

**I ILLINOIS**

CSL | Coordinated  
Science Lab

COLLEGE OF ENGINEERING

Saurabh Jha<sup>1</sup>, Archit Patke<sup>1</sup>, Jim Brandt<sup>2</sup>, Ann  
Gentile<sup>2</sup>, Benjamin Lim<sup>1</sup>, Mike Showerman<sup>1,3</sup>, Greg  
Bauer<sup>1,3</sup>, Larry Kaplan<sup>4</sup>, Zbigniew T. Kalbarczyk<sup>1</sup> and  
William T. Kramer<sup>1,3</sup>, Ravishankar K. Iyer<sup>1,3</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign    <sup>2</sup>Sandia National Laboratories

<sup>3</sup>National Center for Supercomputing Applications    <sup>4</sup>Cray Inc.

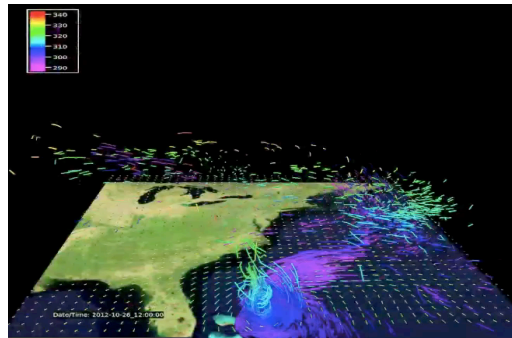
# Measuring Congestion in High-Performance Datacenter Interconnects



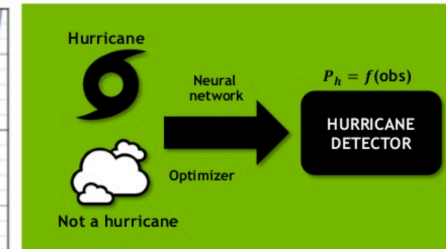
[csl.illinois.edu](http://csl.illinois.edu)

# High-Performance Computing (HPC)

HPC solves critical science, finance, AI, and other problems



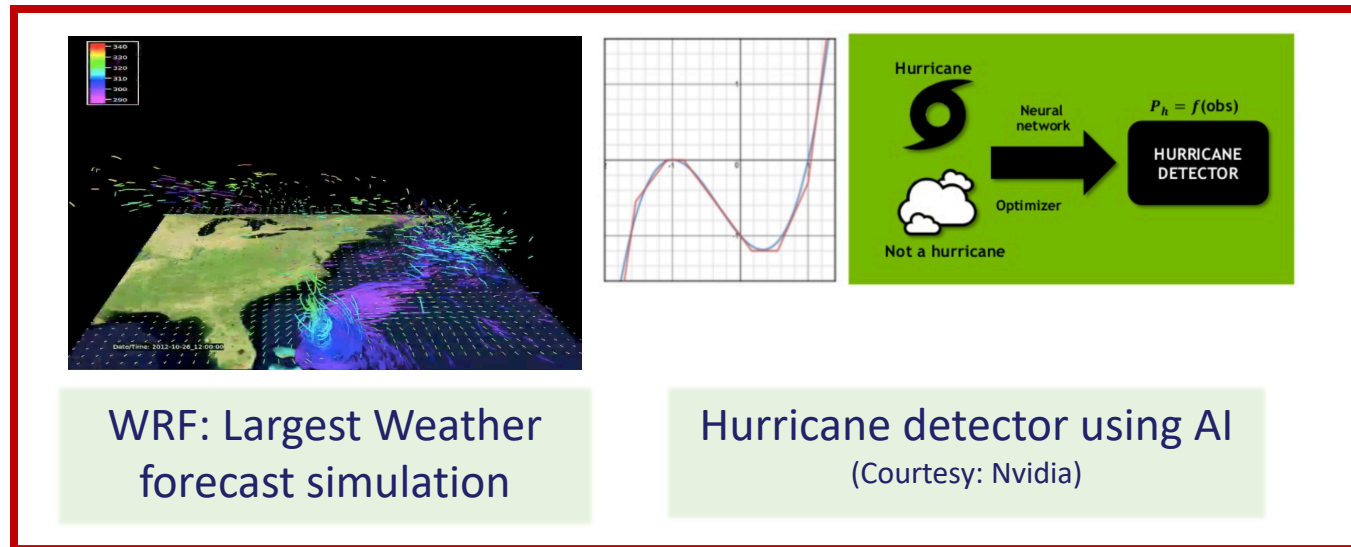
WRF: Largest Weather forecast simulation



Hurricane detector using AI  
(Courtesy: Nvidia)

# High-Performance Computing (HPC)

HPC solves critical science, finance, AI, and other problems



The composite image contains three parts: 1) A 3D visualization of a weather forecast simulation with a color-coded legend on the left and a 'DataTime: 2012-10-26 12:00:00' label at the bottom. 2) A line graph with two curves, one red and one blue, plotted on a grid. 3) A diagram showing a 'Hurricane' icon and a 'Not a hurricane' icon with a cloud, both pointing to a 'Neural network' box. An arrow labeled 'Optimizer' points from the neural network to a 'HURRICANE DETECTOR' box, which outputs the equation  $P_h = f(\text{obs})$ .

WRF: Largest Weather forecast simulation

Hurricane detector using AI  
(Courtesy: Nvidia)

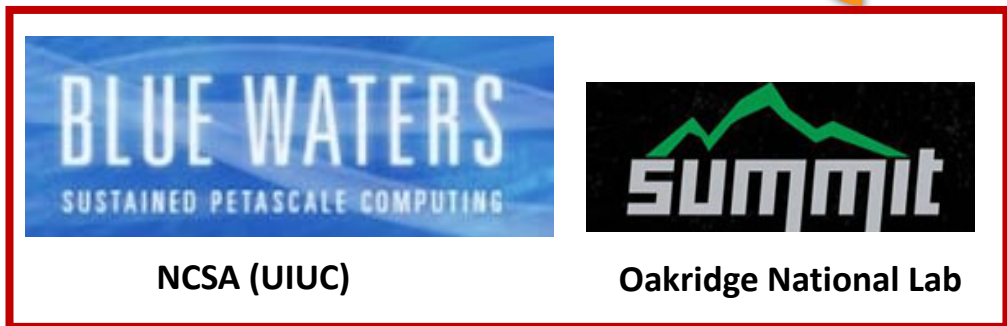
## HPC on Cloud



The illustration shows server racks with blue lights, a large multi-colored cloud icon, and a Microsoft Azure logo. Below the illustration is a dark box with the text 'Cray in Azure' and 'A dedicated supercomputer on your virtual network'.

Cray in Azure  
A dedicated supercomputer on your virtual network

## HPC in Academic and National Labs



The block contains two logos: 'BLUE WATERS SUSTAINED PETASCALE COMPUTING' and 'Summit'. Below the logos are the labels 'NCSA (UIUC)' and 'Oakridge National Lab'.

BLUE WATERS  
SUSTAINED PETASCALE COMPUTING

Summit

NCSA (UIUC)

Oakridge National Lab

# High-Performance Computing (HPC)



## High-speed Networks (HSN)

- Low per-hop latency [1][2]
- Low tail-latency variation
- High bisection bandwidth

[1] <https://www.nextplatform.com/2018/03/27/in-modern-datacenters-the-latency-tail-wags-the-network-dog/>

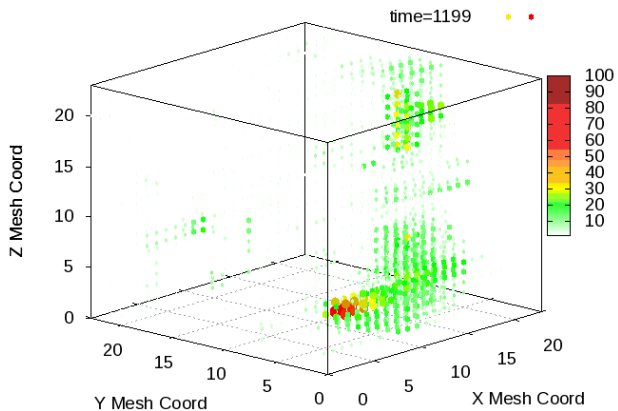
[2] <https://blog.mellanox.com/2017/05/microsoft-enhanced-azure-cloud-efficiency/>

# Networking and Performance Variation

Despite the low-latency, high-speed networks (HSN) are susceptible to high congestion

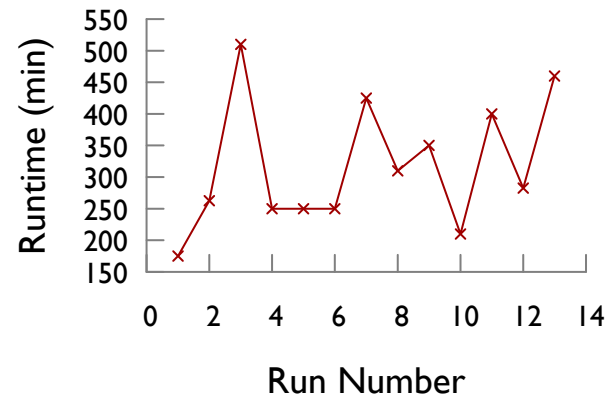
Such congestion can cause up to 2-4X application performance variation in production settings

X+ Gemini Link: Percent Time Spent in Credit Stalls (1 min intervals)



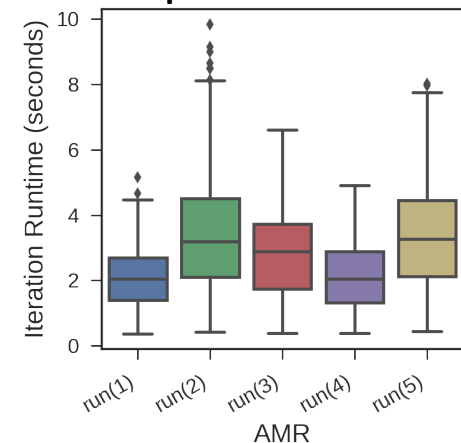
## 1000-node production molecular dynamics code.

Up to 1.89X slowdown compared to median runtime of 282 minutes



## 256-node benchmark app (AMR)

Up to 4X slowdown compared to the median loop iteration time of 2.5 sec



# Networking and Performance Variation

Despite low-latency, high-speed networks (HSN) susceptible to high congestion

Such congestion can lead to up to 2-3X application performance variation in production settings

## Questions:

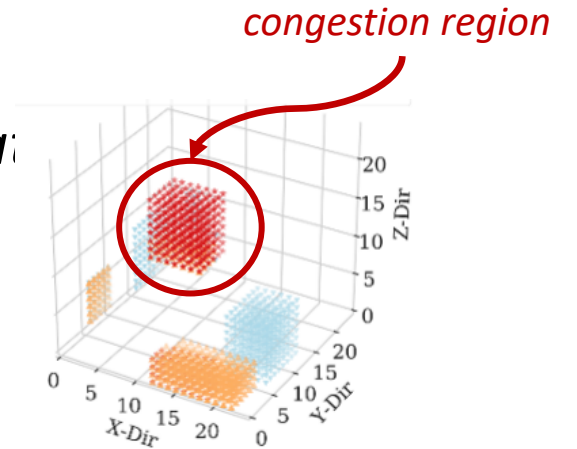
- How often system/applications are experiencing congestion ? [Characterization]
  - What are the culprits behind congestion? [Diagnostics]
- How to avoid and mitigate effects of congestion ? [Network and System Design]

# Highlights

- *Created data mining and ML-driven methodology and associated framework for*
  - Characterizing network design and congestion problems using empirical data
  - Identifying factors leading to the congestion on a live system
  - Checking if the application slowdown was indeed due to congestion
- *Empirical evaluation of a real-world large-scale supercomputer: Blue Waters at NCSA*
  - *Largest 3D Torus network in the world*
  - *5 months* of operational data
  - *815,006* unique application runs
  - *70 PB* of data injected into the network
- Largest dataset on congestion (first on HPC networks)
  - Dataset (*51* downloads and counting!) and code released

# Key Findings

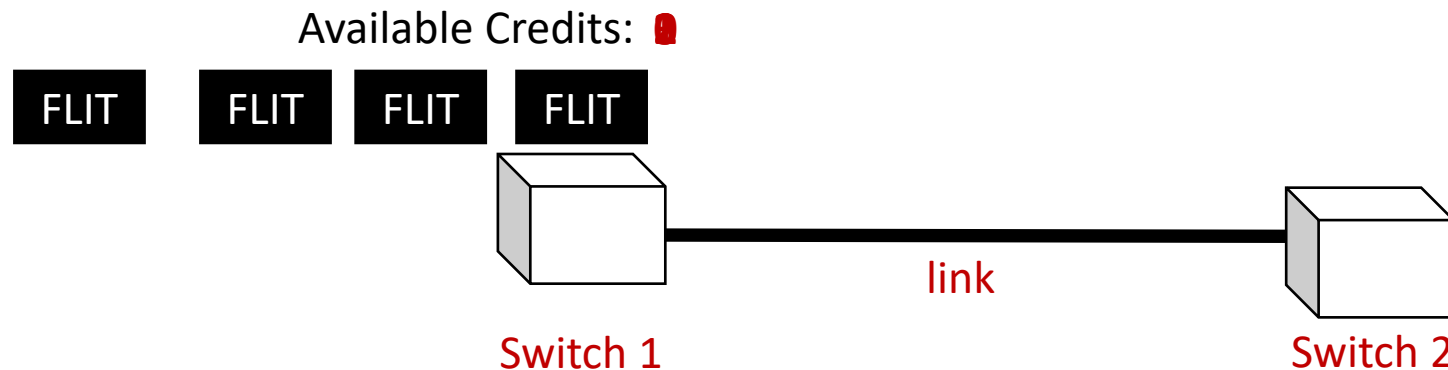
- *HSN congestion is the biggest contributor to app performance variation*
  - *Continuous presence of high congestion regions*
  - *Long lived congestion (may persist for >23 hours)*
- *Default congestion mitigation mechanism have limited efficacy*
  - *Only 8 % (261 of 3390 cases) of high congestion cases found using our framework were detected and acted by default congestion mitigation algorithm*
  - *In ~30% of the cases the default congestion mitigation algorithm was unable to alleviate congestion*
- *Congestion patterns and their tracking enables identification of culprits behind congestion*
  - *critical to system and application performance improvements*
  - *E.g., intra-app congestion can be fixed by changing allocation and mapping strategies*





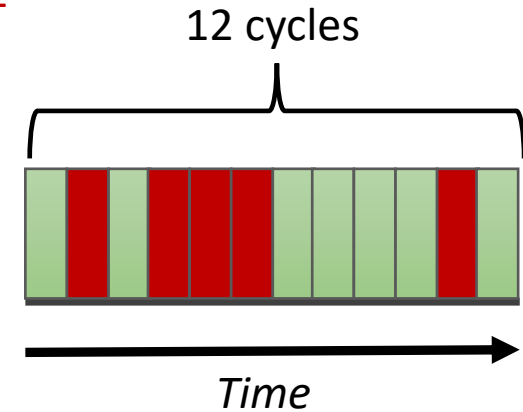
# Congestion in credit-based flow control Network


- Focus on evaluation of credit-based flow control transmission protocol
- Flit is the smallest unit of datum that can be transferred
- Flits are not dropped during congestion
- Backpressure (credits) provides congestion control




If credit  $> 0$ , flit ~~cannot~~ ~~is sent~~

# Measuring Congestion



 Indicates flit waiting (no credit available, allocated buffer full)

 Indicates link is transmitting

Congestion measured using **Percent time stalled** ( $P_{Ts}$ )

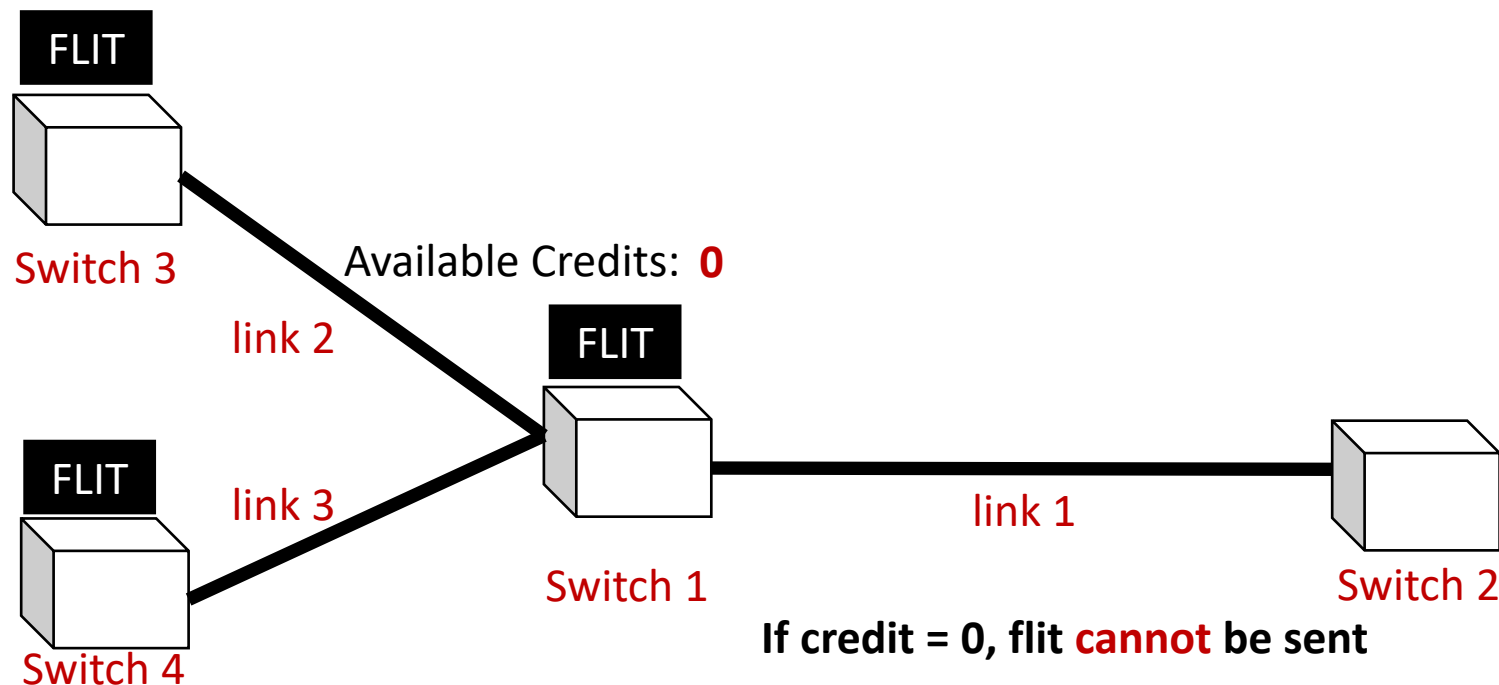
$$P_{Ts}^i = 100 \times \frac{T_s^i}{T^i} = 100 \times \frac{5}{12} = 41.67 \%$$

$T^i$ : # network cycles in  $i^{th}$  measurement interval (fixed value)

$T_s^i$ : # total cycles the link was stalled in  $T^i$  (i.e., flit was ready to be sent but no credits available.)

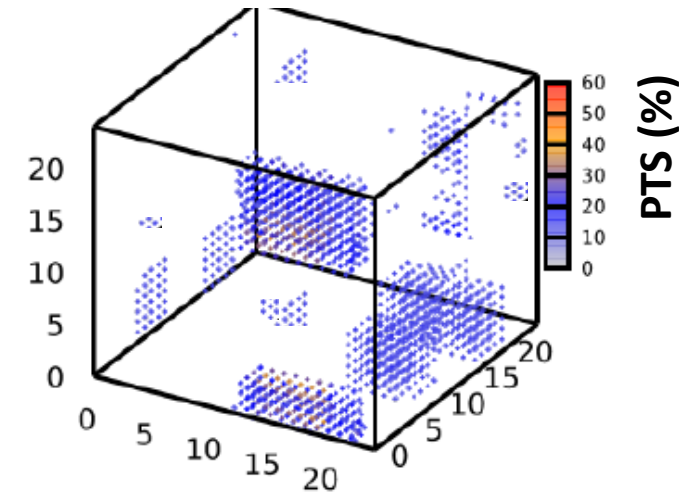
# Congestion in credit-based flow control Network

Insight: Congestion spreads locally (i.e., fans out from an origin point to other senders).

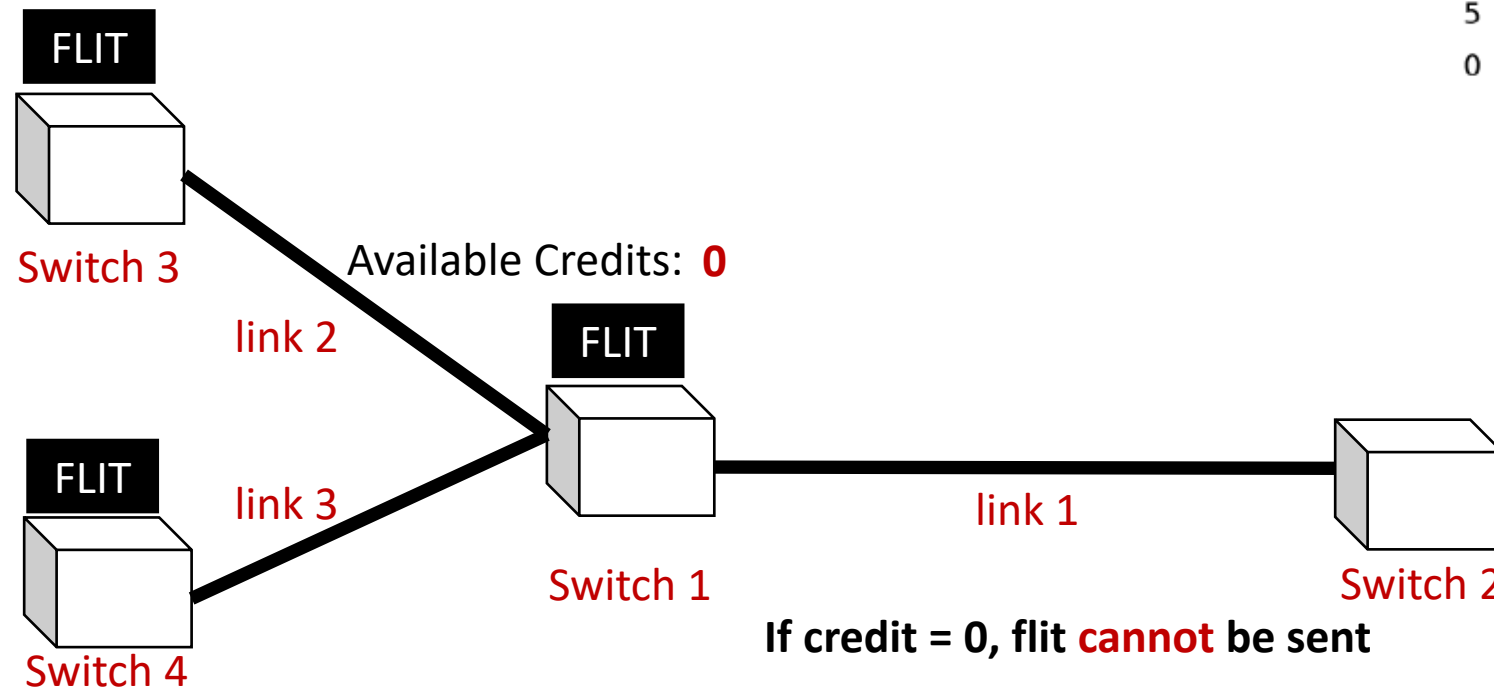


# Congestion in credit-based flow control Network

Insight: Congestion spreads locally (i.e., fans out from an origin point to other senders).

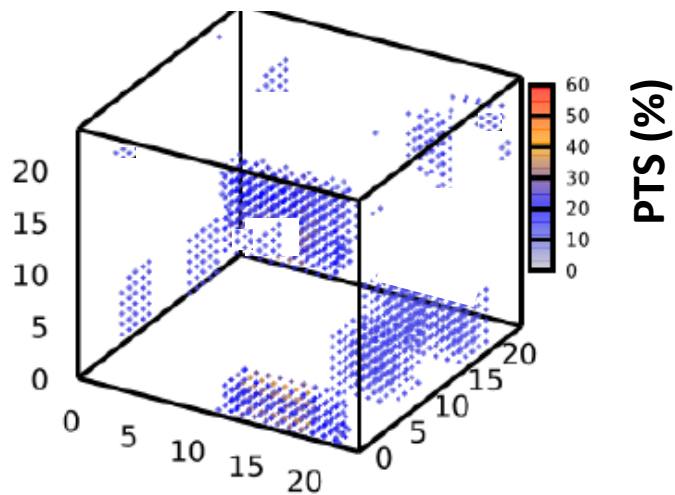


**Congestion Visualization**



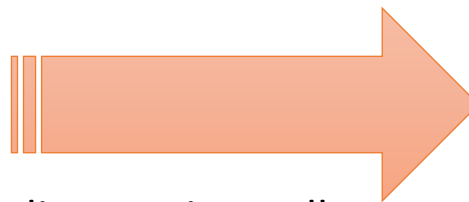
# New unit for measuring congestion

Measure congestion in terms of regions, their size and severity



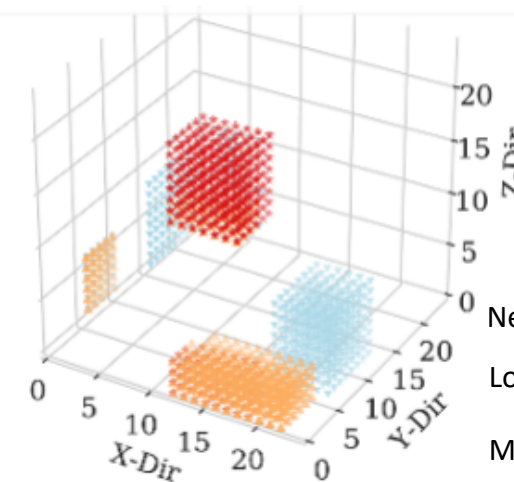
**Raw Congestion Visualization**

*Unsupervised clustering*



distance is small:  
 $d_{\delta}(x,y) \leq \delta$   
 stall difference is small:  
 $d_{\lambda}(x_s - y_s) \leq \theta p$

• Low (1074 links) • High (222 links)  
 ▲ Med. (243 links)



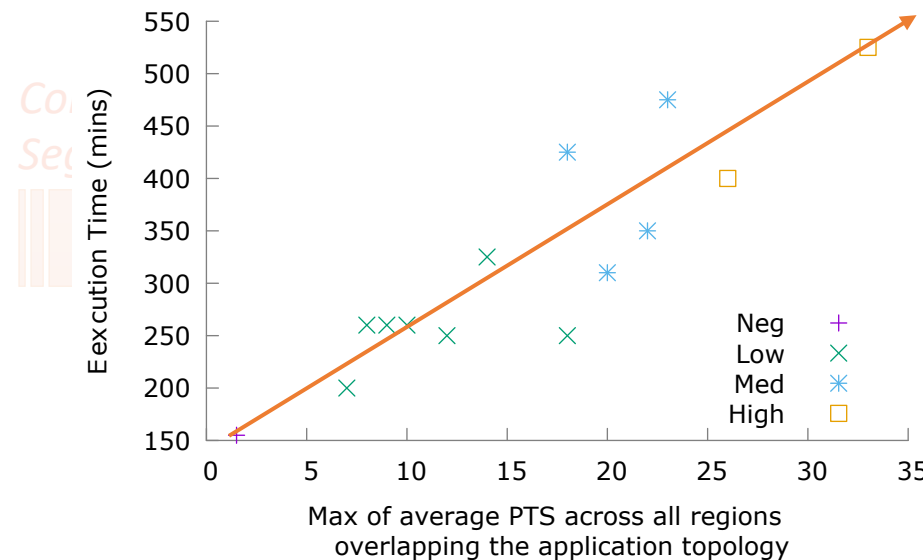
Neg:  $0\% < P_{Ts} \leq 5\%$   
 Low:  $5\% < P_{Ts} \leq 15\%$   
 Med:  $15\% < P_{Ts} \leq 25\%$   
 High:  $25\% < P_{Ts}$

# Congestion Regions Proxy for Performance Evaluation

Congestion Regions (CRs) captures relation between congestion severity and application slowdown and therefore can be used for live forensics and debugging!

(details in paper)

1000-node production molecular dynamics (NAMD) code.

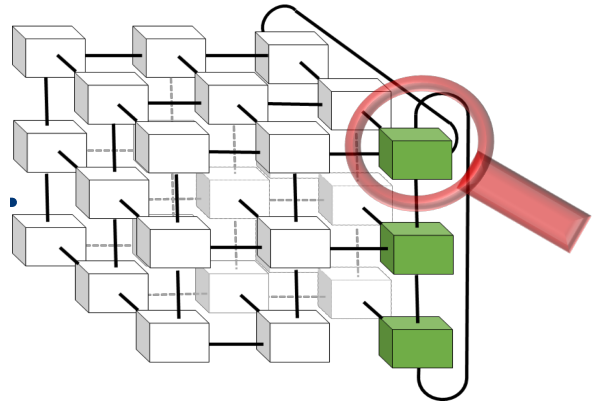


Low:  $5\% < P_{Ts} \leq 15\%$

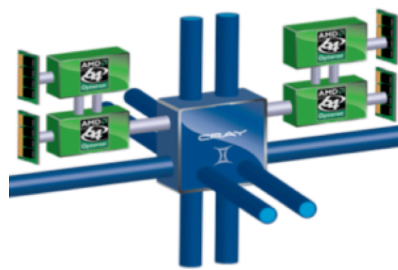
Med:  $15\% < P_{Ts} \leq 25\%$

High:  $25\% < P_{Ts}$

# System, Monitors, and Datasets



**Blue Waters Networks**



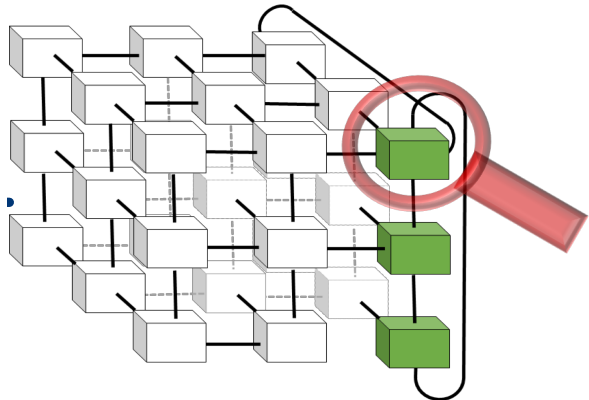
**Cray Gemini Switch**

*Courtesy: Cray Inc. (HP)*

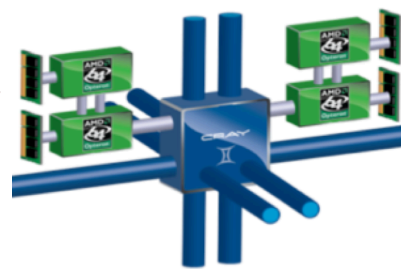
- **Topology:** 3D Torus (24x24x24)
- **Compute nodes :** 28K nodes
- **Avg. Bisection Bandwidth:** 17550 GB/sec
- **Per hop latency:** 105 ns [1]

[1] <https://wiki.alcf.anl.gov/parts/images/2/2c/Gemini-whitepaper.pdf>

# System, Monitors, and Datasets



**Blue Waters Networks**

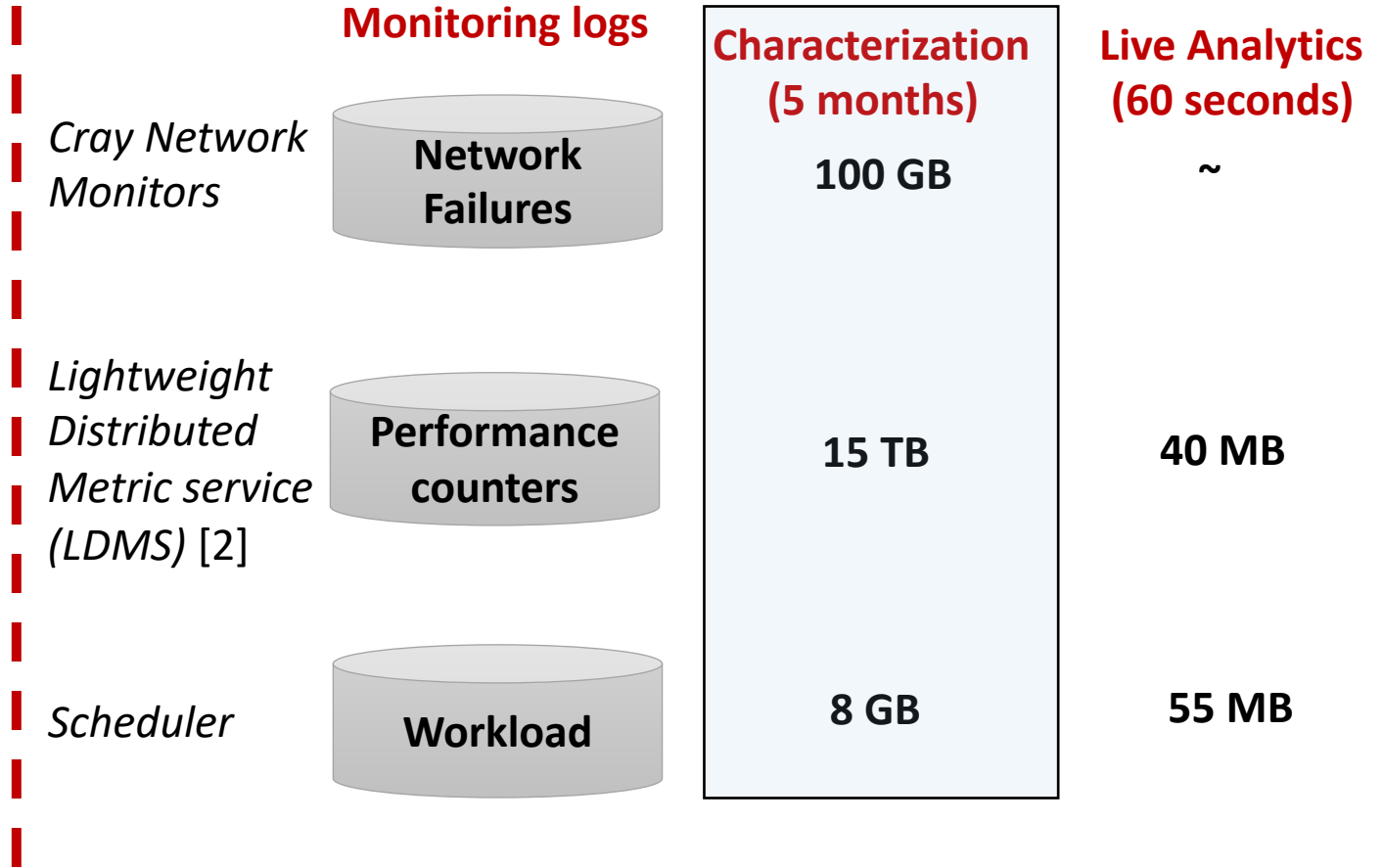


**Cray Gemini Switch**

*Courtesy: Cray Inc. (HP)*

- **Topology:** 3D Torus (24x24x24)
- **Compute nodes :** 28K nodes
- **Avg. Bisection Bandwidth:** 17550 GB/sec
- **Per hop latency:** 105 ns [1]

[1] <https://wiki.alcf.anl.gov/parts/images/2/2c/Gemini-whitepaper.pdf>



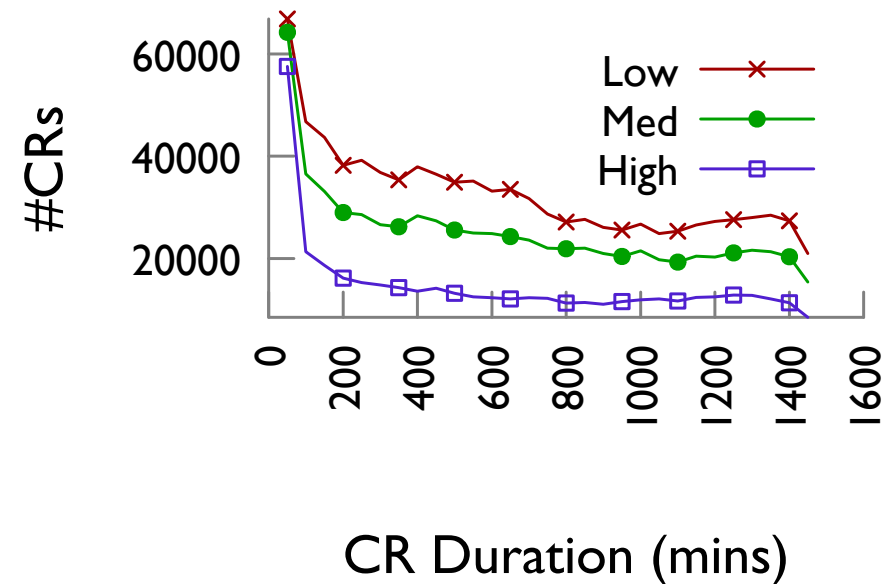
[2] A. Agelastos et al. Lightweight Distributed Metric Service: A Scalable Infrastructure for Continuous Monitoring of Large-scale Computing Systems and Applications. In *SC14: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 154–165, 2014.



# 1. Congestion is the biggest contributor to app performance variation

## Long-lived congestion

- *Congestion Region* can persist up to ~24 hours (median: 9.7 hours)
- *Congestion Region* count decreases with increasing duration



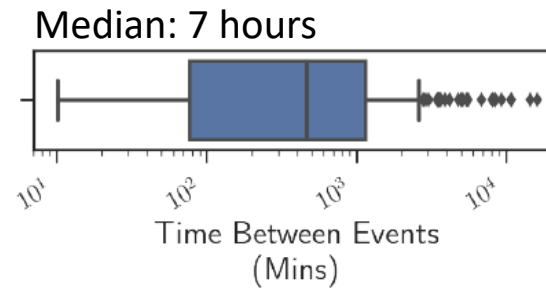
Low:  $5\% < P_{Ts} \leq 15\%$   
 Med:  $15\% < P_{Ts} \leq 25\%$   
 High:  $25\% < P_{Ts}$

## 2. Limited efficacy of default congestion detection and mitigation mechanisms

- #congestion mitigating triggered : **261**
- Median time between events: **7 hours**
- Failed to alleviate congestion in **29.8% cases**

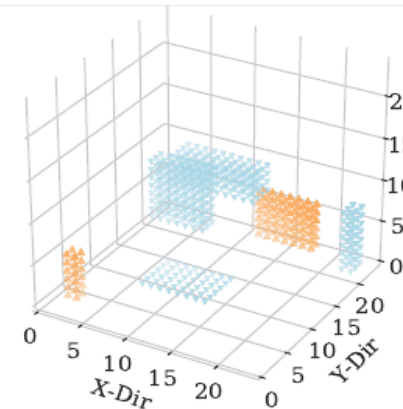
Default mitigation throttles all NICs such that

aggregate traffic injection bandwidth across all nodes < single node bandwidth ejection



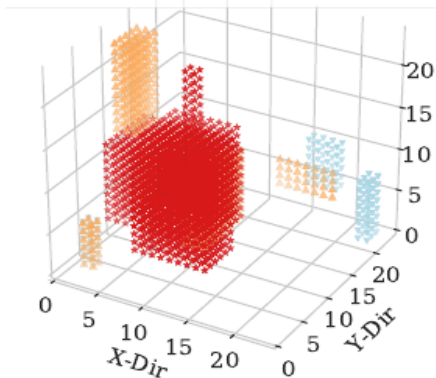
Default system congestion detection and mitigation

Low (347 links) Med. (77 links)



Before congestion mitigation

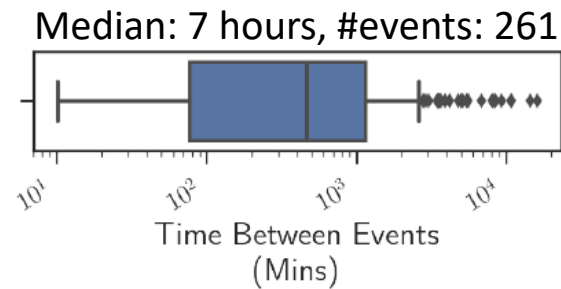
Low (174 links) High (969 links) Med. (150 links)



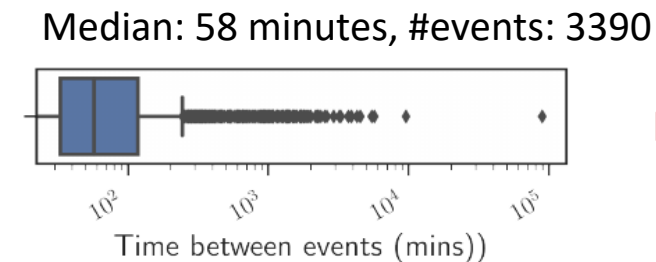
After congestion mitigation

## 2. Limited efficacy of default congestion detection and mitigation algorithms

- *Default congestion mitigating triggered : 261*
- *Median time between events: 7 hours*
- *Failed to alleviate congestion in 29.8% of the cases*



Default system congestion detection and mitigation

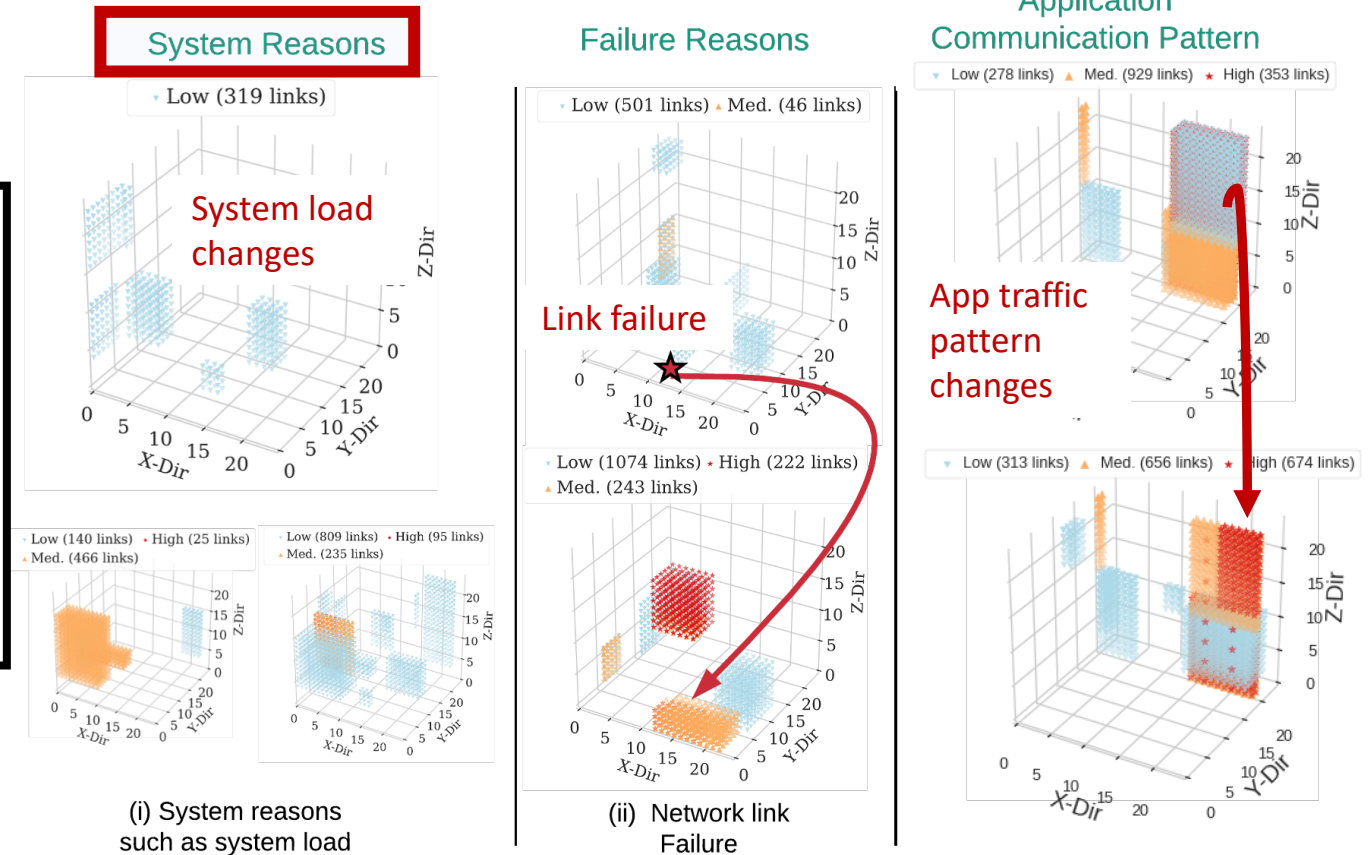


Monet detection

*Only 8 % (261 of 3390 cases) of high congestion cases found by Monet were detected and acted by default congestion mitigation algorithm*

## 3. Congestion patterns and their tracking enables identification of culprits behind congestion

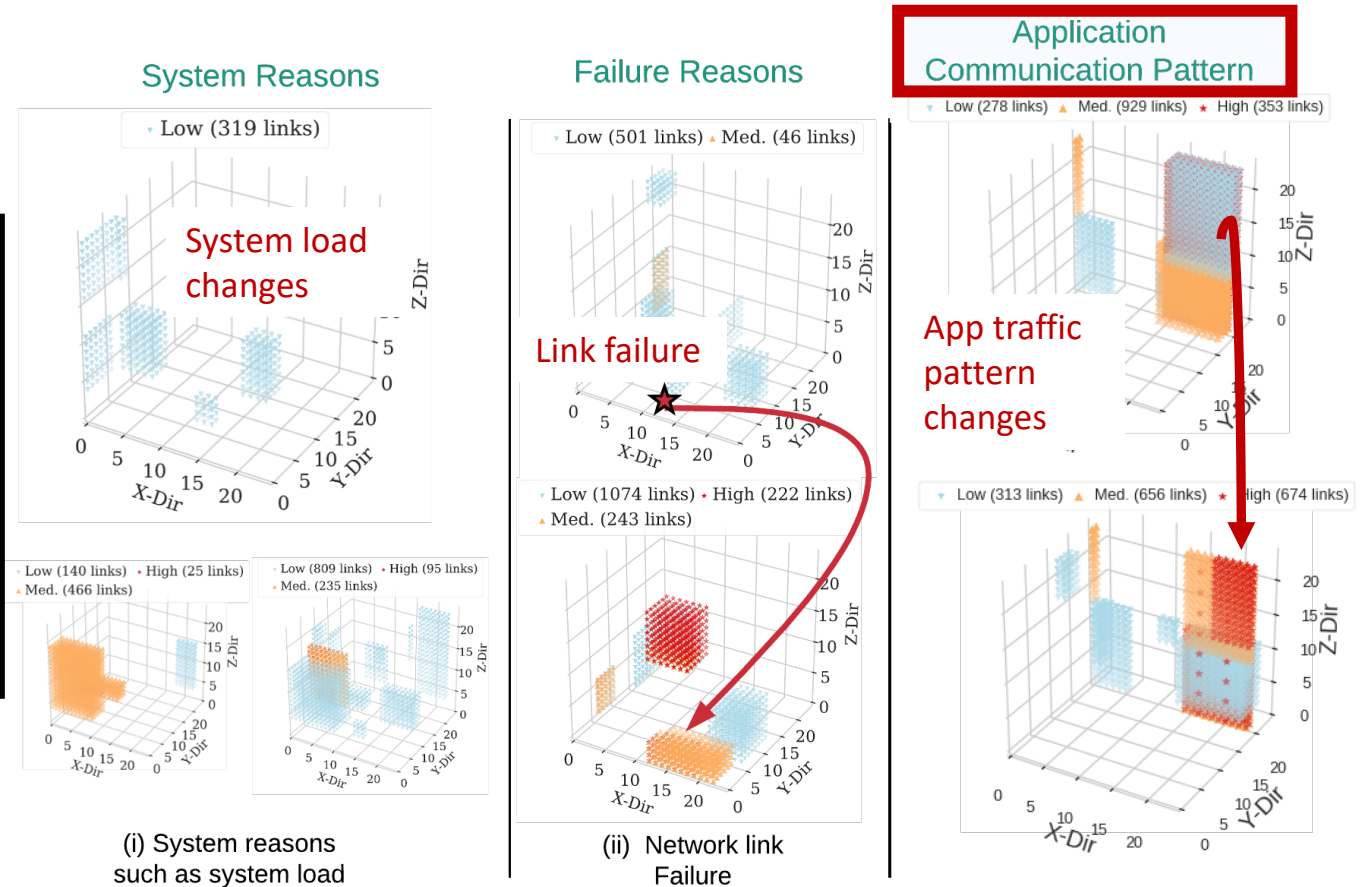
- Network design and congestion-aware scheduling
- E.g., topology-aware scheduling [1] improved system throughput by **56%** by tuning resource allocation strategies



[1] J Enos et al. Topology-aware job scheduling strategies for torus networks. In *Proc. Cray User Group*, 2014.

## 3. Congestion patterns and their tracking enables identification of culprits behind congestion

- Node mapping within the allocation reduces intra-app congestion
- E.g., TopoMapping [2] for finding optimal process rank mapping for the allocated resource



[2] Galvez et al. Automatic topology mapping of diverse large-scale parallel applications. In *Proceedings of the International Conference on Supercomputing, ICS '17*, pages 17:1–17:10, New York, NY, USA, 2017. ACM.

# Conclusion

- Developed and validated the proposed methodology on production datasets
- Code and dataset online (51 downloads and counting!)
  - <https://databank.illinois.edu/datasets/IDB-2921318>
  - <https://github.com/CSLDepend/monet>

# Future Work

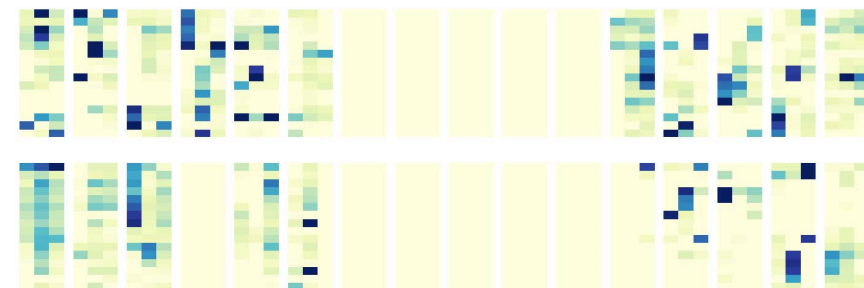
Congestion avoidance and mitigation is an ongoing problem !

Developing workload-aware high-speed networks

- Inferring and meeting application demands
- Optimizing congestion control and routing strategies

**Meet us at the poster session!**

Wednesday 6:30 PM - 8:00 PM  
Cypress Room



**Congestion Visualization on a production Cray Aries (DragonFly Network)**



**Questions?**