#### Advanced Computer Networks

# Flow Scheduling in Data Center Networks (I)

https://pages.cs.wisc.edu/~mgliu/CS740/F25/index.html

Ming Liu mgliu@cs.wisc.edu

#### Outline

- Last lecture
  - Addressing and routing in data center networks (II)

- Today
  - Flow scheduling in data center networks (I)

- Announcements
  - Project proposal due 10/02/2025 11:59 PM
  - Lab1 due 10/08/2025 11:59 PM

#### Where we are?

**Jata Center Network** 

Multiple communication paths exist when accessing and traversing data center networks!

Physical Connectivity + Networking Architecture (L1, L2, L3)

#### Where we are?

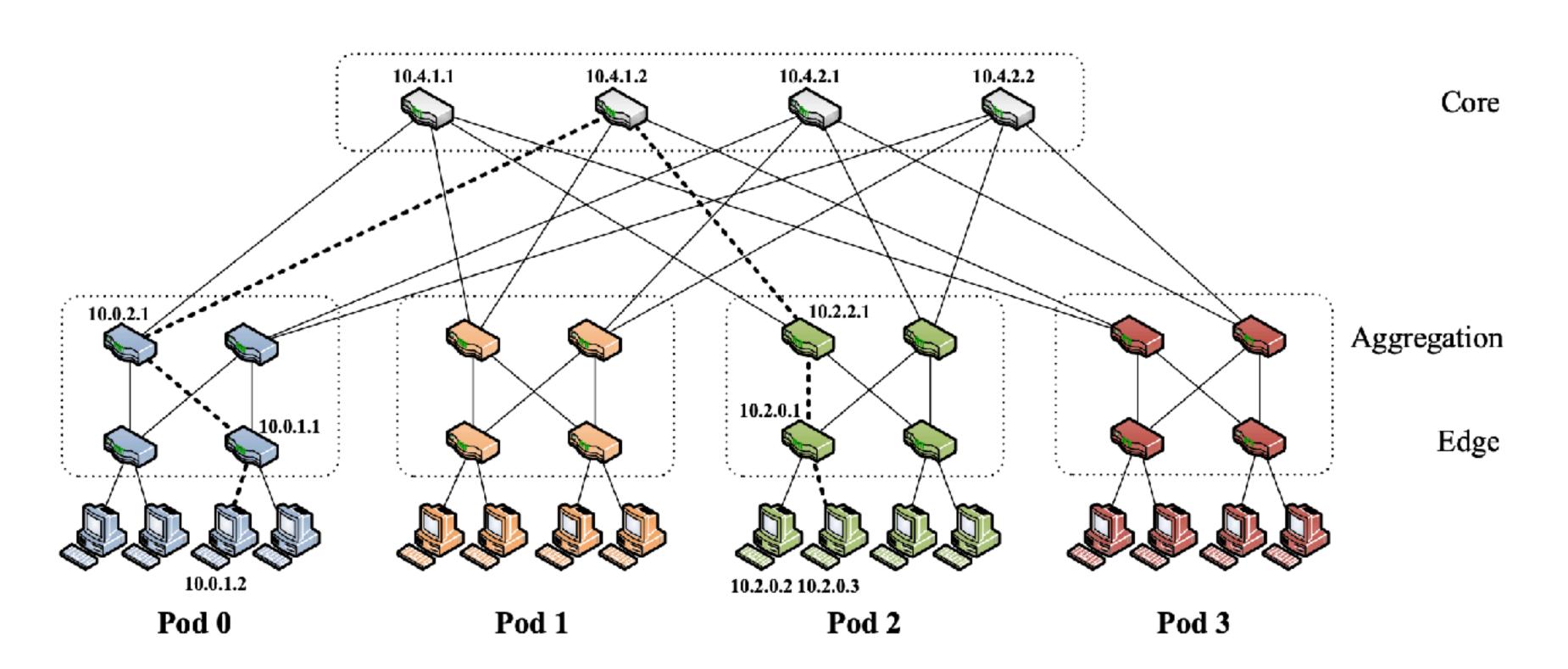
Jata Center Network

The forwarding (destination) address and routing table determine how packets are forwarded!

Addressing and Routing (L4, L5)

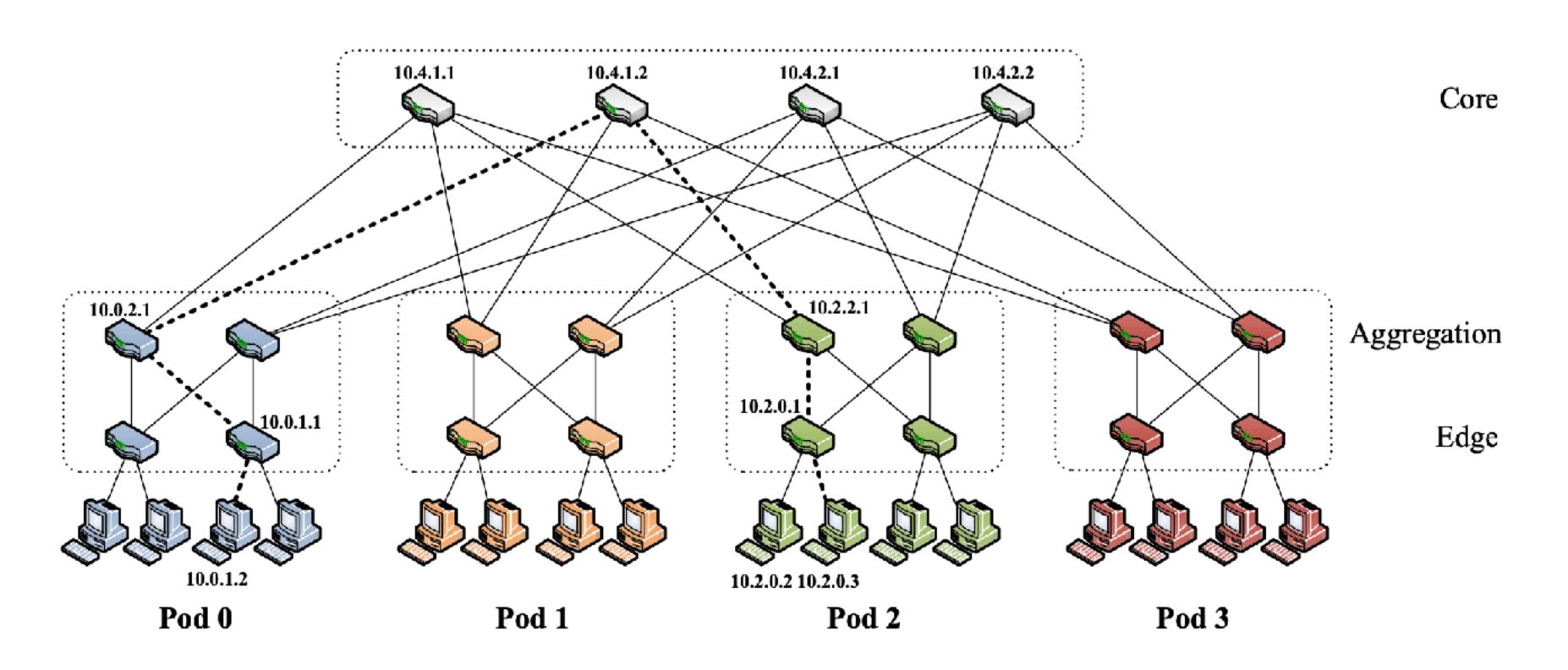
Physical Connectivity + Networking Architecture (L1, L2, L3)

# Which path should a network flow take?



#### What is a "network flow"?

# Which path should a network flow take?



A unique identifier in the data center networks

- A unique identifier in the data center networks: 5-Tuple
  - Source IP
  - Destination IP
  - Source Port
  - Destination Port
  - Protocol Number

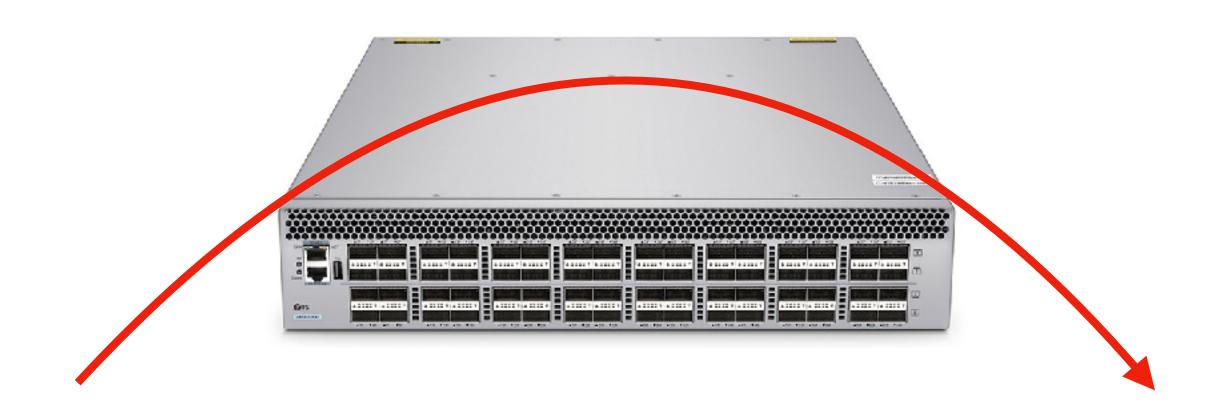
Host-Host communication channel

App-App communication channel

- A unique identifier in the data center networks: 5-Tuple
  - Source IP
  - Destination IP
  - Source Port
  - Destination Port
  - Protocol Number

Host-Host communication channel

App-App communication channel



- A unique identifier in the data center networks: 5-Tuple
  - Source IP
  - Destination IP
  - Source Port
  - Destination Port
  - Protocol Number

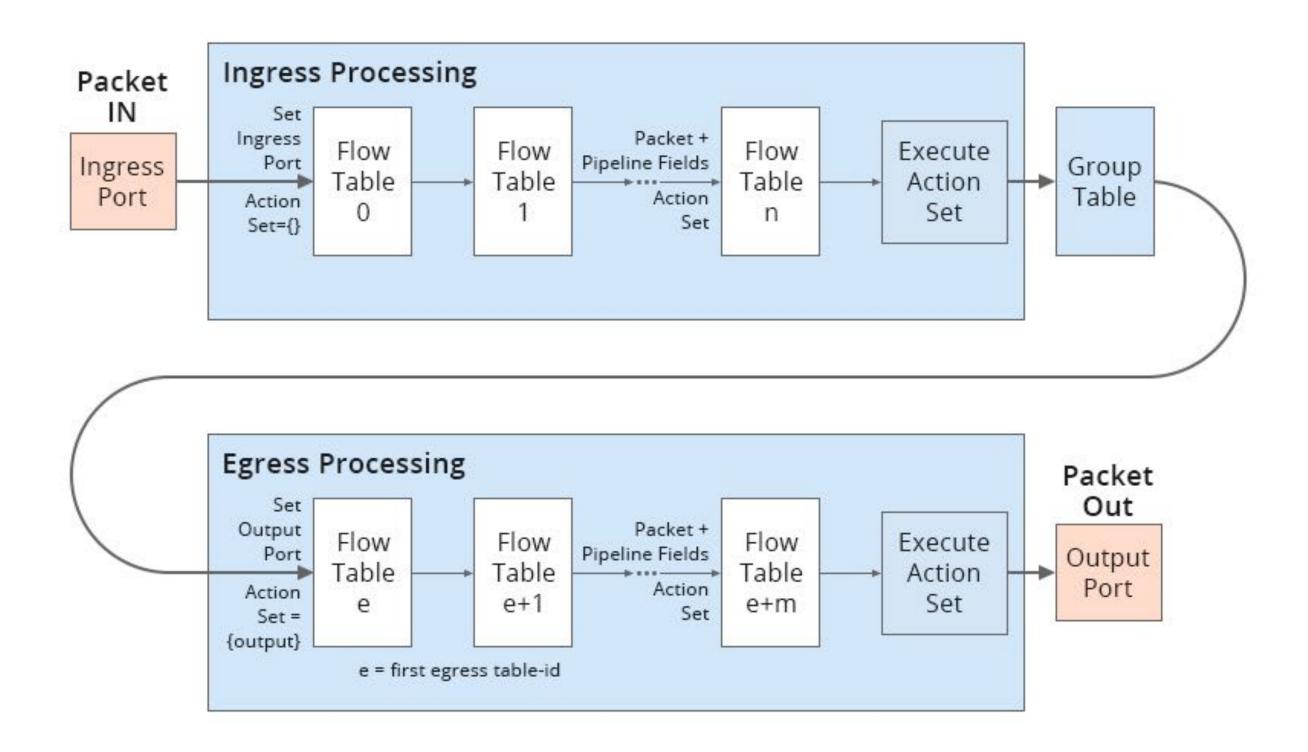
Host-Host communication channel

App-App communication channel

How to represent a network flow in a switch?

#### Some Switch Internals

- Switch ASIC: per-packet processing under high bandwidth!
  - Ingress/egress pipeline: table lookup, header modifications
  - Traffic manager: packet scheduling, rate limiting, and mirroring



## Query 5-Tuple is costly!

- High memory footprint and compute cycles
  - 13B=src\_ip(4B)+dst\_ip(4B)+src\_port(2B)+dst\_port(2B)+procol\_num(1B)
  - 3 SRAM lookups

#### Query 5-Tuple is costly!

- High memory footprint and compute cycles
  - 13B=src\_ip(4B)+dst\_ip(4B)+src\_port(2B)+dst\_port(2B)+procol\_num(1B)
  - 3 SRAM lookups

## Let's do hashing!

#### Network Flow Representation in Commodity Switches

- FlowID=Hash(src\_ip, dst\_ip, src\_port, dst\_port, proto\_num)
  - The hashing algorithm is mostly proprietary
  - Perform at the switch ingress side before moving to the pipeline
  - FlowID becomes a N-bit metadata field (N=16, for example)

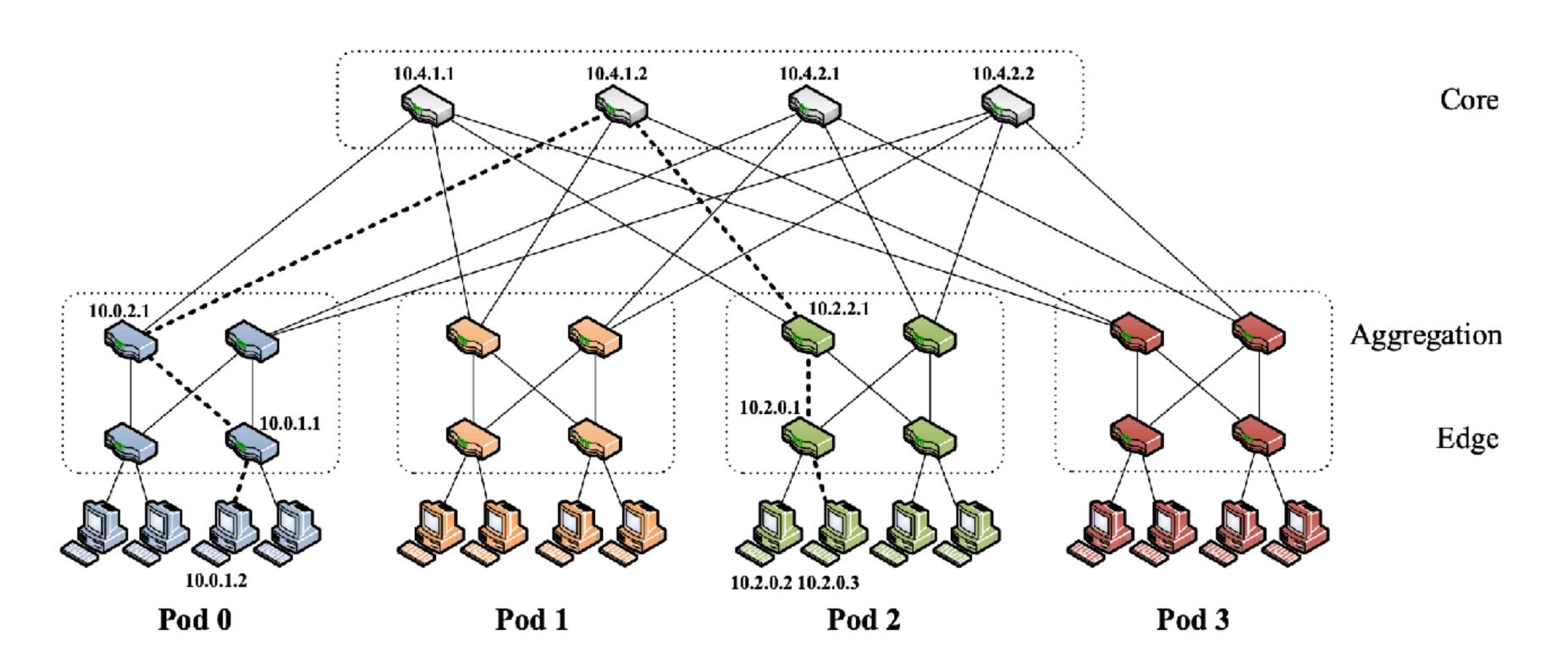
## Is FlowID unique?

#### Network Flow Representation in Commodity Switches

- FlowID=Hash(src\_ip, dst\_ip, src\_port, dst\_port, proto\_num)
  - The hashing algorithm is mostly proprietary
  - Perform at the switch ingress side before moving to the pipeline
  - FlowID becomes a N-bit metadata field (N=16, for example)

# Is FlowID unique? No, because of hash collosion

# Which path should a network flow take?



- Equal-Cost Multi-Path Fordwaring
  - chosen\_path = FlowID mod (path #)

- Equal-Cost Multi-Path Fordwaring
  - chosen\_path = FlowID mod (path #)
- Cost = The number of hops
- Simple and it works!

- Equal-Cost Multi-Path Fordwaring
  - chosen\_path = FlowID mod (path #)
- Cost = The number of hops
- Simple and it works!

Is this aligned with our discussion last week?

- Equal-Cost Multi-Path Fordwaring
  - chosen\_path = FlowID mod (path #)
- Cost = The number of hops
- Simple and it works!

Is this aligned with our discussion last week?

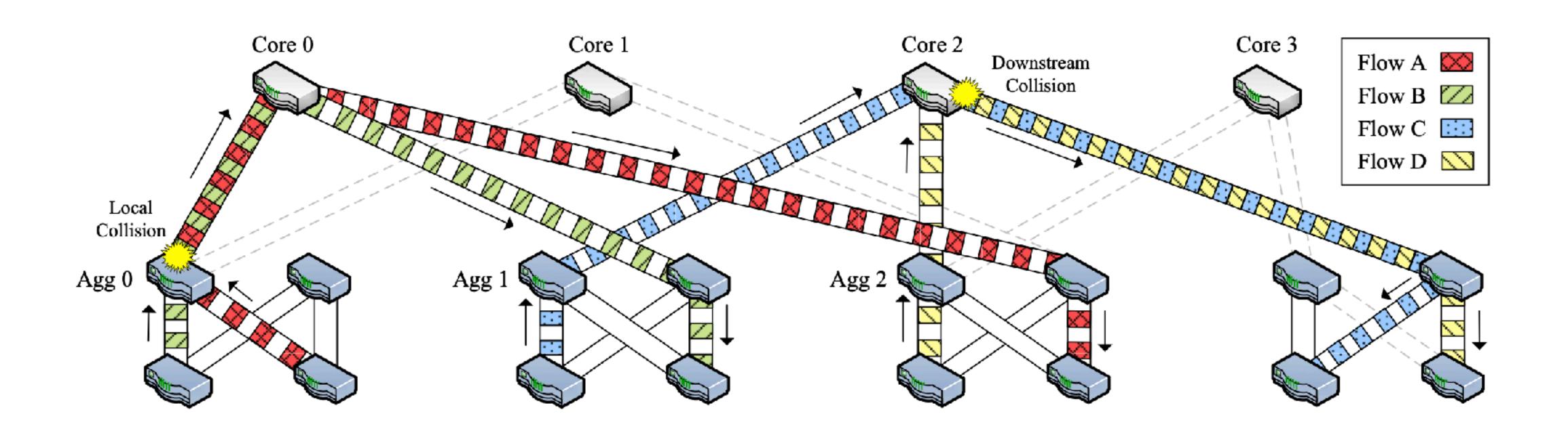
No. In the scalable DCNet, the host ID

determines the routing path!

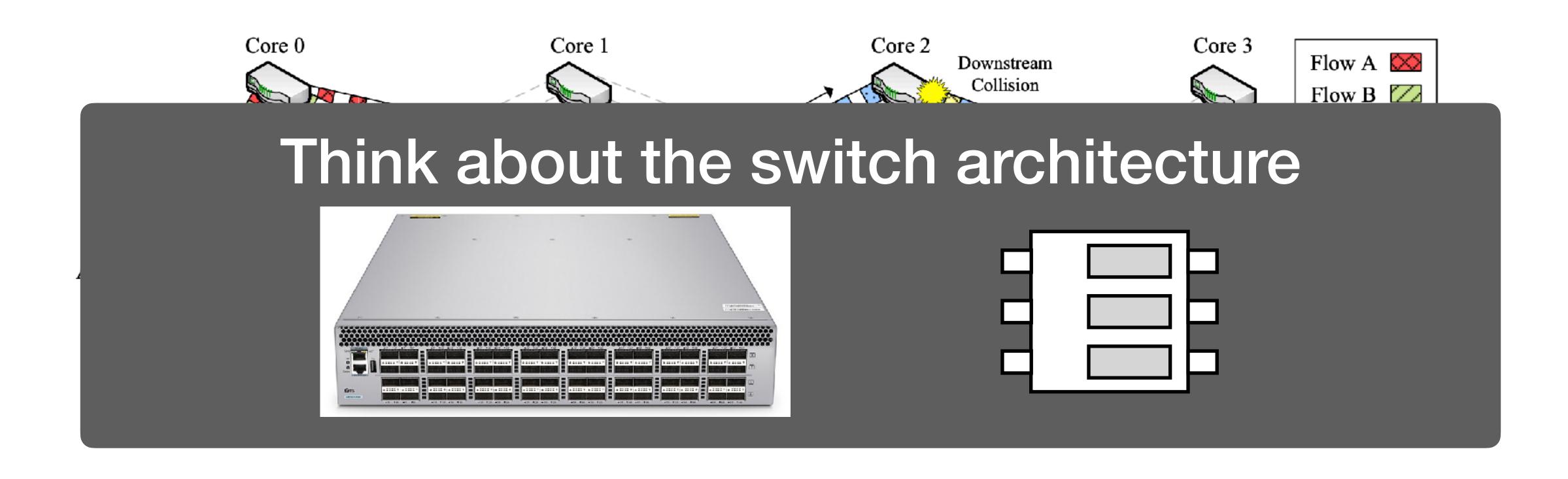
- Equal-Cost Multi-Path Fordwaring
  - chosen\_path = FlowID mod (path #)
- Cost = The number of hops
- Simple and it works!

What are the issues of ECMP?

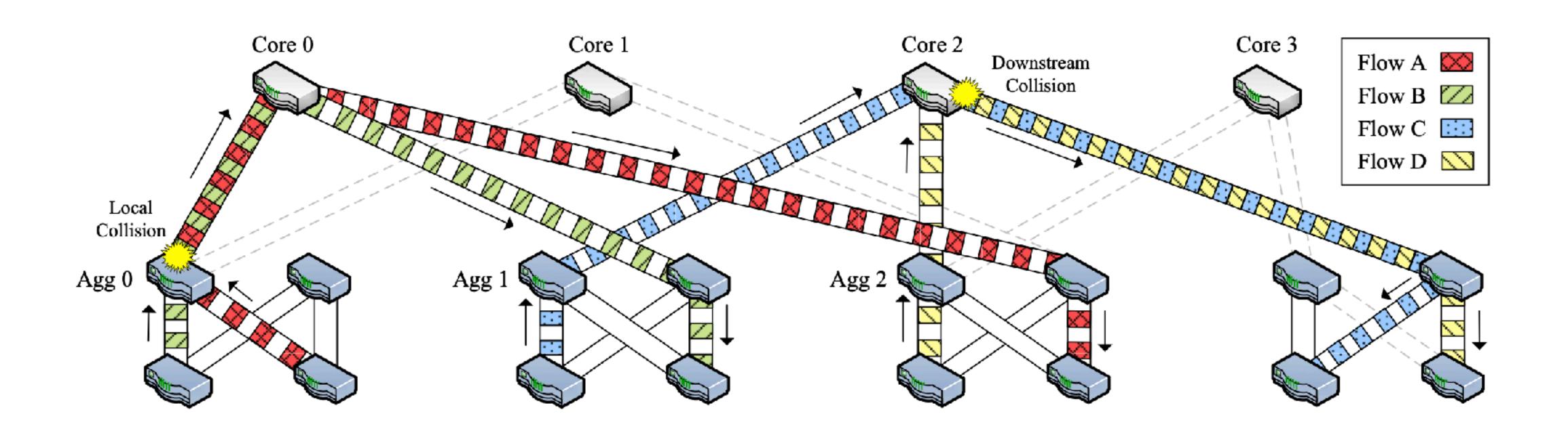
If incoming bandwidth <= outgoing bandwidth</li>



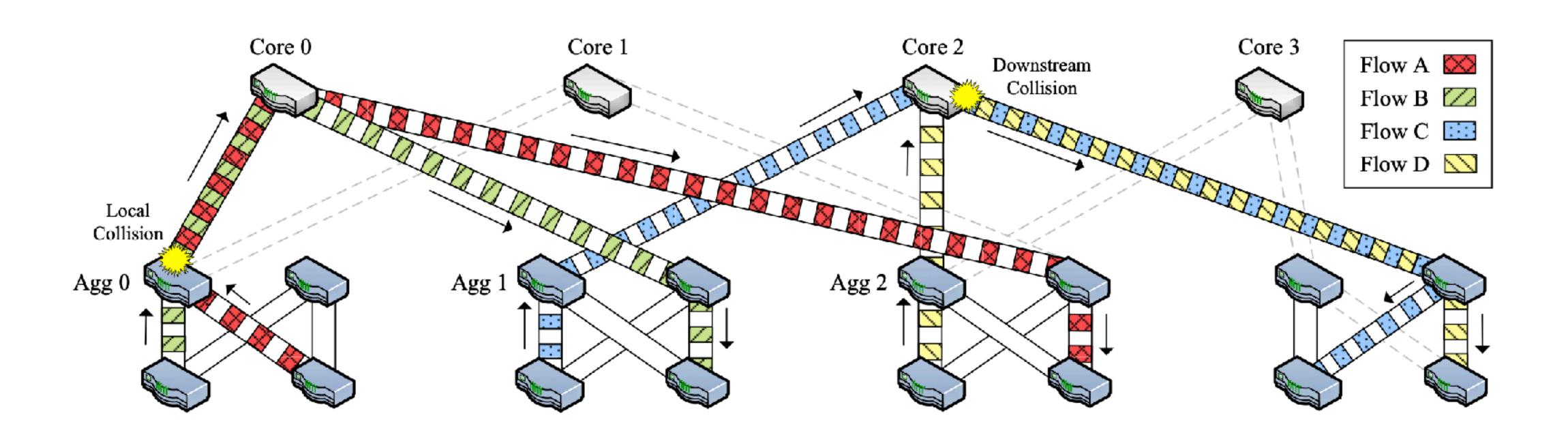
If incoming bandwidth <= outgoing bandwidth</li>



If incoming bandwidth <= outgoing bandwidth Okay!</li>

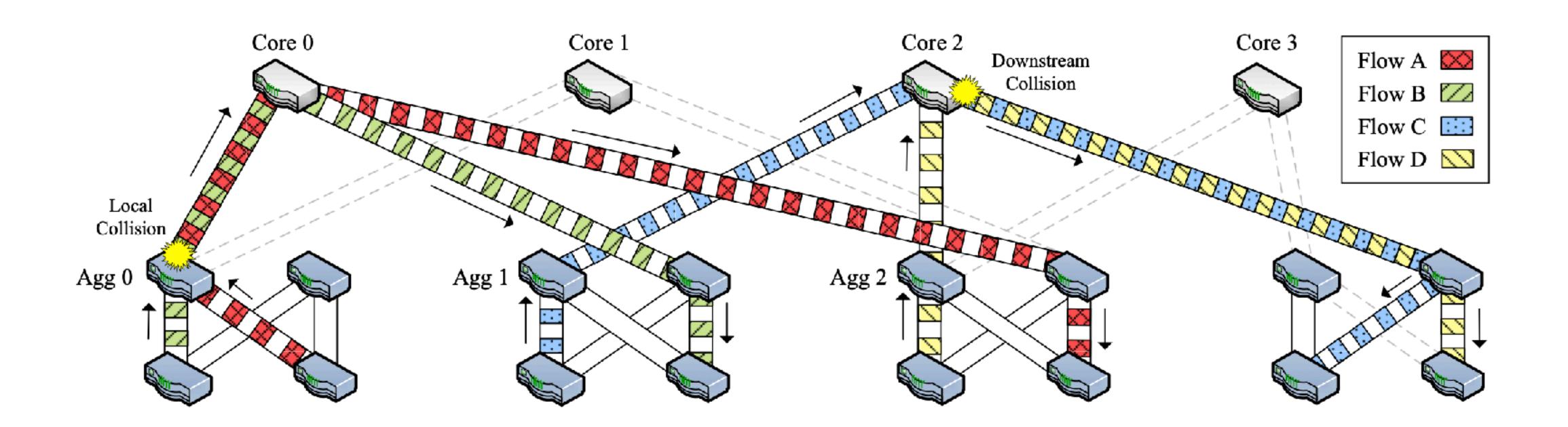


- If incoming bandwidth <= outgoing bandwidth Okay!</li>
- If incoming bandwidth > outgoing bandwidth

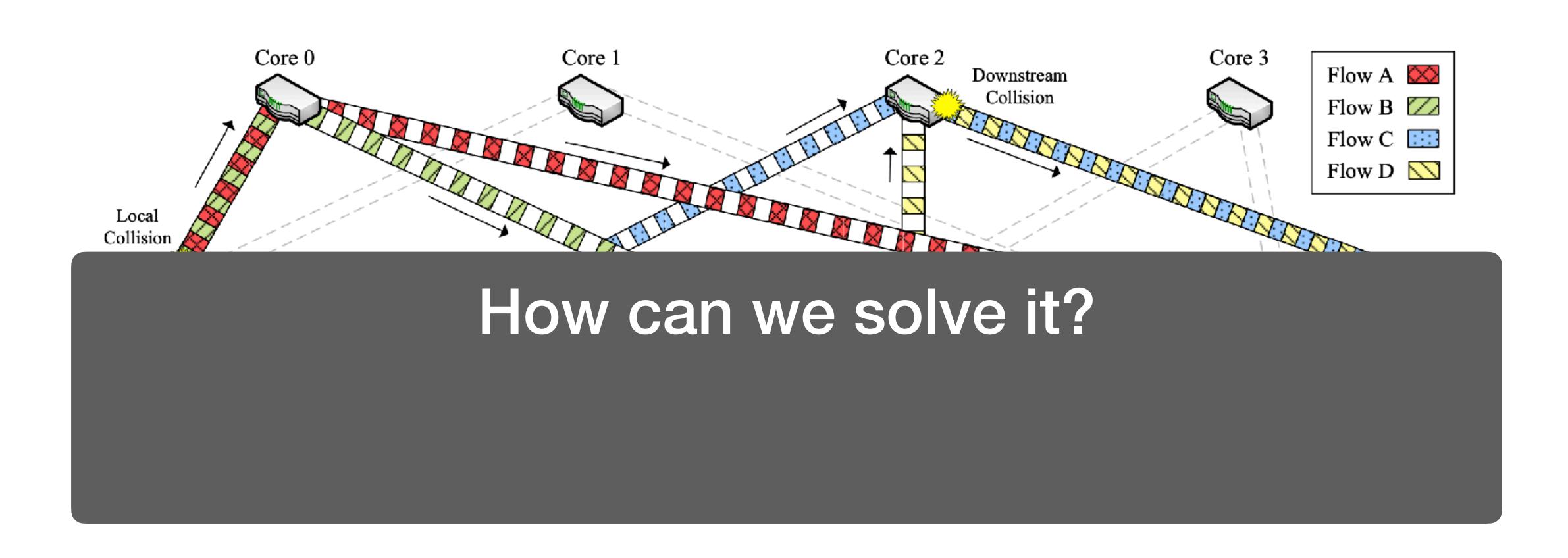


If incoming bandwidth <= outgoing bandwidth Okay!</li>

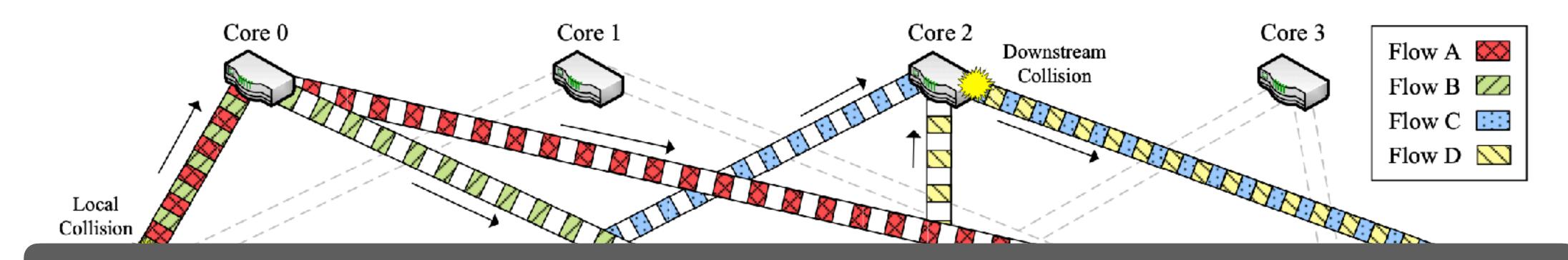
If incoming bandwidth > outgoing bandwidth Bad!



- If incoming bandwidth <= outgoing bandwidth Okay!</li>
- If incoming bandwidth > outgoing bandwidth Bad!



- If incoming bandwidth <= outgoing bandwidth Okay!</li>
- If incoming bandwidth > outgoing bandwidth Bad!



How can we solve it?
Mitigate bandwidth oversubscription by leveraging more paths!

#### Problem Formulation

- Multi-commodity flow problem
  - NP-complete for integer flows
- Approximation hint
  - Match flow demand and link bandwidth supply

#### Problem Formulation

- Multi-commodity flow problem
  - NP-complete for integer flows
- Approximation hint
  - Match flow demand and link bandwidth supply

- How can we know the flow demand?
- How can we know the bandwidth supply?
- Which path should a flow choose?

## Proposal #2: Hedra (NSDI'10)

- Let's start with a simple example
- Network-limited flows

#### An example:

- H0 sends 1 flow to H1, H2, and H3
- H1 sends 2 flows to H0 and 1 flow to H2
- H2 sends 1 flow each to H0 and H3
- H3 sends 2 flows to H1

## Demand Estimation Algorithm

- Traffic matrix: MxN
  - M: the number of senders
  - N: the number of receivers
- Each element in the matrix
  - The number of flows from host i to host j
  - The estimated demand of each of the flows from host i to host j
  - A converged flag that marks flows whose demands have converged

dF-> The "converged" demand
dU-> The number of unconverged flows
eS-> The computed equal share rate
dT-> The total demand
dS-> The sender limited demand
f.rl-> A flag for a receiver limited flow
nR-> The number of receiver limited flow

```
ESTIMATE-DEMANDS()

1 for all i, j

2 M_{i,j} \leftarrow 0

3 do

4 foreach h \in H do EST-SRC(h)

5 foreach h \in H do EST-DST(h)

6 while some M_{i,j}.demand changed

7 return M
```

- dF-> The "converged" demand
- dU-> The number of unconverged flows
- eS-> The computed equal share rate
- dT-> The total demand
- dS-> The sender limited demand
- f.rl-> A flag for a receiver limited flow
- nR-> The number of receiver limited flow

```
ESTIMATE-DEMANDS()

1 for all i, j

2 M_{i,j} \leftarrow 0

3 do

4 foreach h \in H do EST-SRC(h)

5 foreach h \in H do EST-DST(h)

6 while some M_{i,j}.demand changed

7 return M
```

```
EST-SRC(src: host)

1 d_F \leftarrow 0

2 n_U \leftarrow 0

3 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle do

4 if f.\operatorname{converged} then

5 d_F \leftarrow d_F + f.\operatorname{demand}

6 else

7 n_U \leftarrow n_U + 1

8 e_S \leftarrow \frac{1.0 - d_F}{n_U}

9 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle and not f.\operatorname{converged} do

10 M_{f.\operatorname{Src},f.\operatorname{dst}}.\operatorname{demand} \leftarrow e_S
```

```
    dF-> The "converged" demand
```

- dU-> The number of unconverged flows
- eS-> The computed equal share rate
- dT-> The total demand
- dS-> The sender limited demand
- f.rl-> A flag for a receiver limited flow
- nR-> The number of receiver limited flow

```
ESTIMATE-DEMANDS()

1 for all i, j

2 M_{i,j} \leftarrow 0

3 do

4 foreach h \in H do EST-SRC(h)

5 foreach h \in H do EST-DST(h)

6 while some M_{i,j}.demand changed

7 return M
```

```
EST-SRC(src: host)

1 d_F \leftarrow 0

2 n_U \leftarrow 0

3 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle do

4 if f.converged then

5 d_F \leftarrow d_F + f.demand

6 else

7 n_U \leftarrow n_U + 1

8 e_S \leftarrow \frac{1.0 - d_F}{n_U}

9 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle and not f.converged do

10 M_{f.\operatorname{src},f.\operatorname{dst}}.demand \leftarrow e_S
```

```
EST-DST(dst: host)
         d_T, d_S, n_R \leftarrow 0
         foreach f \in \langle \operatorname{src} \to \operatorname{dst} \rangle
                  f.rl \leftarrow true
                 d_T \leftarrow d_T + f.demand
         if d_T \leq 1.0 then
                 return
                 foreach f \in \langle \operatorname{src} \to \operatorname{dst} \rangle and f.\operatorname{rl} do
                         if f.demand < e_S then
                                 d_S \leftarrow d_S + f.demand
                                 f.rl \leftarrow false
                         else
                                n_R \leftarrow n_R + 1
                e_S \leftarrow \frac{1.0 - d_S}{2}
         while some f.rl was set to false
         foreach f \in \langle \operatorname{src} \to \operatorname{dst} \rangle and f.\operatorname{rl} do
                 M_{f. \mathrm{src}, f. \mathrm{dst}}. \mathrm{demand} \leftarrow e_S
                 M_{f. \mathrm{src}, f. \mathrm{dst}}. \mathrm{converged} \leftarrow \mathrm{true}
```

- dF-> The "converged" demand
- dU-> The number of unconverged flows
- eS-> The computed equal share rate
- dT-> The total demand

```
EST-SRC(src: host)

1 d_F \leftarrow 0

2 n_U \leftarrow 0

3 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle do

4 if f.converged then

5 d_F \leftarrow d_F + f.demand

6 else

7 n_U \leftarrow n_U + 1
```

```
\begin{bmatrix} & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ (\frac{1}{3})_2 & & (\frac{1}{3})_1 & 0_0 \\ (\frac{1}{2})_1 & 0_0 & & (\frac{1}{2})_1 \\ 0_0 & (\frac{1}{2})_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_2 & & (\frac{1}{3})_1 & 0_0 \\ (\frac{1}{3})_1 & 0_0 & & (\frac{1}{2})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_2 & & (\frac{1}{3})_1 & 0_0 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_2 & & (\frac{1}{3})_1 & 0_0 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 & (\frac{1}{3})_1 \\ [\frac{1}{3}]_1 & 0_0 & & (\frac{2}{3})_1 \\ 0_0 & [\frac{1}{3}]_2 & 0_0 \end{bmatrix} \Rightarrow \begin{bmatrix} & [\frac{1}{3}]_1 & (\frac{1}{3})_1 & (\frac{1}{3})_
```

```
ESTIMATE-DEMANDS()

1 for all i, j

2 M_{i,j} \leftarrow 0

3 do

4 foreach h \in H do EST-SRC(h)

5 foreach h \in H do EST-DST(h)

6 while some M_{i,j}.demand changed

7 return M
```

```
d_T \leftarrow d_T + f. 	ext{demand}
n_R \leftarrow n_R + 1

if d_T \leq 1.0 then
return
e_S \leftarrow \frac{1.0}{n_R}

do

n_R \leftarrow 0

foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle and f. \operatorname{rl} do

if f. 	ext{demand} < e_S then

d_S \leftarrow d_S + f. 	ext{demand}

f. \operatorname{rl} \leftarrow \operatorname{false}

else

n_R \leftarrow n_R + 1

e_S \leftarrow \frac{1.0 - d_S}{n_R}

while some f. \operatorname{rl} was set to false
foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle and f. \operatorname{rl} do

M_{f. \operatorname{src}, f. \operatorname{dst}} \cdot \operatorname{demand} \leftarrow e_S

M_{f. \operatorname{src}, f. \operatorname{dst}} \cdot \operatorname{converged} \leftarrow \operatorname{true}
```

- dF-> The "converged" demand
- dU-> The number of unconverged flows
- eS-> The computed equal share rate
- dT-> The total demand
- dS-> The sender limited demand
- f.rl-> A flag for a receiver limited flow
- nR-> The number of receiver limited flow

```
EST-SRC(src: host)

1 d_F \leftarrow 0

2 n_U \leftarrow 0

3 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle do

4 if f.\operatorname{converged} then

5 d_F \leftarrow d_F + f.\operatorname{demand}

6 else

7 n_U \leftarrow n_U + 1

8 e_S \leftarrow \frac{1.0 - d_F}{n_U}

9 foreach f \in \langle \operatorname{src} \rightarrow \operatorname{dst} \rangle and not f.\operatorname{converged} do

10 M_{f.\operatorname{src}, f.\operatorname{dst}}.\operatorname{demand} \leftarrow e_S
```

```
EST-DST(dst: host)

1 d_T, d_S, n_R \leftarrow 0
```

- SRC host: max out the bandwidth supply
- DST host: restrict the bandwidth demand
- Maxmin fairness at a global scale

```
while some f.rl was set to false foreach f \in (\operatorname{src} \to \operatorname{dst}) and f.rl do M_{f.\operatorname{src},f.\operatorname{dst}}.demand \leftarrow e_S M_{f.\operatorname{src},f.\operatorname{dst}}.converged \leftarrow true
```

#### Path Selection

- Global First Fit
  - Search all possible paths and choose the first fit one

- Simulated Annealing
  - Choose the one that satisfies the energy function
  - Pruning the search space is key
  - Assign a single core switch for each destination host

#### Hedra in Practice

- #1: Central scheduler
  - Generate the forwarding table
  - Distribute to all switches
- #2: Flow demand estimator
  - Max-min fairness based on the traffic matrix
- #3: Approximate path selection
  - Employ heuristic algorithms

#### Hedra in Practice

- #1: Central scheduler
  - Generate the forwarding table
  - Distribute to all switches
- #2: Flow demand estimator
  - Max-min fairness based on the traffic matrix
- #3: Approximate path selection
  - Employ heuristic algorithms

## Does Herdra always work?

## Challenges in Deployment

• #1: Networking bandwidth capacity (supply) is ephemeral!

- #2: Flow demand and flow # are unpredictable!
  - Hedra targets long-lived flows
- #3: Real-time global view is hard to maintain!

#### Summary

- Today
  - Flow scheduling in data center networks (I)

- Next
  - Flow scheduling in data center networks (II)
  - Conga (SIGCOMM'14)