

# Priority-Aware Scheduling for Packet-Switched Optical Networks in Datacenter

Lin Wang\*, Xinbo Wang\*, **Massimo Tornatore**\*<sup>†</sup>, Kwang joon Kim<sup>‡</sup>, Sun Me Kim<sup>‡</sup>, Dae-Ub Kim<sup>‡</sup>, Kyeong-Eun Han<sup>‡</sup>, and Biswanath Mukherjee\*

\*University of California Davis, USA †Politecnico di Milano, Italy ‡Electronics and Telecommunications Research Institute, Korea











#### **Packet-Switched Optical Network (PSON) Architecture**

- Top-of-rack (ToR) switch connected with servers
- Ingress module/egress module with framer/de-framer per switch
- A space switch connects with *m* Arrayed Waveguide Grating Routers (AWGR)
- Scalable core switching consists of AWGRs
- Control interface management





# **PSON data plane (with optical switch fabric)**



Space switch: Optical path switch to switch optical signal between module and AWGR at sub micro-second speed.





## **Mice VS. Elephant Flow**





- Small size packet
- Large number
- Short flow
- Short-lived

transactional traffic, web browsing, search queries (≈ 80% traffic)



- Large size packet
- Small number
- Large volume flow
- Long-lasting

bulk data transfer, data backup, virtual machine migration (≈ 20% traffic)



# **Mice VS. Elephant Flow**

ERSITY OF CALIFORN

 If schedule first queues with larger number of packets, elephant flows will be penalized



 If schedule first queues with larger total size of packets, mice flows will be penalized



#### Iterative Round Robin algorithm

- ✓ Step 1: Request
- ✓ Step 2: Grant
- ✓ Step 3: Accept



#### • Step 1. Request

- ✓ Each ingress module maintains a queue for each egress module
- $\checkmark$  Send request, if the queue at ingress module is not empty





#### • Step 2. Grant

- ✓ Output maintains a *grant* pointer
- ✓ Output chooses requesting input closest to the grant pointer
- ✓ *Grant* pointer moves





#### • Step 3. Accept

- ✓ Input maintains a *accept* pointer
- ✓ Input accepts granting output closest to the accept pointer
- ✓ Accept pointer moves





- · Pro's
  - Simple
  - Fast
- · Con's
  - $\cdot\,$  does not consider traffic characteristics
  - $\cdot\,$  results in high delays and unfairness



## A Priority-Aware Scheduling Algorithm for PSON



✓ Send only first k VoQs with highest weight



# **Scheduling Algorithm for PSON**





#### **Scheduling Algorithm for PSON**





# **Topology for Simulation**

- **PSON:** *n*=80 ToR switches, 40×40 AWGRs, and 1×2 space switches.
- Each ToR receives input traffic generated by 36 servers.
- Classify 36 servers of each ToR into three groups:
  - 1st group: generates packets which contain 80% elephant flow and 20% mice flow;
  - · 2nd group generates packets with 20% elephant flow and 80% mice flow;
  - 3rd group generates packets with 50% probability.



# **Traffic Generation (1)**

 Packet length follows a bimodal distribution around 40 bytes and 1500 bytes





#### **Traffic Generation (2)**

- Packet arrival times are modeled matching ON/OFF periods
- ON/OFF periods follows Pareto distribution







Average delay of PA and RR algorithms



#### **Results: Packet Loss Ratio**



Packet-loss ratio of PA and RR algorithms



#### **Results: Average Delay (comparison of ω)**



Influence of weight factors on average delay of PA algorithms



#### **Results: Packet Loss Ratio (comparison of ω)**



Influence of weight factors on packet-loss ratio of PA algorithms



#### Conclusion

- We studied scheduling for PSON architecture
- Proposed a Priority-aware (PA) scheduling algorithm
  - $\cdot$  Extension of Round Robin to account for «mice vs elephant» flows
- Performed performance evaluation of PA algorithm
- Significant savings: up to one order of magnitude on blocking and up to 50% on delay





#### amlwang@ucdavis.edu



#### **Scheduling Algorithm for PSON**



