatPaths outperform fat trees

NDP: a very recent baseline for fat trees

# EVALUATION

$N \approx 10,000$ ; comparable cost; random uniform traffic; “full TCP” setting

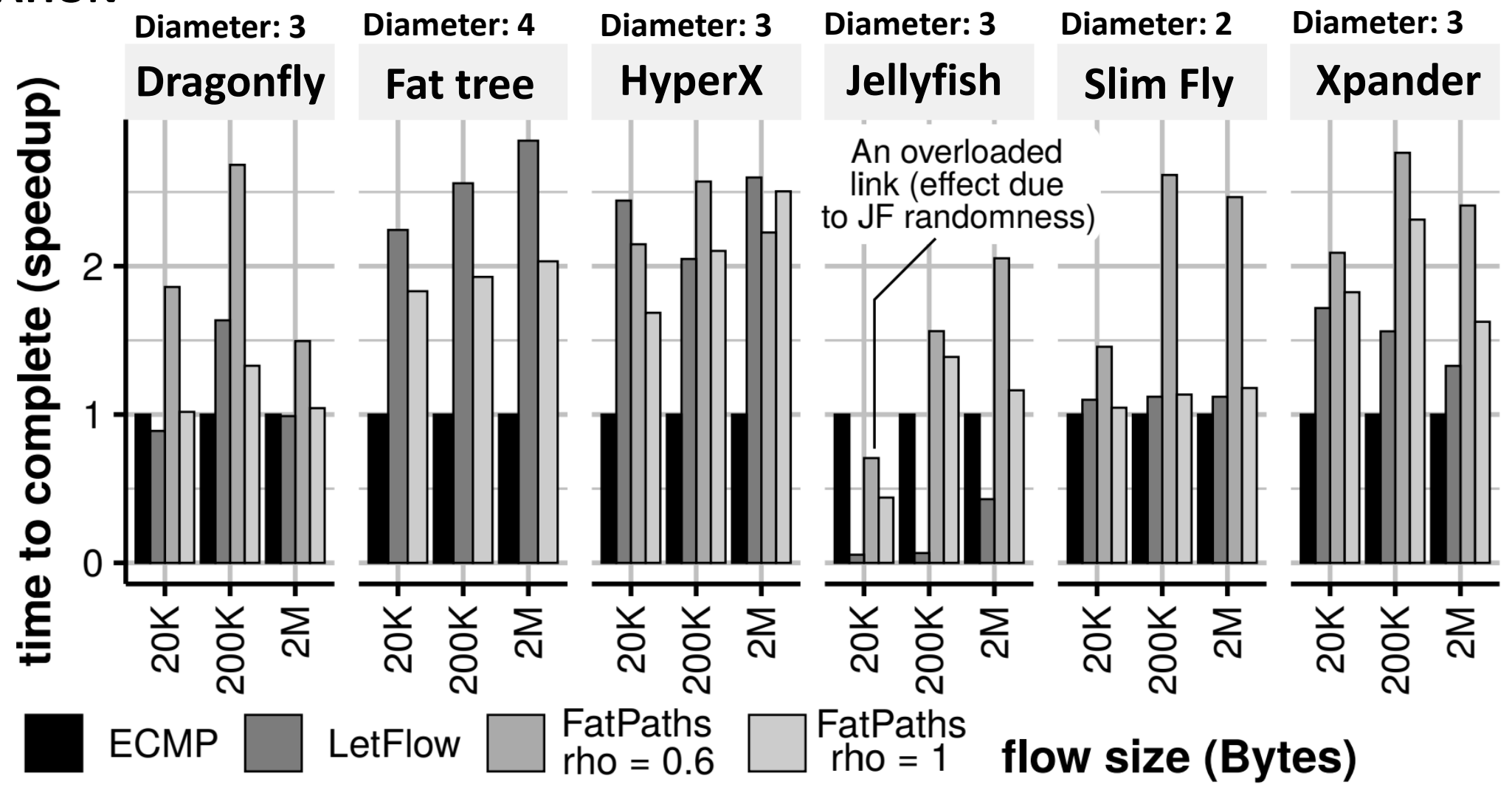
Diameter: 3	Diameter: 4	Diameter: 3	Diameter: 3	Diameter: 2	Diameter: 3
<b>Dragonfly</b>	<b>Fat tree</b>	<b>HyperX</b>	<b>Jellyfish</b>	<b>Slim Fly</b>	<b>Xpander</b>

**ECMP**: traditional static load balancing  
**LetFlow**: recent adaptive load balancing

**Speedups** are measured over **ECMP**  
**rho**: fraction of links kept in a layer

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

# EVALUATION

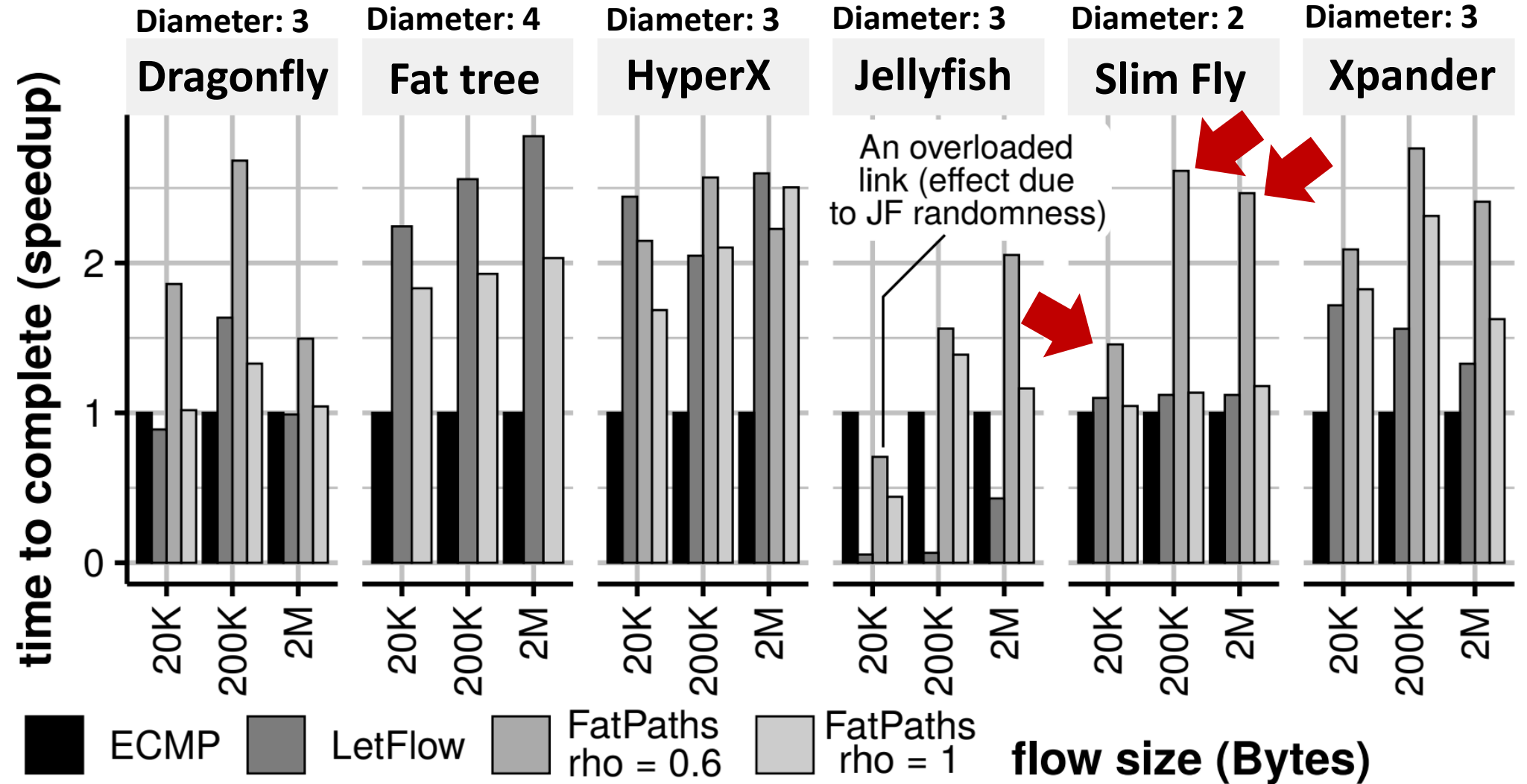


**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

# EVALUATION

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

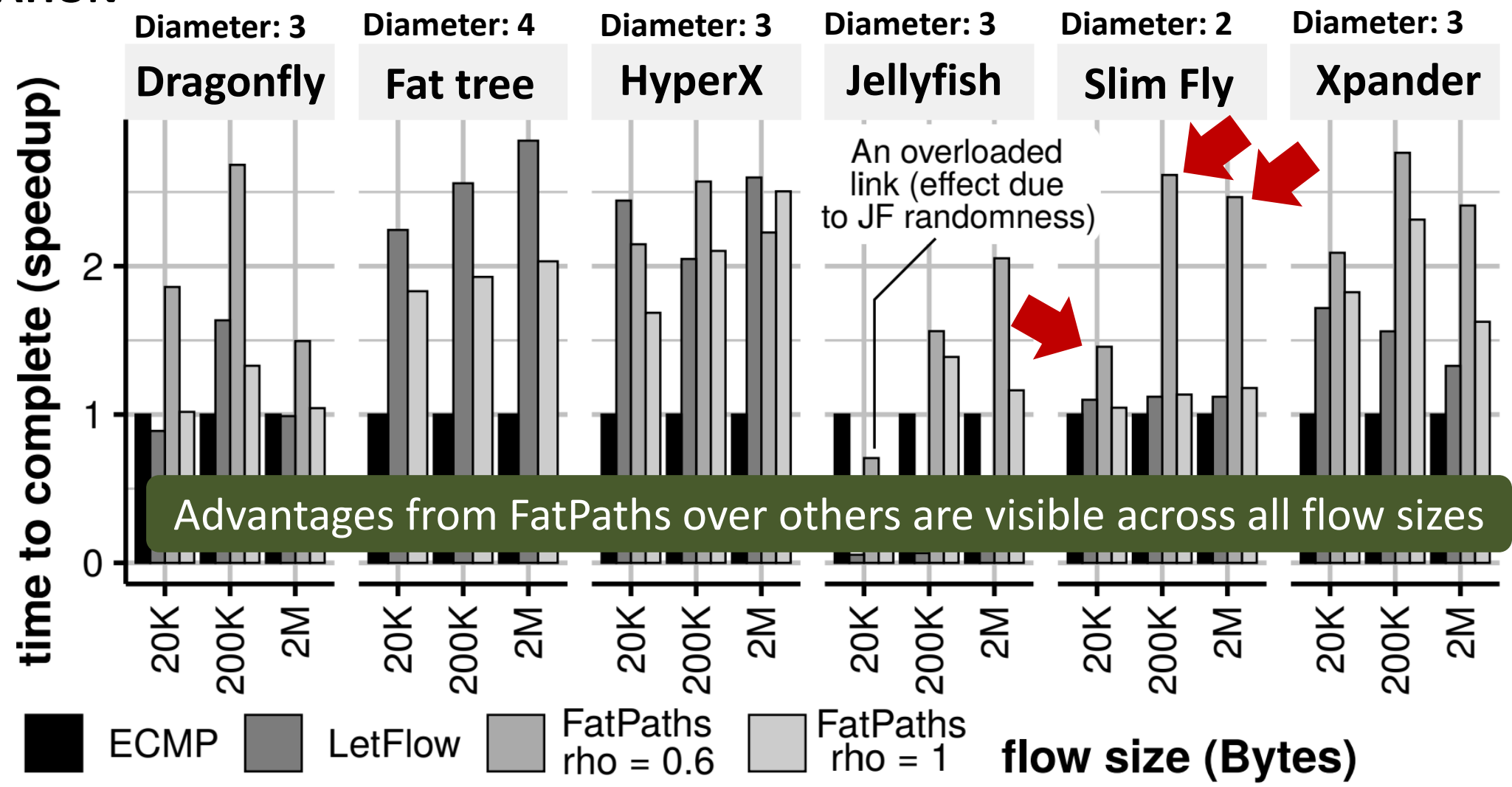


**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

N ≈ 10,000; comparable cost; random uniform traffic; “full TCP” setting

# EVALUATION



**ECMP:** traditional static load balancing  
**LetFlow:** recent adaptive load balancing

**Speedups are measured over ECMP**  
**rho:** fraction of links kept in a layer

# SURVEY OF ROUTING PROTOCOLS

# SURVEY OF ROUTING PROTOCOLS

**Focus: path  
diversity in routing**

# SURVEY OF ROUTING PROTOCOLS

**Focus: path diversity in routing**

Routing Scheme (Name, Abbreviation, Reference)	Stack Layer	Supported path diversity aspect						
		SP	NP	SM	MP	DP	ALB	AT
<b>(1) SIMPLE ROUTING PROTOCOLS (often used as building blocks):</b>								
Valiant load balancing (VLB) [17]	L2–L3	👍	👍	👍	👍	👍	👍	👍
Simple Spanning Tree (ST) [18]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple routing, e.g., OSPF [19]–[22]	L2, L3	👍	👍	👍	👍	👍	👍	👍
UGAL [23]	L2–L3	👍	👍	👍	👍	👍	👍	👍
ECMP [11], OMP [24], Pkt. Spraying [25]	L2, L3	👍	👍	👍	👍	👍	👍	👍
<b>(2) ROUTING ARCHITECTURES:</b>								
DCell [26]	L2–L3	👍	👍	👍	👍	👍	👍	👍
Monsoon [27]	L2, L3	👍	👍	👍	👍	👍	👍	👍
PortLand [9]	L2	👍	👍	👍	👍	👍	👍	👍
DRILL [28], LocalFlow [29], DRB [30]	L2	👍	👍	👍	👍	👍	👍	👍
VL2 [31]	L3	👍	👍	👍	👍	👍	👍	👍
Architecture by Al-Fares et al. [32]	L2–L3	👍	👍	👍	👍	👍	👍	👍
BCube [33]	L2–L3	👍	👍	👍	👍	👍	👍	👍
SEATTLE [34], others* [35]–[38]	L2	👍	👍	👍	👍	👍	👍	👍
VIRO [39]	L2–L3	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Ethernet on Air [40]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍 <sup>R</sup>	👍	👍	👍
PAST [41]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
MLAG, MC-LAG, others [42]	L2	👍	👍	👍	👍 <sup>R</sup>	👍	👍	👍
MOOSE [43]	L2	👍	👍	👍	👍	👍	👍	👍
MPA [44]	L3	👍	👍	👍	👍	👍	👍	👍
AMP [45]	L3	👍	👍	👍	👍	👍	👍	👍
MSTP [46], GOE [47], Viking [48]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
SPB [49], TRILL [50], Shadow MACs [51]	L2	👍	👍 <sup>R</sup>	👍	👍	👍	👍	👍
SPAIN [52]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍
<b>(3) Schemes for exposing/encoding paths (can be combined with FatPaths):</b>								
XPath [53]	L3	👍	👍	👍	👍	👍	👍	👍
Source routing for flexible DC fabric [54]	L3	👍	👍 <sup>R</sup>	👍 <sup>R</sup>	👍	👍	👍	👍 <sup>†</sup>
<b>(3) FatPaths [This work]</b>								
	L2–L3	👍	👍	👍	👍	👍	👍	👍



# SURVEY OF ROUTING PROTOCOLS

**Focus: path diversity in routing**

**(1) Taxonomy of „path diversity support“ in routing**

Routing Scheme (Name, Abbreviation, Reference)	Stack Layer	Supported path diversity aspect						
		SP	NP	SM	MP	DP	ALB	AT
<b>(1) SIMPLE ROUTING PROTOCOLS (often used as building blocks):</b>								
Valiant load balancing (VLB) [17]	L2–L3	👍	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple Spanning Tree (ST) [18]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple routing, e.g., OSPF [19]–[22]	L2, L3	👍	👍	👍	👍	👍	👍	👍
UGAL [23]	L2–L3	👍	👍	👍	👍	👍	👍	👍
ECMP [11], OMP [24], Pkt. Spraying [25]	L2, L3	👍	👍	👍	👍	👍	👍	👍
<b>(2) ROUTING ARCHITECTURES:</b>								
DCell [26]	L2–L3	👍	👍	👍	👍	👍	👍	👍
Monsoon [27]	L2, L3	👍	👍	👍	👍	👍	👍	👍
PortLand [9]	L2	👍	👍	👍	👍	👍	👍	👍
DRILL [28], LocalFlow [29], DRB [30]	L2	👍	👍	👍	👍	👍	👍	👍
VL2 [31]	L3	👍	👍	👍	👍	👍	👍	👍
Architecture by Al-Fares et al. [32]	L2–L3	👍	👍	👍	👍	👍	👍	👍
BCube [33]	L2–L3	👍	👍	👍	👍	👍	👍	👍
SEATTLE [34], others* [35]–[38]	L2	👍	👍	👍	👍	👍	👍	👍
VIRO [39]	L2–L3	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Ethernet on Air [40]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍 <sup>R</sup>	👍	👍	👍
PAST [41]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
MLAG, MC-LAG, others [42]	L2	👍	👍	👍	👍 <sup>R</sup>	👍	👍	👍
MOOSE [43]	L2	👍	👍	👍	👍	👍	👍	👍
MPA [44]	L3	👍	👍	👍	👍	👍	👍	👍
AMP [45]	L3	👍	👍	👍	👍	👍	👍	👍
MSTP [46], GOE [47], Viking [48]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
SPB [49], TRILL [50], Shadow MACs [51]	L2	👍	👍 <sup>R</sup>	👍	👍	👍	👍	👍
SPAIN [52]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍
<b>(3) Schemes for exposing/encoding paths (can be combined with FatPaths):</b>								
XPath [53]	L3	👍	👍	👍	👍	👍	👍	👍
Source routing for flexible DC fabric [54]	L3	👍	👍 <sup>R</sup>	👍 <sup>R</sup>	👍	👍	👍	👍 <sup>†</sup>
<b>(3) FatPaths [This work]</b>								
	L2–L3	👍	👍	👍	👍	👍	👍	👍

# SURVEY OF ROUTING PROTOCOLS

**Focus: path diversity in routing**

**(1) Taxonomy of „path diversity support” in routing**

**(2) Analysis of routing protocols**

Routing Scheme (Name, Abbreviation, Reference)	Stack Layer	Supported path diversity aspect						
		SP	NP	SM	MP	DP	ALB	AT
<b>(1) SIMPLE ROUTING PROTOCOLS (often used as building blocks):</b>								
Valiant load balancing (VLB) [17]	L2–L3	👍	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple Spanning Tree (ST) [18]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple routing, e.g., OSPF [19]–[22]	L2, L3	👍	👍	👍	👍	👍	👍	👍
UGAL [23]	L2–L3	👍	👍	👍	👍	👍	👍	👍
ECMP [11], OMP [24], Pkt. Spraying [25]	L2, L3	👍	👍	👍	👍	👍	👍	👍
<b>(2) ROUTING ARCHITECTURES:</b>								
DCell [26]	L2–L3	👍	👍	👍	👍	👍	👍	👍
Monsoon [27]	L2, L3	👍	👍	👍	👍	👍	👍	👍
PortLand [9]	L2	👍	👍	👍	👍	👍	👍	👍
DRILL [28], LocalFlow [29], DRB [30]	L2	👍	👍	👍	👍	👍	👍	👍
VL2 [31]	L3	👍	👍	👍	👍	👍	👍	👍
Architecture by Al-Fares et al. [32]	L2–L3	👍	👍	👍	👍	👍	👍	👍
BCube [33]	L2–L3	👍	👍	👍	👍	👍	👍	👍
SEATTLE [34], others* [35]–[38]	L2	👍	👍	👍	👍	👍	👍	👍
VIRO [39]	L2–L3	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Ethernet on Air [40]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍 <sup>R</sup>	👍	👍	👍
PAST [41]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
MLAG, MC-LAG, others [42]	L2	👍	👍	👍	👍 <sup>R</sup>	👍	👍	👍
MOOSE [43]	L2	👍	👍	👍	👍	👍	👍	👍
MPA [44]	L3	👍	👍	👍	👍	👍	👍	👍
AMP [45]	L3	👍	👍	👍	👍	👍	👍	👍
MSTP [46], GOE [47], Viking [48]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
SPB [49], TRILL [50], Shadow MACs [51]	L2	👍	👍 <sup>R</sup>	👍	👍	👍	👍	👍
SPAIN [52]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍
<b>(3) Schemes for exposing/encoding paths (can be combined with FatPaths):</b>								
XPath [53]	L3	👍	👍	👍	👍	👍	👍	👍
Source routing for flexible DC fabric [54]	L3	👍	👍 <sup>R</sup>	👍 <sup>R</sup>	👍	👍	👍	👍 <sup>†</sup>
<b>(3) FatPaths [This work]</b>	L2–L3	👍	👍	👍	👍	👍	👍	👍

# SURVEY OF ROUTING PROTOCOLS

**Focus: path diversity in routing**

**(1) Taxonomy of „path diversity support” in routing**

**(2) Analysis of routing protocols**

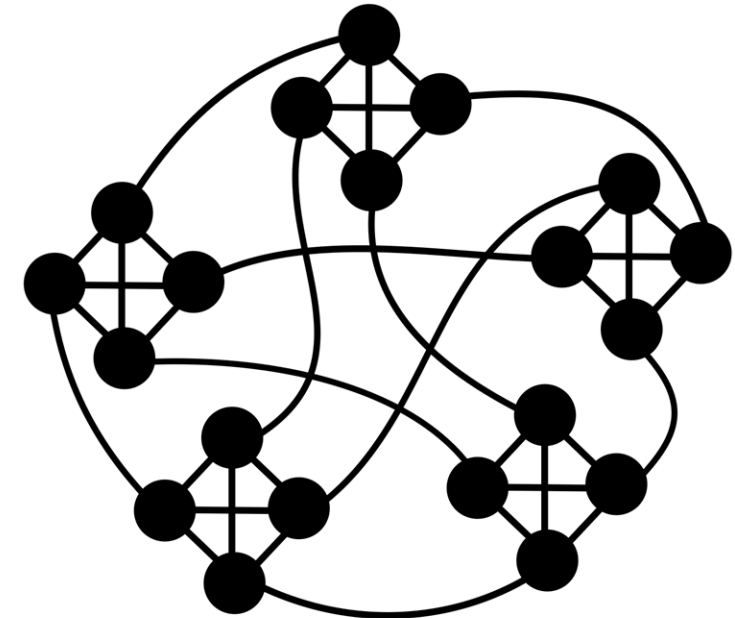
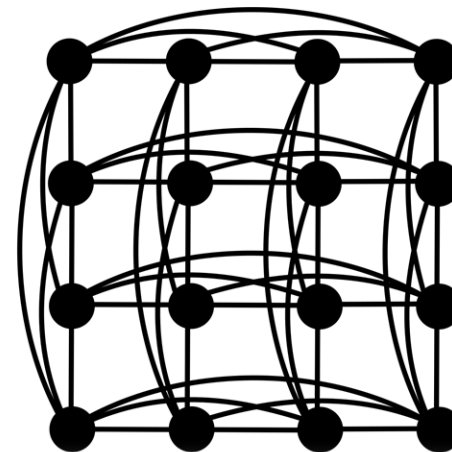
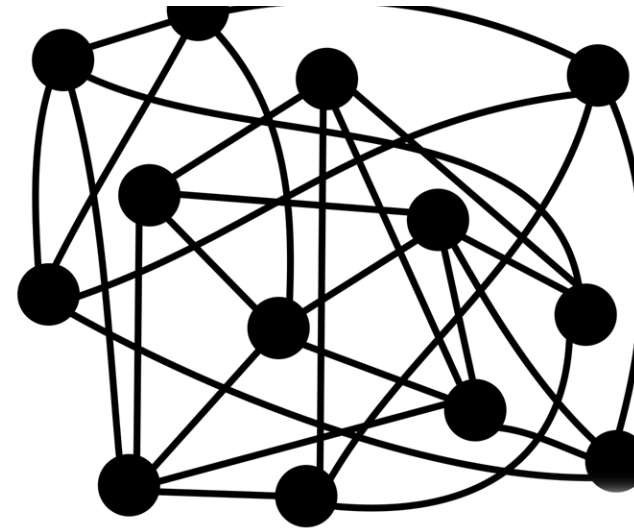
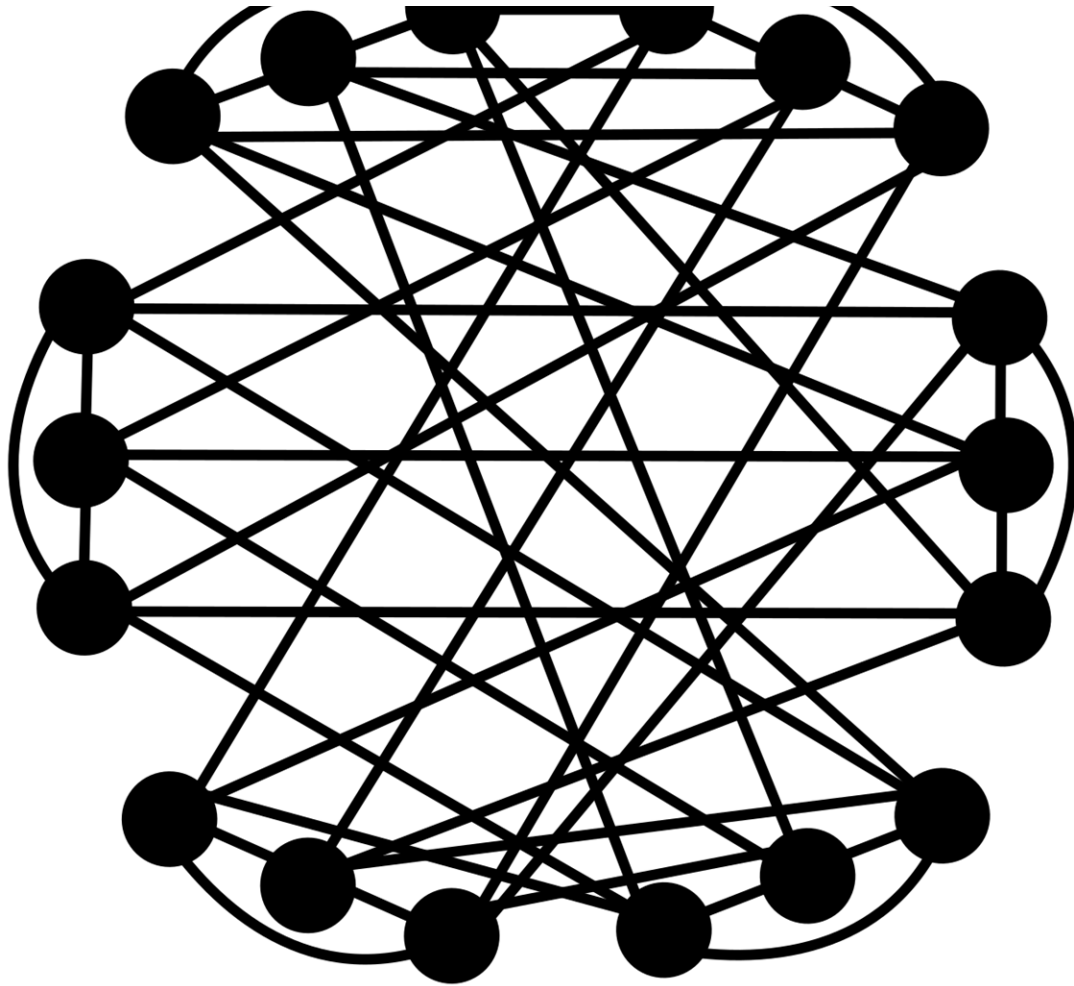
1

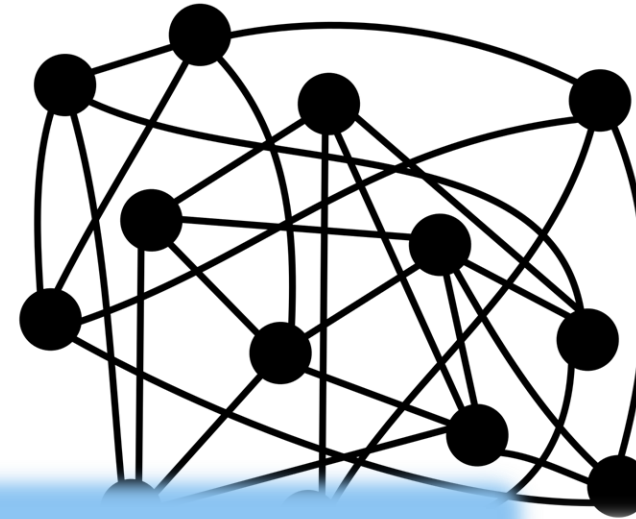
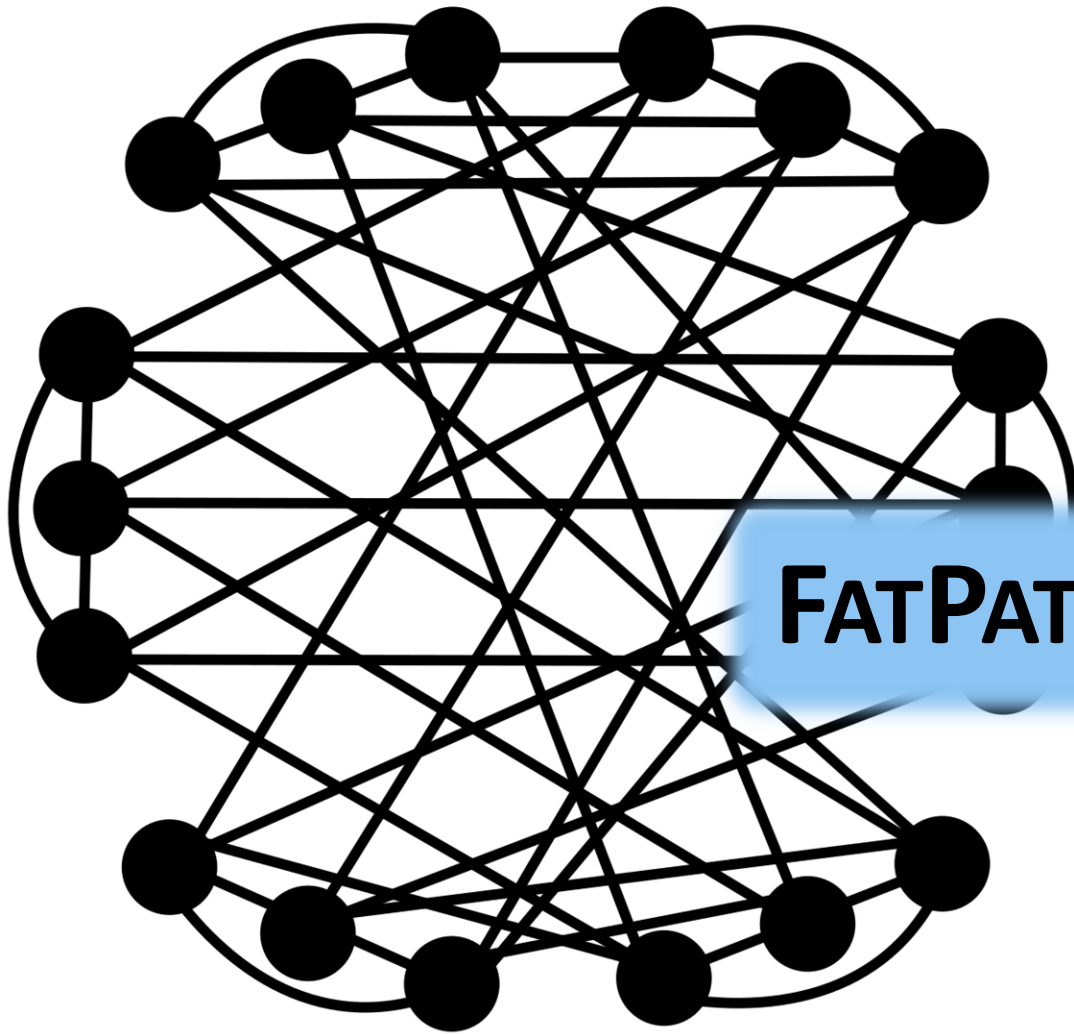
## High-Performance Routing with Multipathing and Path Diversity in Supercomputers and Data Centers

Maciej Besta<sup>1</sup>, Jens Domke<sup>2</sup>, Marcel Schneider<sup>1</sup>, Marek Konieczny<sup>3</sup>,  
 Salvatore Di Girolamo<sup>1</sup>, Timo Schneider<sup>1</sup>, Ankit Singla<sup>1</sup>, Torsten Hoeffler<sup>1</sup>

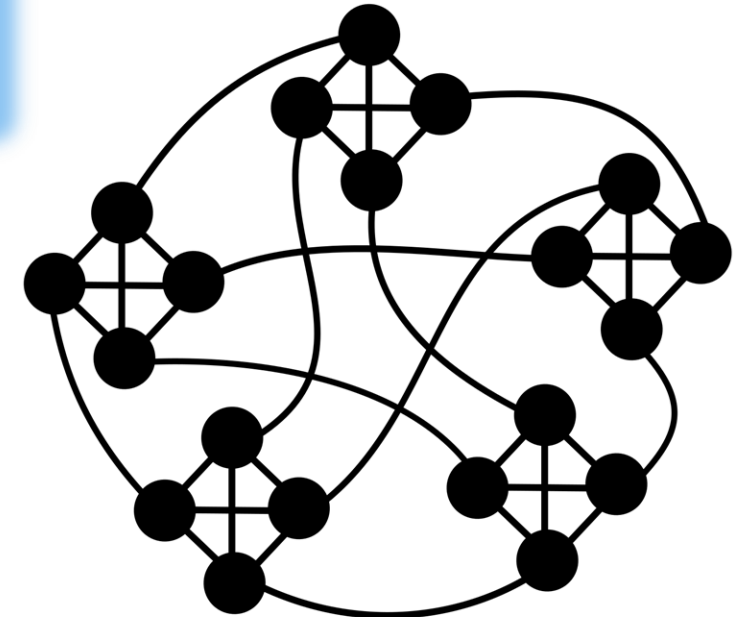
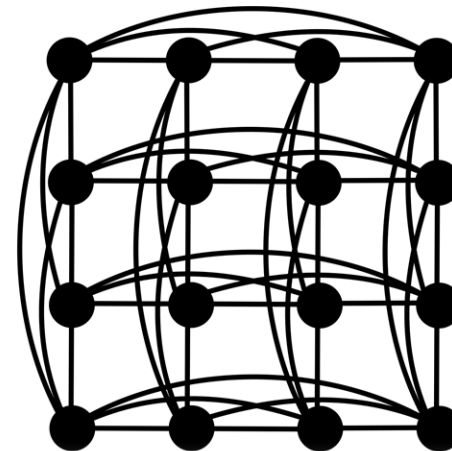
<sup>1</sup>Department of Computer Science, ETH Zurich; <sup>2</sup>RIKEN Center for Computational Science (R-CCS)  
<sup>3</sup>Department of Computer Science, Electronics and Telecommunications; AGH-UST

Routing Scheme (Name, Abbreviation, Reference)	Stack Layer	Supported path diversity aspect						
		SP	NP	SM	MP	DP	ALB	AT
<b>(1) SIMPLE ROUTING PROTOCOLS (often used as building blocks):</b>								
Valiant load balancing (VLB) [17]	L2–L3	👍	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple Spanning Tree (ST) [18]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Simple routing, e.g., OSPF [19]–[22]	L2, L3	👍	👍	👍	👍	👍	👍	👍
UGAL [23]	L2–L3	👍	👍	👍	👍	👍	👍	👍
ECMP [11], OMP [24], Pkt. Spraying [25]	L2, L3	👍	👍	👍	👍	👍	👍	👍
<b>(2) ROUTING ARCHITECTURES:</b>								
DCell [26]	L2–L3	👍	👍	👍	👍	👍	👍	👍
Monsoon [27]	L2, L3	👍	👍	👍	👍	👍	👍	👍
PortLand [9]	L2	👍	👍	👍	👍	👍	👍	👍
DRILL [28], LocalFlow [29], DRB [30]	L2	👍	👍	👍	👍	👍	👍	👍
VL2 [31]	L3	👍	👍	👍	👍	👍	👍	👍
Architecture by Al-Fares et al. [32]	L2–L3	👍	👍	👍	👍	👍	👍	👍
BCube [33]	L2–L3	👍	👍	👍	👍	👍	👍	👍
SEATTLE [34], others* [35]–[38]	L2	👍	👍	👍	👍	👍	👍	👍
VIRO [39]	L2–L3	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
Ethernet on Air [40]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍 <sup>R</sup>	👍	👍	👍
PAST [41]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
MLAG, MC-LAG, others [42]	L2	👍	👍	👍	👍 <sup>R</sup>	👍	👍	👍
MOOSE [43]	L2	👍	👍	👍	👍	👍	👍	👍
MPA [44]	L3	👍	👍	👍	👍	👍	👍	👍
AMP [45]	L3	👍	👍	👍	👍	👍	👍	👍
MSTP [46], GOE [47], Viking [48]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍	👍
SPB [49], TRILL [50], Shadow MACs [51]	L2	👍	👍 <sup>R</sup>	👍	👍	👍	👍	👍
SPAIN [52]	L2	👍 <sup>S</sup>	👍 <sup>S</sup>	👍 <sup>S</sup>	👍	👍	👍	👍
<b>(3) Schemes for exposing/encoding paths (can be combined with FatPaths):</b>								
XPath [53]	L3	👍	👍	👍	👍	👍	👍	👍
Source routing for flexible DC fabric [54]	L3	👍	👍 <sup>R</sup>	👍 <sup>R</sup>	👍	👍	👍	👍 <sup>†</sup>
<b>(3) FatPaths [This work]</b>								
	L2–L3	👍	👍	👍	👍	👍	👍	👍





## FATPATHS OVERVIEW



## A ROUTING SCHEME FOCUSING ON LOW-DIAMETER TOPOLOGIES

LOW-DIAMETER NETWORK TOPOLOGIES

How to route modern low diameter topologies?

LAYERED ROUTING How to encode & use diversity of shortest and non-minimal paths?

(a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

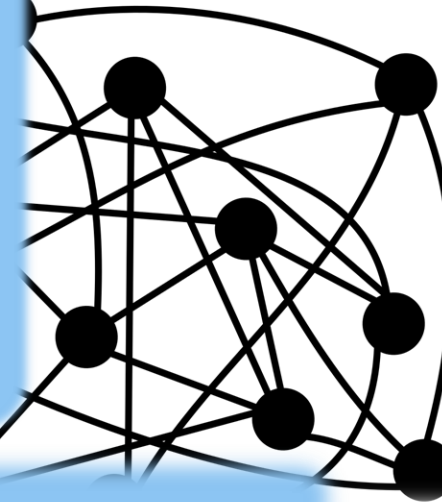
(b) Divide one flow into subflows and send subflows across different layers

(c) Route minimally in each layer

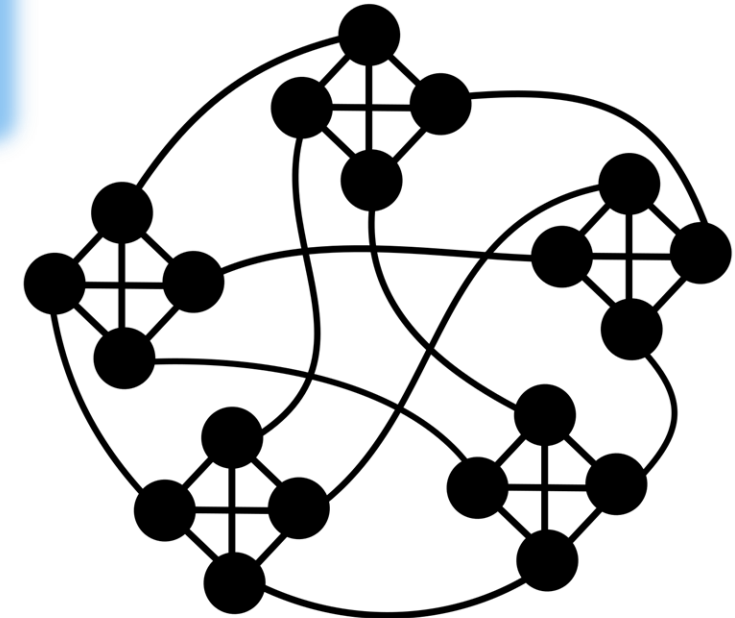
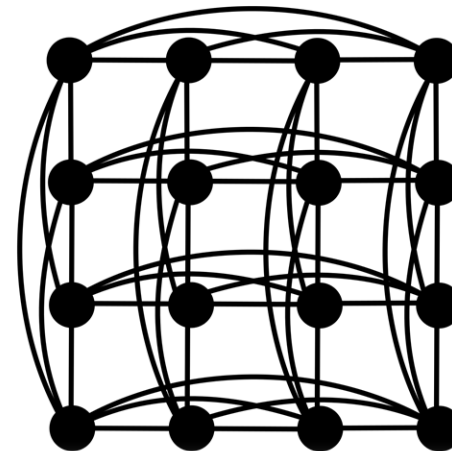
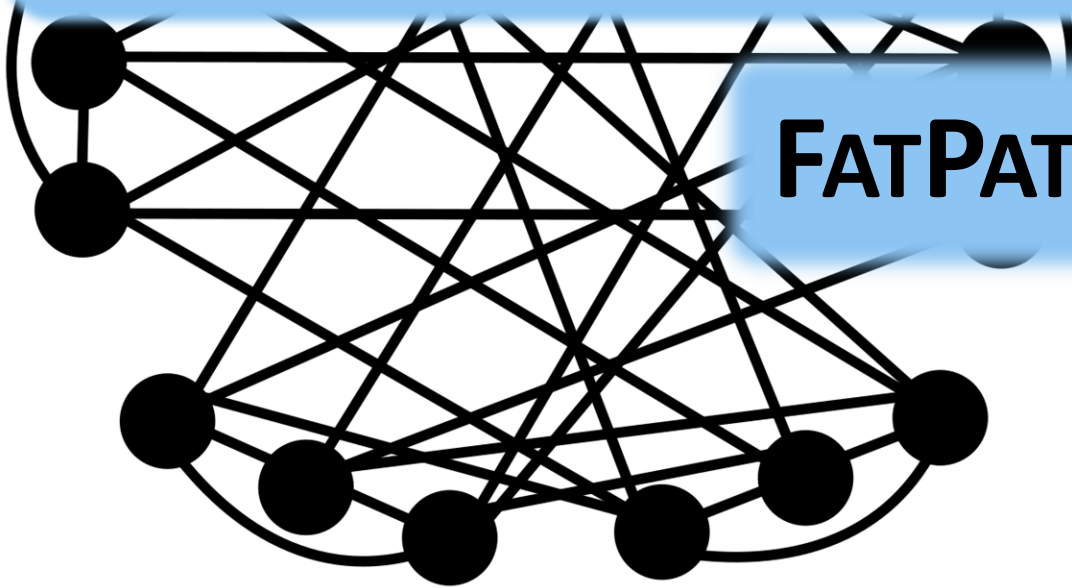
Default topology  
Shortest path: 2 hops

Layer 1: "Almost"-shortest path: 3 hops

Layer 2: "Almost"-shortest path: 3 hops



# FATPATHS OVERVIEW



## A ROUTING SCHEME FOCUSING ON LOW-DIAMETER TOPOLOGIES

**LOW-DIAMETER NETWORK TOPOLOGIES**

How to route modern low diameter topologies?

**LAYERED ROUTING** How to encode & use diversity of shortest and non-minimal paths?

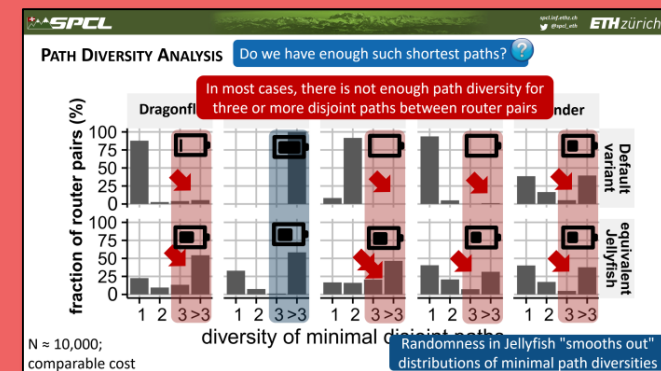
- (a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links
- (a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)
- (b) Divide one flow into subflows and send subflows across different layers
- (c) Route minimally in each layer

Default topology  
Shortest path: 2 hops

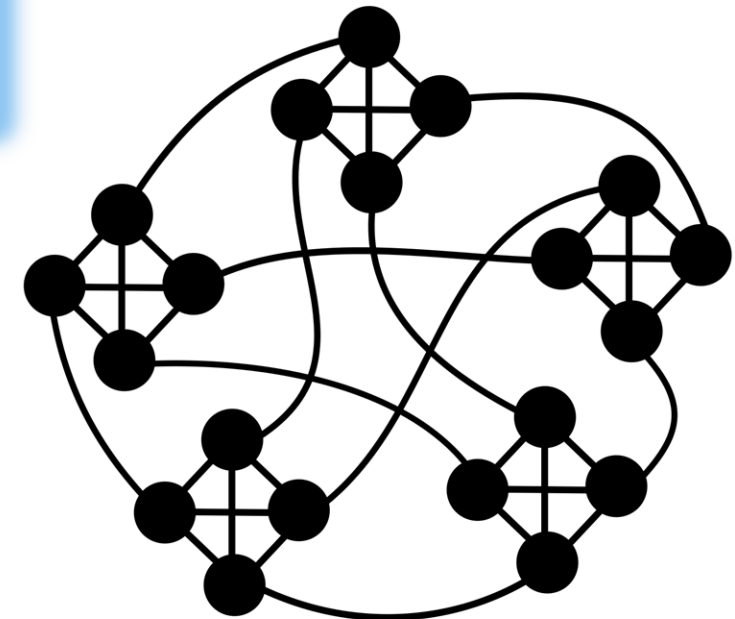
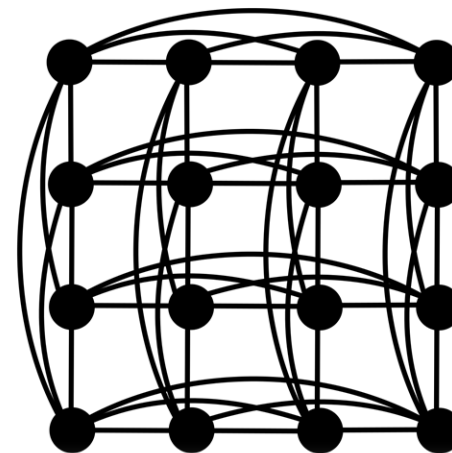
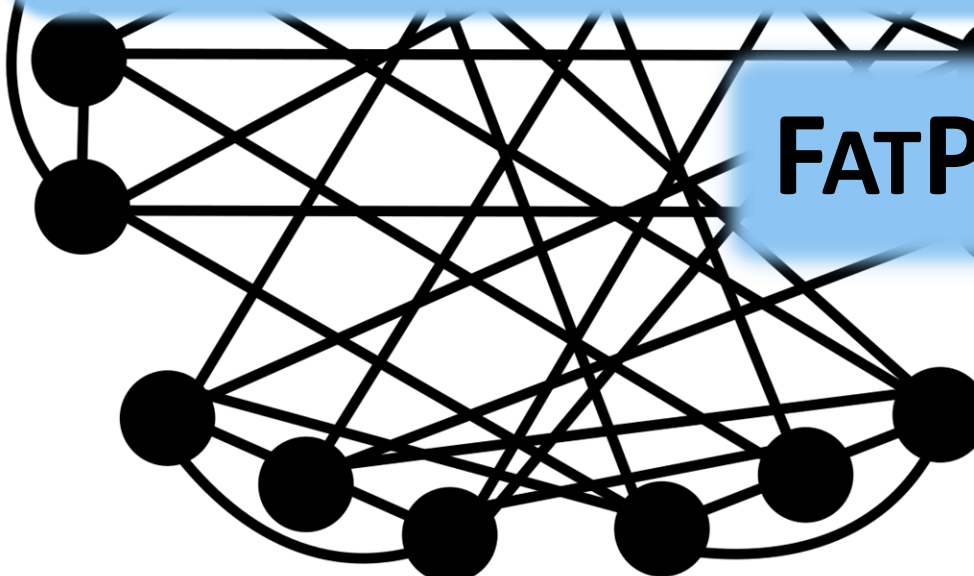
Layer 1: "Almost"-shortest path: 3 hops

Layer 2: "Almost"-shortest path: 3 hops

## BROAD ANALYSIS OF PATH DIVERSITY



# FATPATHS OVERVIEW



## A ROUTING SCHEME FOCUSING ON LOW-DIAMETER TOPOLOGIES

**LOW-DIAMETER NETWORK TOPOLOGIES**

How to route modern low diameter topologies?

**LAYERED ROUTING** How to encode & use diversity of shortest and non-minimal paths?

- (a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links
- (a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)
- (b) Divide one flow into subflows and send subflows across different layers
- (c) Route minimally in each layer

Default topology: Shortest path: 2 hops

Layer 1: "Almost"-shortest path: 3 hops

Layer 2: "Almost"-shortest path: 3 hops

## BROAD ANALYSIS OF PATH DIVERSITY

**PATH DIVERSITY ANALYSIS** Do we have enough such shortest paths?

In most cases, there is not enough path diversity for three or more disjoint paths between router pairs.

fraction of router pairs (%)

diversity of minimal disjoint paths

Default variant, Jellyfish, Randomness in Jellyfish "smooths out" distributions of minimal path diversities

N ≈ 10,000; comparable cost

# FATPATHS OVERVIEW

## HIGH-PERFORMANCE NETWORKING ARCHITECTURE

**TRANSPORT DESIGN** How to maximize performance of the transport layer?

Adapt recent flow control for fat trees [1]

Prioritize packets with dropped payload and retransmitted packets

Key design choice: Drop only payload if router buffers fill up

Remaining header, Router buffer, Dropped payload, A special queue for packets that dropped payload and for retransmitted packets

[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

**LOAD BALANCING** How to load balance traffic over shortest and non-minimal paths?

Use flowlet switching

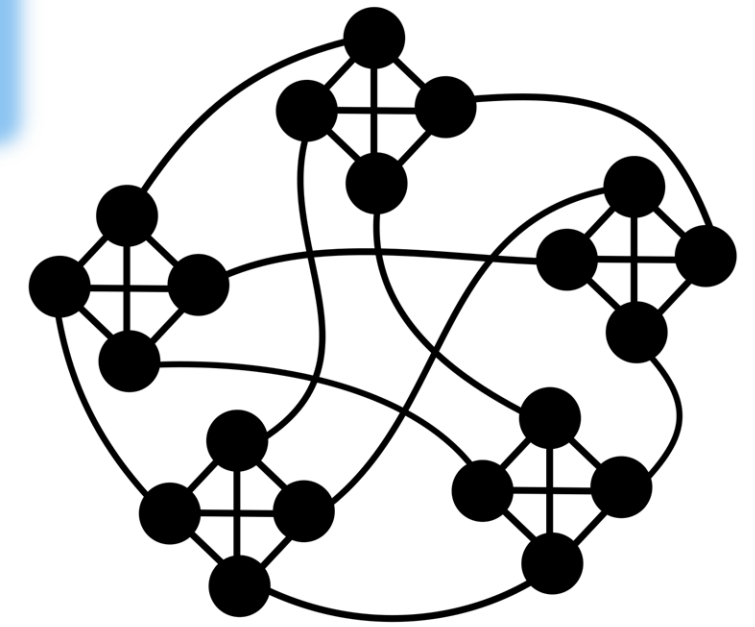
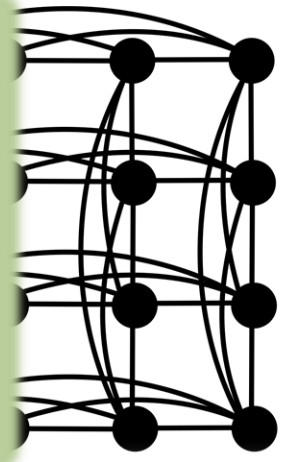
A flowlet [1] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.

Very simple load balancing: a router just picks a random path (layer) for each flowlet, without any probing for congestion.

Size of flowlets changes automatically based on conditions in the network.

On longer paths, flowlets tend to be smaller, and vice versa.

[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR, 2007.





## A ROUTING SCHEME FOCUSING ON LOW-DIAMETER TOPOLOGIES

**LOW-DIAMETER NETWORK TOPOLOGIES**

How to route modern low diameter topologies?

**LAYERED ROUTING** How to encode & use diversity of shortest and non-minimal paths?

(a.1) Divide links into subsets (layers). A minimal route in one layer is usually non-minimal when considering all links

(a.2) Create a layer by removing a fraction of edges (e.g., random uniform sampling)

(b) Divide one flow into subflows and send subflows across different layers

(c) Route minimally in each layer

**Default topology**  
Shortest path: 2 hops

**Layer 1: "Almost"-shortest path: 3 hops**

**Layer 2: "Almost"-shortest path: 3 hops**

## BROAD ANALYSIS OF PATH DIVERSITY

**PATH DIVERSITY ANALYSIS** Do we have enough such shortest paths?

In most cases, there is not enough path diversity for three or more disjoint paths between router pairs under

fraction of router pairs (%)

Default variant Jellyfish

diversity of minimal disjoint paths

Randomness in Jellyfish "smooths out" distributions of minimal path diversities

N ≈ 10,000; comparable cost

# FATPATHS OVERVIEW

## HIGH-PERFORMANCE NETWORKING ARCHITECTURE

**TRANSPORT DESIGN** How to maximize performance of the transport layer?

Adapt recent flow control for fat trees [1]

Prioritize packets with dropped payload and retransmitted packets

Key design choice: Drop only payload if router buffers fill up

Remaining header  
Router buffer  
Dropped payload

A special queue for packets that dropped payload and for retransmitted packets

**LOAD BALANCING** How to load balance traffic over shortest and non-minimal paths?

Use flowlet switching

A flowlet [1] is a sequence of packets within one flow, separated from other flowlets by sufficient time gaps, which prevents packet reordering at the receiver.

Very simple load balancing: a router just picks a random path (layer) for each flowlet, without any probing for congestion.

Size of flowlets changes automatically based on conditions in the network.

On longer paths, flowlets tend to be smaller, and vice versa.

## RICH EVALUATION, THEORETICAL ANALYSIS, TAXONOMY

**EVALUATION**

Number of router pairs (%)

Number of paths (%)

Number of endpoints

SP DF HX3 XP FT3 SP-JF

Routing Scheme (Name, Abbreviation, Reference)	Stack	Supported path diversity aspect
	Layer	SP NP SM MP DP ALB AT
<b>(1) SIMPLE ROUTING PROTOCOLS (often used as building blocks):</b>		
Variant load balancing (VLB) [17]	L2-L3	○ ● ● ● ● ● ● ●
Simple Spanning Tree (ST) [18]	L2	○ ● ● ● ● ● ● ●
Simple routing, e.g., OSPF [19-22]	L2, L3	○ ● ● ● ● ● ● ●
USUAL [20]	L2-L3	○ ● ● ● ● ● ● ●
ECMP [11] OMP [24], Pkt. Spraying [25]	L2, L3	○ ● ● ● ● ● ● ●
<b>(2) ROUTING ARCHITECTURES:</b>		
DCoil [26]	L2-L3	○ ● ● ● ● ● ● ●
Moonshot [27]	L2, L3	○ ● ● ● ● ● ● ●
Portland [9]	L2	○ ● ● ● ● ● ● ●
DRILL [28], LocalFlow [29], DRB [30]	L2	○ ● ● ● ● ● ● ●
VL2 [21]	L2	○ ● ● ● ● ● ● ●
Architecture by Al-Fares et al. [32]	L2-L3	○ ● ● ● ● ● ● ●
ESCube [33]	L2-L3	○ ● ● ● ● ● ● ●
SEATTLE [34], others [35]-[38]	L2	○ ● ● ● ● ● ● ●
VRD [39]	L2-L3	○ ● ● ● ● ● ● ●
Ethernet on Air [40]	L2	○ ● ● ● ● ● ● ●
PAST [41]	L2	○ ● ● ● ● ● ● ●
MLAG, MLAG, others [42]	L2	○ ● ● ● ● ● ● ●
MOOSE [43]	L2	○ ● ● ● ● ● ● ●
NPA [44]	L3	○ ● ● ● ● ● ● ●
ZAMP [45]	L2	○ ● ● ● ● ● ● ●
MSTP [46], QDE [47], Viking [48]	L2	○ ● ● ● ● ● ● ●
SPR [49], TRILL [50], Shadow MACs [51]	L2	○ ● ● ● ● ● ● ●
SPAN [52]	L2	○ ● ● ● ● ● ● ●
<b>(3) Schemes for exposing/encoding paths (can be combined with FatPaths):</b>		
XPath [53]	L3	○ ● ● ● ● ● ● ●
Source routing for flexible DC fabric [54]	L3	○ ● ● ● ● ● ● ●
<b>(4) FatPaths (This work)</b>		
	L2-L3	○ ● ● ● ● ● ● ●

[1] M. Handley et al. Re-architecting datacenter networks and stacks for low latency and high performance. SIGCOMM'17.

[1] S. Kandula et al. Dynamic load balancing without packet reordering. ACM CCR, 2007.