# Devesh Tiwari

Assistant Professor
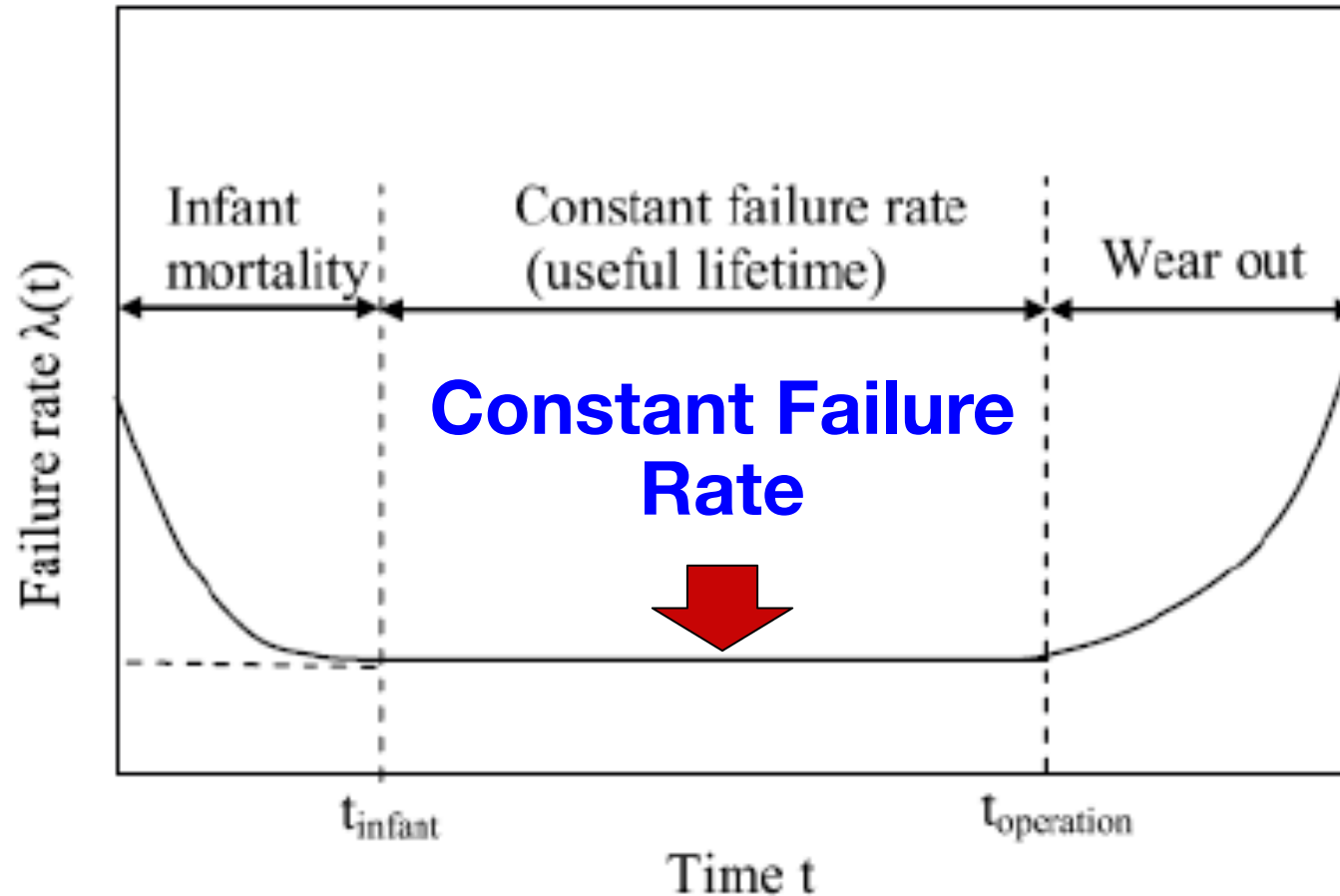Northeastern University
tiwari@northeastern.edu

Research Talk at IFIP WG 10.4 Meeting
Longmont, CO June 2017

In collaboration with Saurabh Gupta, Evgenia Smirni, Christian Engelmann, Sudharshan Vazhkudai, Franck Cappello, Jim Rogers et al.

# Useful Fake News

## Devesh Tiwari

Assistant Professor
Northeastern University
tiwari@northeastern.edu

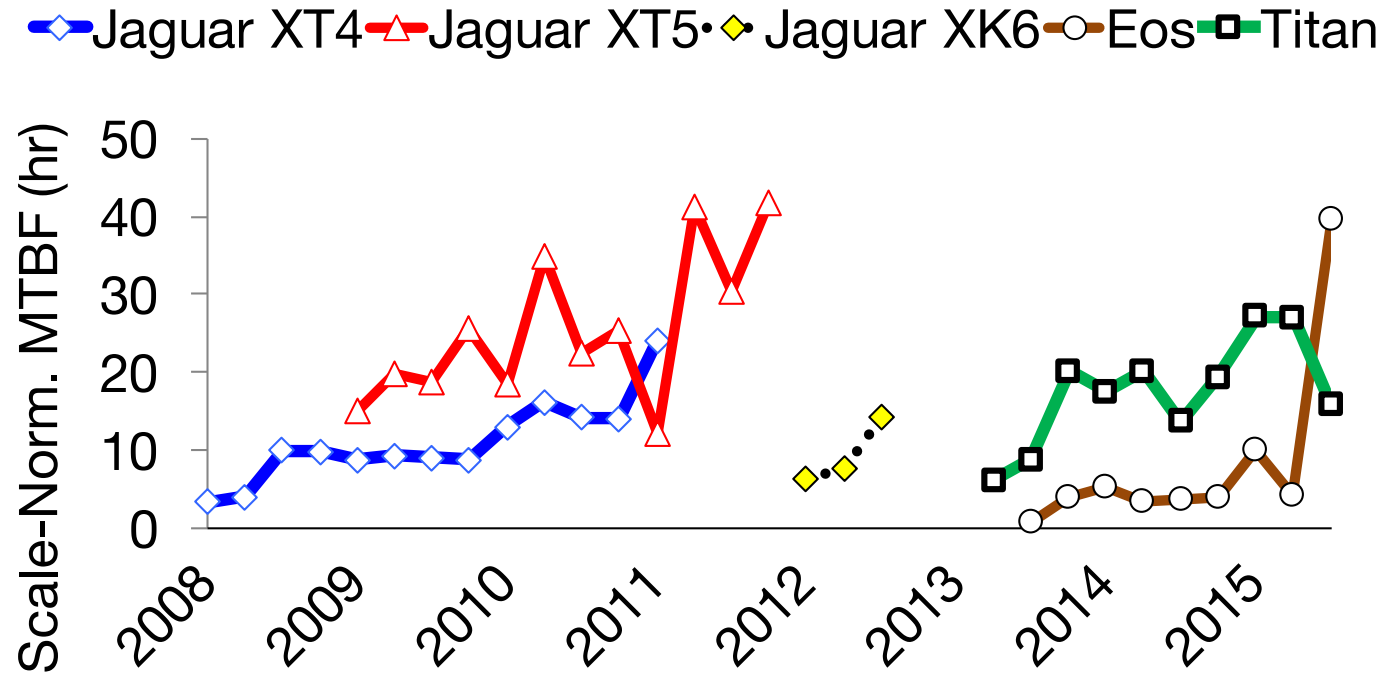Research Talk at IFIP WG 10.4 Meeting
Longmont, CO June 2017

In collaboration with Saurabh Gupta, Evgenia Smirni, Christian Engelmann, Sudharshan Vazhkudai, Franck Cappello, Jim Rogers et al.

# Bathtub Curve



**Constant Failure Rate**

**Uniformly random temporal and spatial distribution**
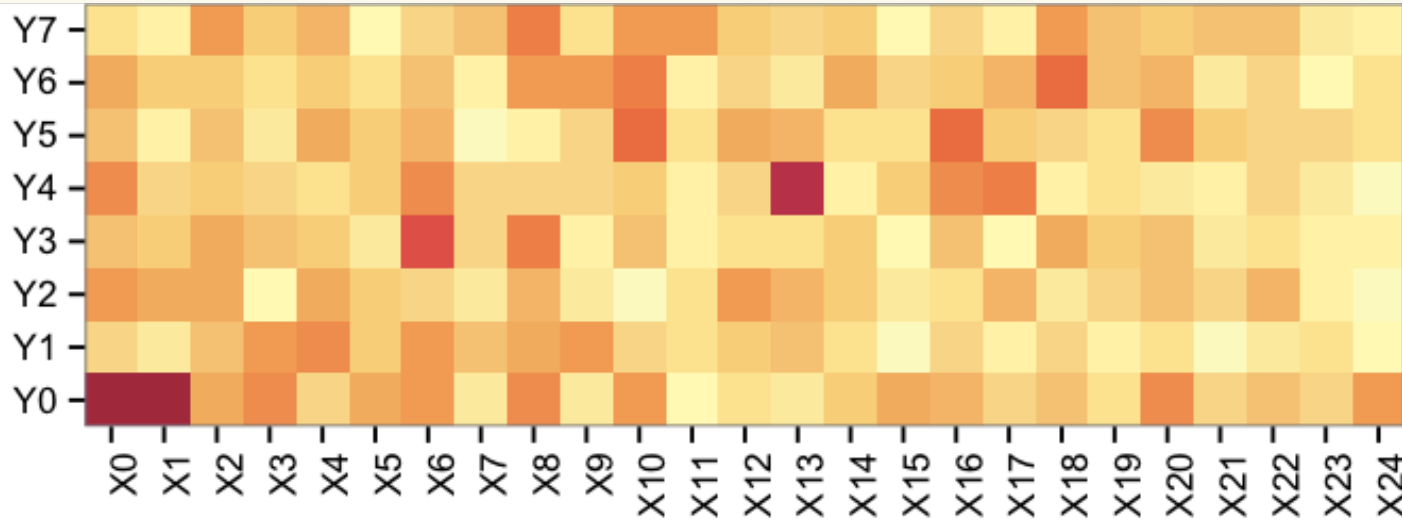
# Bathtub Curve



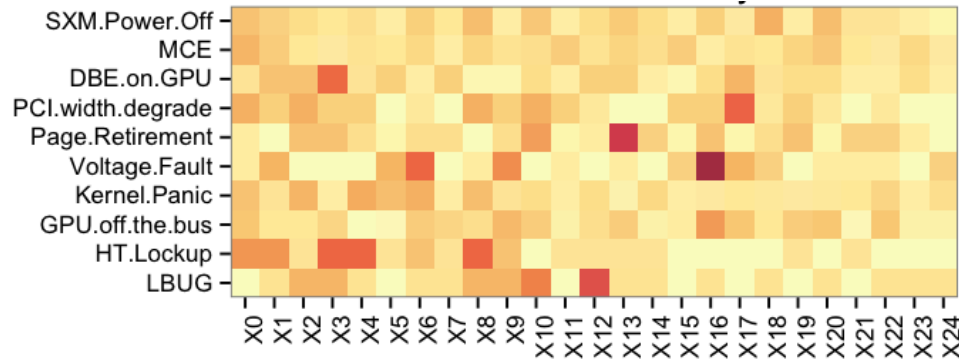**Up to 4x variation in MTBF during useful lifetime**
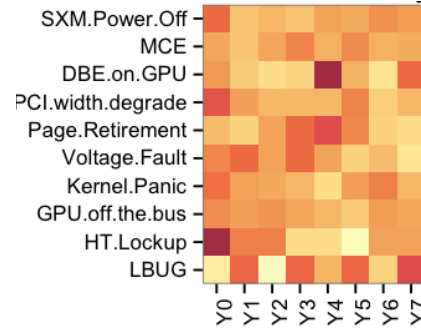
# Spatial Distribution of System Failures
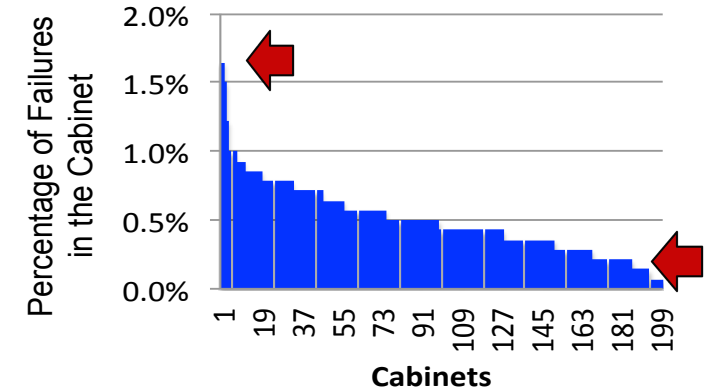


200 Cabinets

Titan supercomputer cabinet columns

600 Cages

Failure type

Cabinet columns

Cabinet rows

**System failures are not uniformly randomly distributed in space.**
**This holds true for individual failure types, different time windows, spatial granularity.**

Double Bit Errors

ECC Page Retirement Errors

Off the Bus Errors

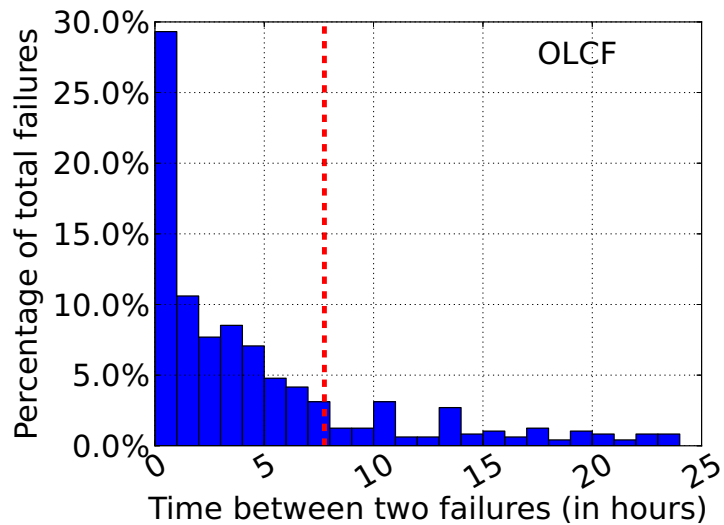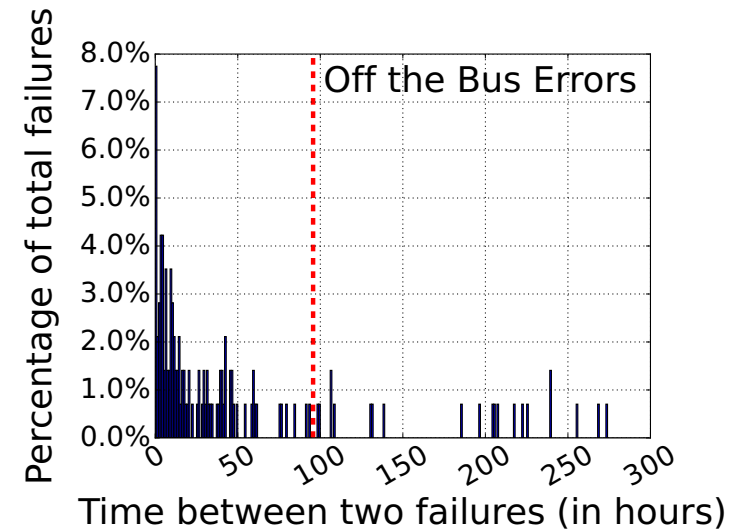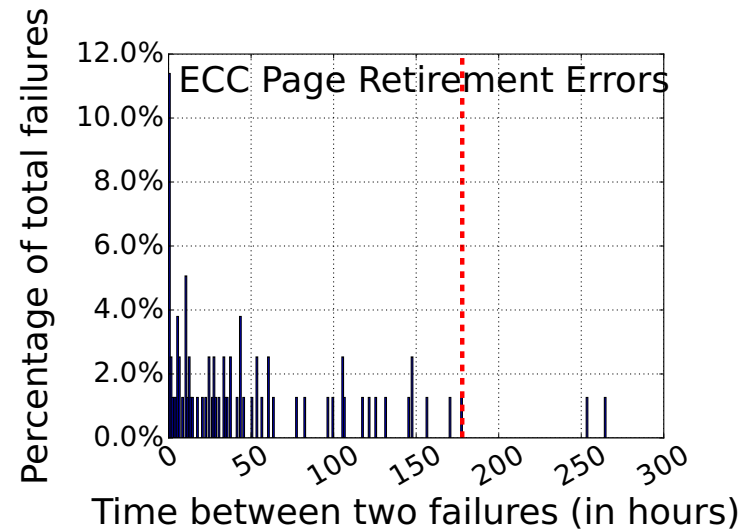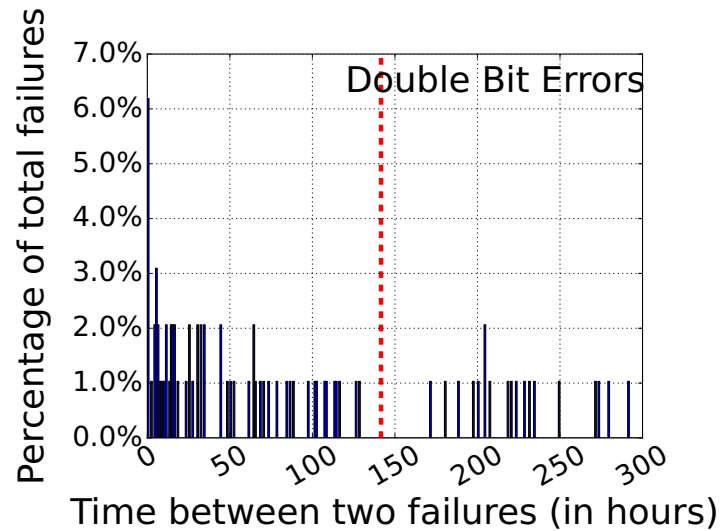System failures exhibit strong temporal locality (not same as the traditional early bath-tub curve).

GPU core utilization or variance in utilization does not necessarily correlate with soft-errors, application and users do!

[HPCA 2016] A Large-Scale Study of Soft-Errors on GPUs in the Field

# Quarantine Job Scheduling



**% Failures Avoided**   **% Quarantine Hours**

Quarantine time duration 48 hours

| | % Failures Avoided | % Quarantine Hours |
|---|---|---|
| Node | 3.85% | 0.02% |
| Blade | 5.07% | 0.09% |
| Cage | **7.21%** | **0.69%** |
| Cabinet | 9.64% | 2.04% |

Percentage of Node-hours used by debug jobs

Mean 1.4%

0.69%

**Significant fraction of failures can be avoided from interrupting production applications**

**Debug or non-production jobs can be scheduled on quarantine nodes**

[DSN 2015] Understanding and Exploiting Spatial Properties of System Failures on Extreme-Scale HPC Systems

# Lazy Checkpointing



performance gain $= \beta e^{-(\frac{t_3}{\lambda})^k}$

performance loss $= (\alpha_{max-oci} - \alpha_{oci})(e^{-(\frac{t_2}{\lambda})^k} - e^{-(\frac{t_4}{\lambda})^k})$
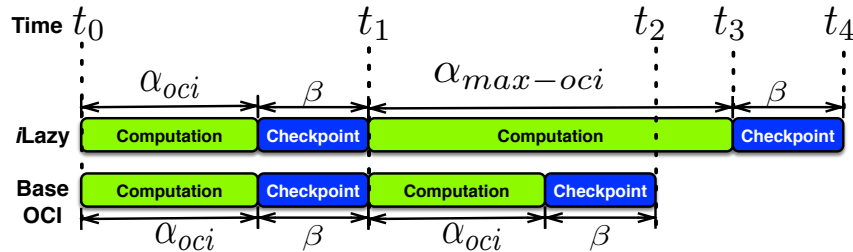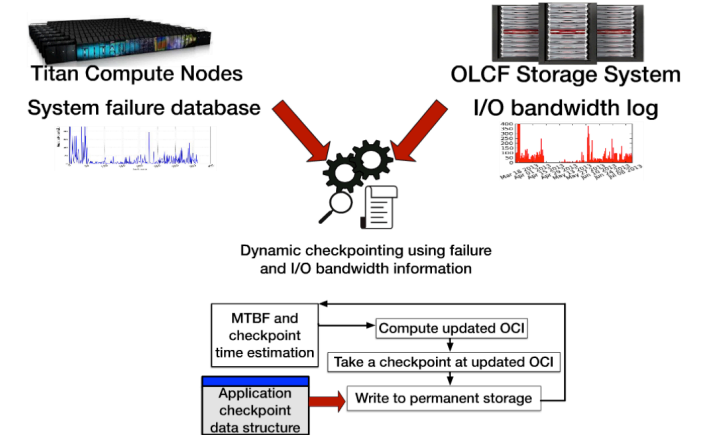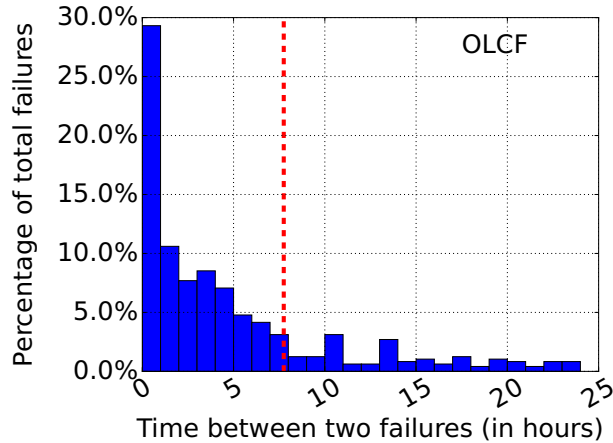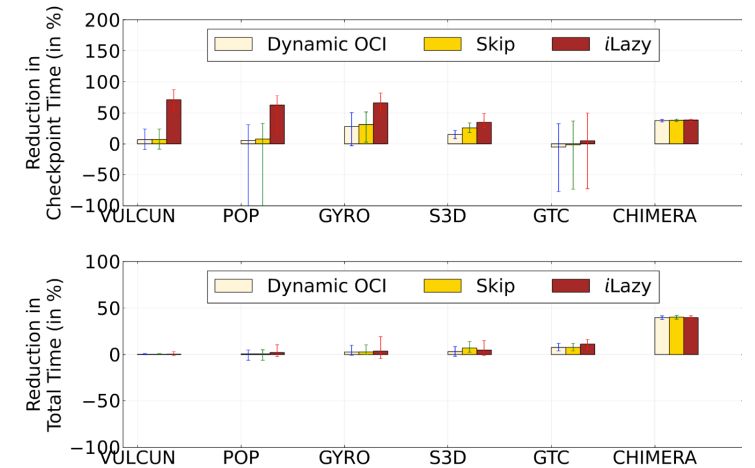
$= (\alpha_{max-oci} - \alpha_{oci})(e^{-(\frac{2(\alpha_{oci}+\beta)}{\lambda})^k} - e^{-(\frac{\alpha_{max-oci}+\alpha_{oci}+2\beta}{\lambda})^k})$

$\beta e^{-(\frac{\alpha_{max-oci}+\alpha_{oci}+\beta}{\lambda})^k} = (\alpha_{max-oci} - \alpha_{oci})e^{-(\frac{2(\alpha_{oci}+\beta)}{\lambda})^k}$

$\quad -(\alpha_{max-oci} - \alpha_{oci})e^{-(\frac{\alpha_{max-oci}+\alpha_{oci}+2\beta}{\lambda})^k}$

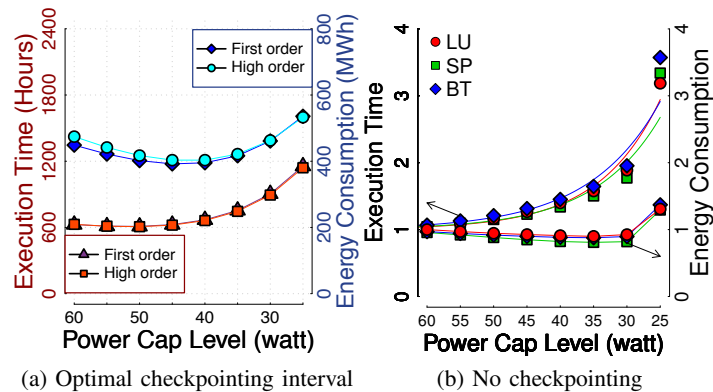## Being Lazy is good (but with some bounds!)

[DSN 2014] Lazy Checkpointing: Exploiting Temporal Locality in Failures to Mitigate
Checkpointing Overheads on Extreme-Scale Systems

# Optimal Checkpointing under Power-Constraints



(a) Xeon E5-2670

(b) Xeon E5-2630

**Under power-capping optimal checkpointing interval is quite different than traditional optimal checkpointing interval**

(a) Optimal checkpointing interval

(b) No checkpointing

**Optimal power capping level is different than no-checkpointing case.**

[DSN 2016] Power-capping Aware Checkpointing: On the Interplay among Power-capping, Temperature, Reliability, Performance, and Energy

# Thanks!

## Devesh Tiwari

tiwari@northeastern.edu

Research Talk at IFIP WG 10.4 Meeting
Longmont, CO June 2017