

Oak Ridge National Laboratory

Computing and Computational Sciences Directorate

High-Performance Data Flows Using Analytical Models and Measurements

Nagi Rao, Qiang Liu
Oak Ridge National Laboratory

Don Towsley, Gayane Vardoyan
University of Massachusetts, Amherst

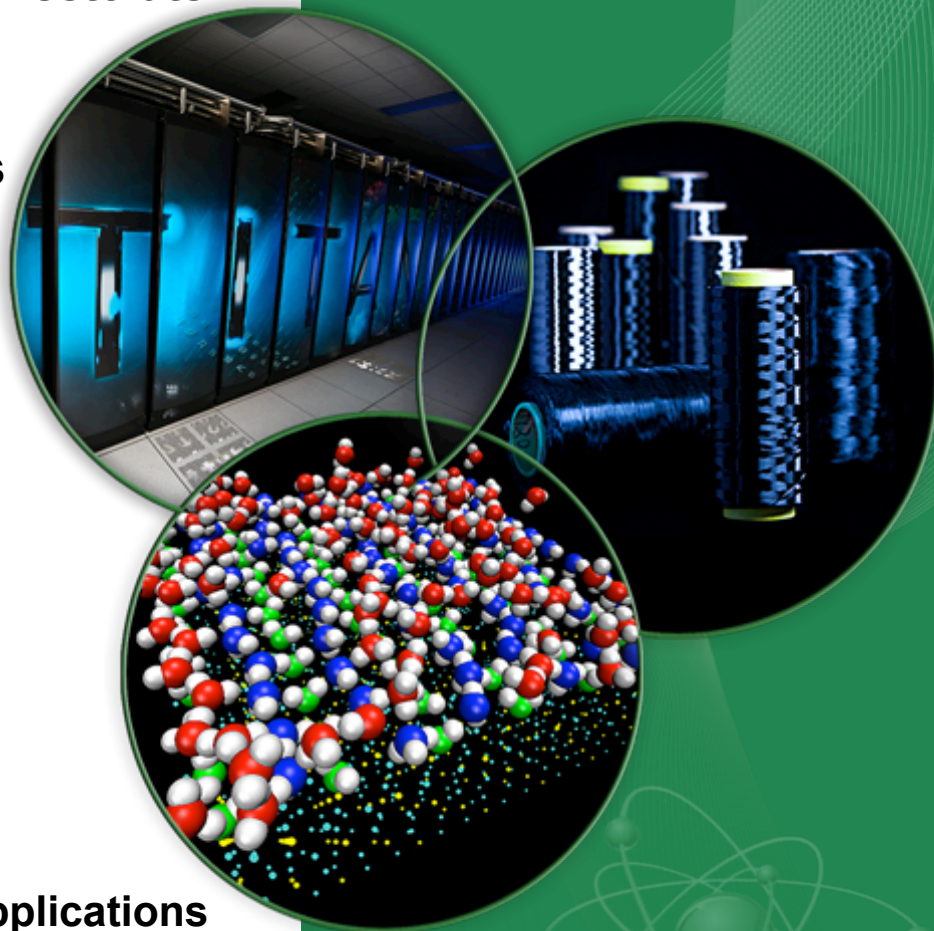
Raj Kettimuthu, Ian Foster
Argonne National Laboratory

Brad Settlemyer
Los Alamos National Laboratory

Workshop on
Modeling & Simulation of Systems and Applications
August 10-12, 2016, Seattle, Washington

Sponsored by
RAMSES Project, ASCR, DOE

ORNL is managed by UT-Battelle
for the US Department of Energy



Outline of Presentation

- **Introduction**
- **Throughput Profiles and Traces**
 - TCP, UDT measurements
- **Throughput Model**
 - Monotonicity and concavity
- **File Transfer Rates**
 - Lustre and xfs
- **Conclusions**

Dedicated Network Connections: Increasing Deployments

Dedicated connections becoming more available

- DOE OSCARS provides ESnet WAN connections
- Google SDN dedicated networks

Desirable Features

- dedicated capacity - no competing traffic
- low loss rates – no induced losses from “other” traffic
- no need for “graceful degradation” to accommodate other traffic

Expectations for data transport methods

- peak throughput easier to achieve using “simple” flow control
- easier to tune parameters – due to predictable and simple dynamics

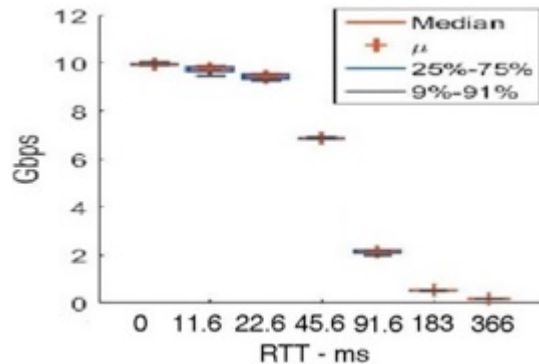
Scenarios:

- memory-to-memory, disk-to-disk, memory-to/from-disk transfers:
 - high bandwidth file transfers
 - data transfers between remote computations
- monitoring and control channels: low loss and jitter requirements
 - computational monitoring and steering
 - remote experimentation – computation controlled

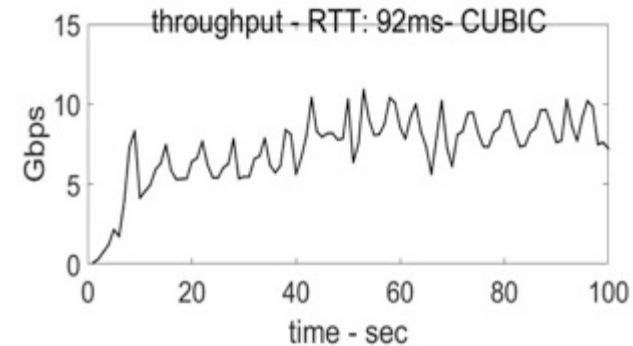
Data Transfers Over Dedicated Network Connections

Performance of data transfer methods

throughput profile: RTT



throughput trace: time

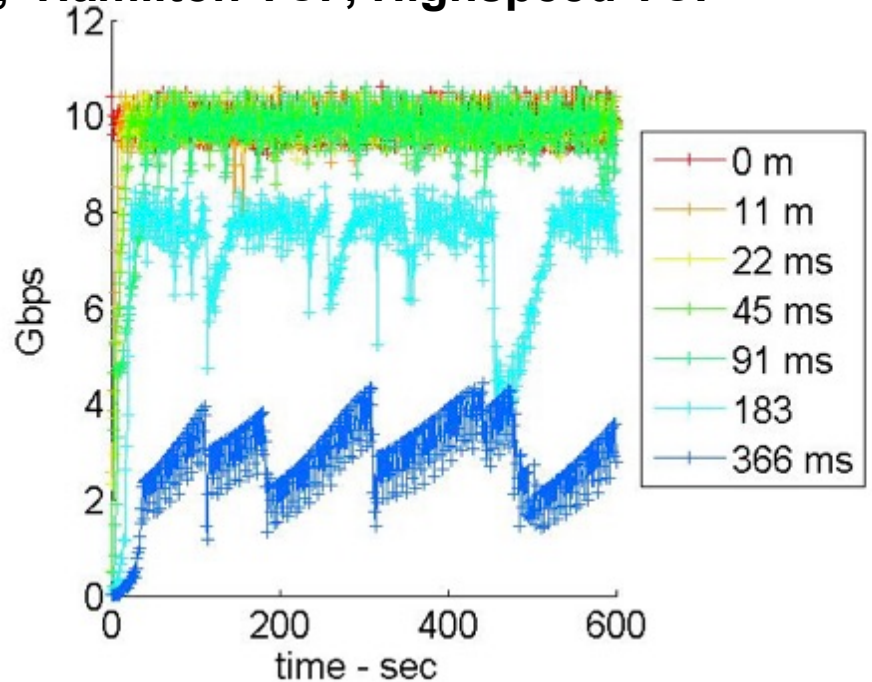
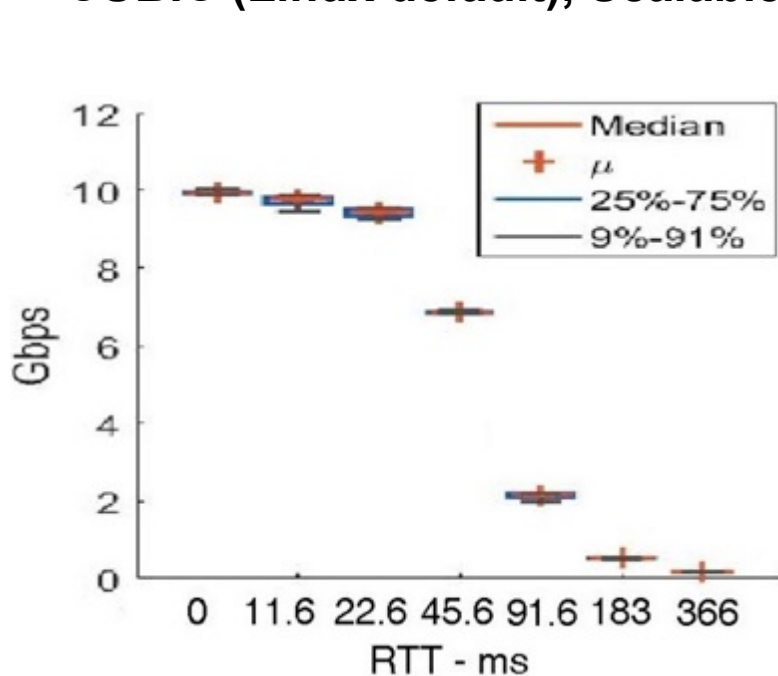


Network Transport Methods

- **TCP** – widely deployed, including over dedicated connections
 - mechanism: slow-start followed by congestion avoidance
 - expected performance:
 - convex throughput profiles
 - slow-start followed by periodic trajectories
- **UDT** – UDP-based, particularly well-suited for dedicated connections
 - mechanism: ramp-up followed by stable flow rate
 - expected performance
 - flat throughput profile
 - ramp-up followed by constant trajectory
- **ASPERA** – commercial UDP-based transport method

TCP Memory-to-Memory Throughput Measurements

Throughput profiles and traces: qualitatively similar across TCP variants
CUBIC (Linux default), Scalable TCP, Hamilton TCP, Highspeed TCP



As expected:

- profile: decreases with RTT;
- trace: sort of periodic in time

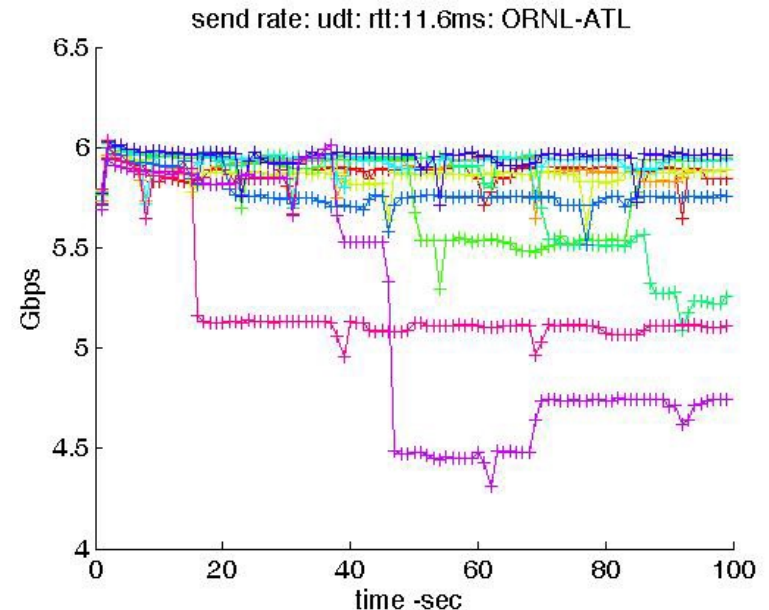
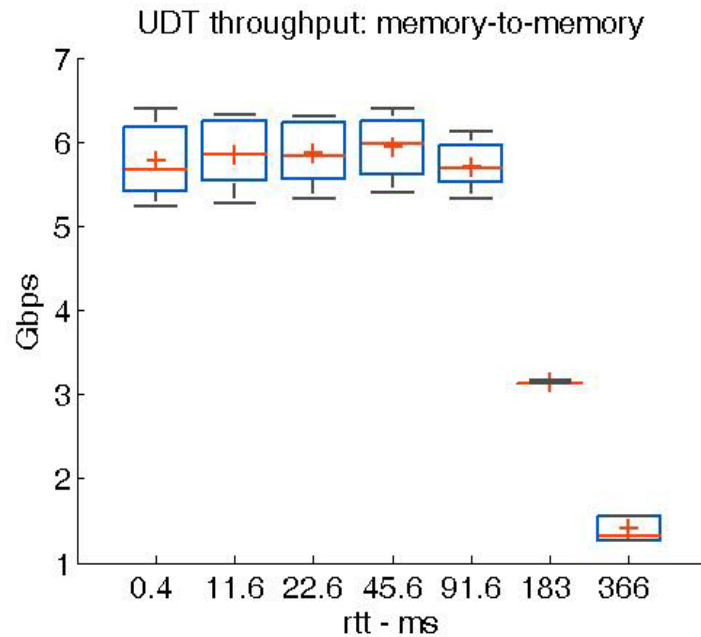
Additional characteristics:

- profile: concave at lower RTT
- trace:
 - significant variations
 - larger variation at higher RTT

UDT Memory-to-Memory Throughput Measurements

Implements flow control and loss recovery over UDP

- application-level – particularly suited for dedicated connections



Analytical models indicate:

- profile: flat with RTT
- trace: smooth rise and constant

Measured characteristics:

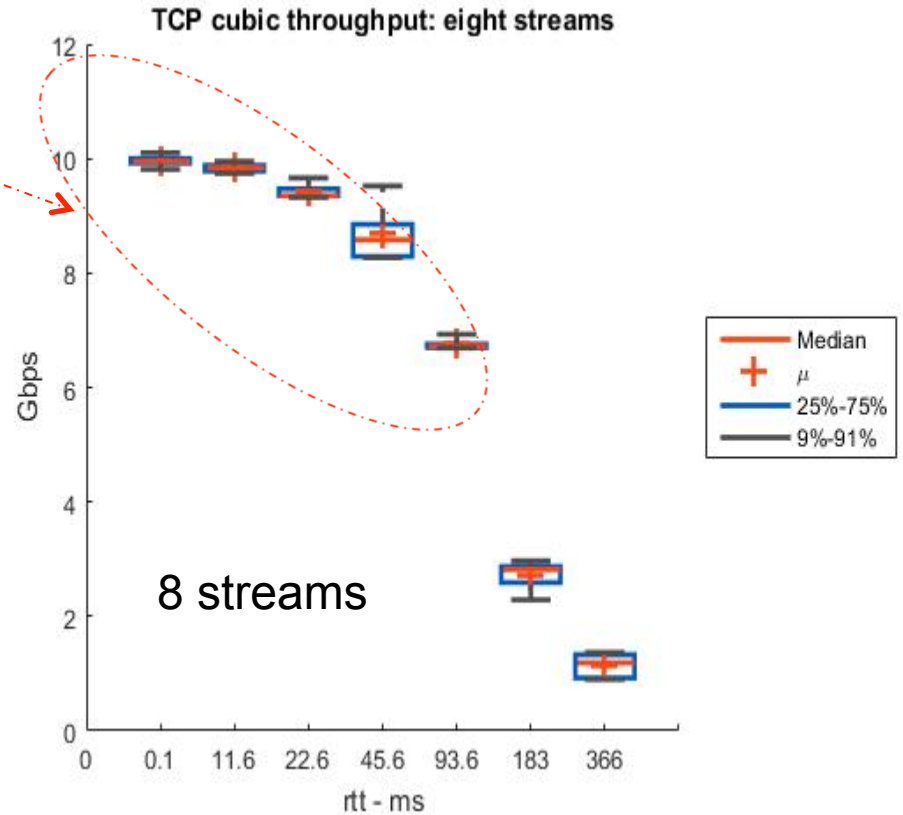
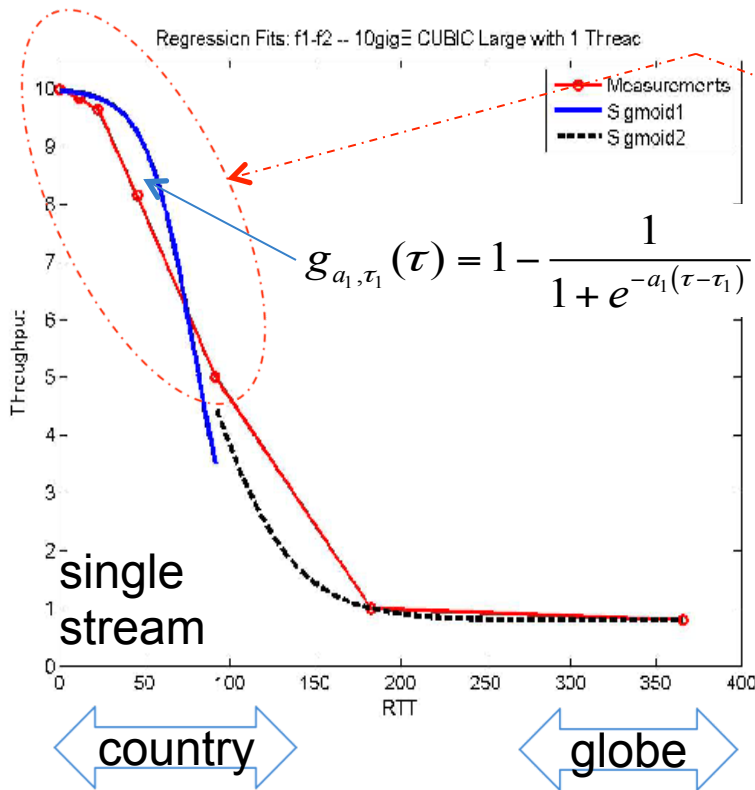
- profile: overall decrease with RTT
- trace:
 - significant variations: same RTT
 - repeated – significant drops

TCP Profiles: memory to memory transfer

10Gbps dedicated connections: bohr03 and bohr04

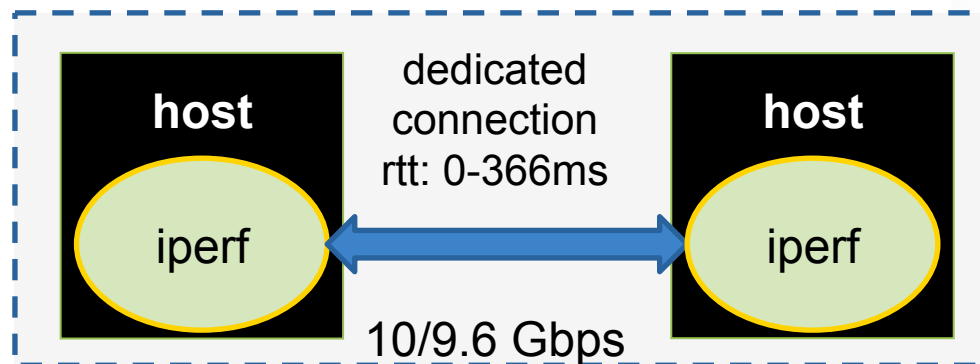
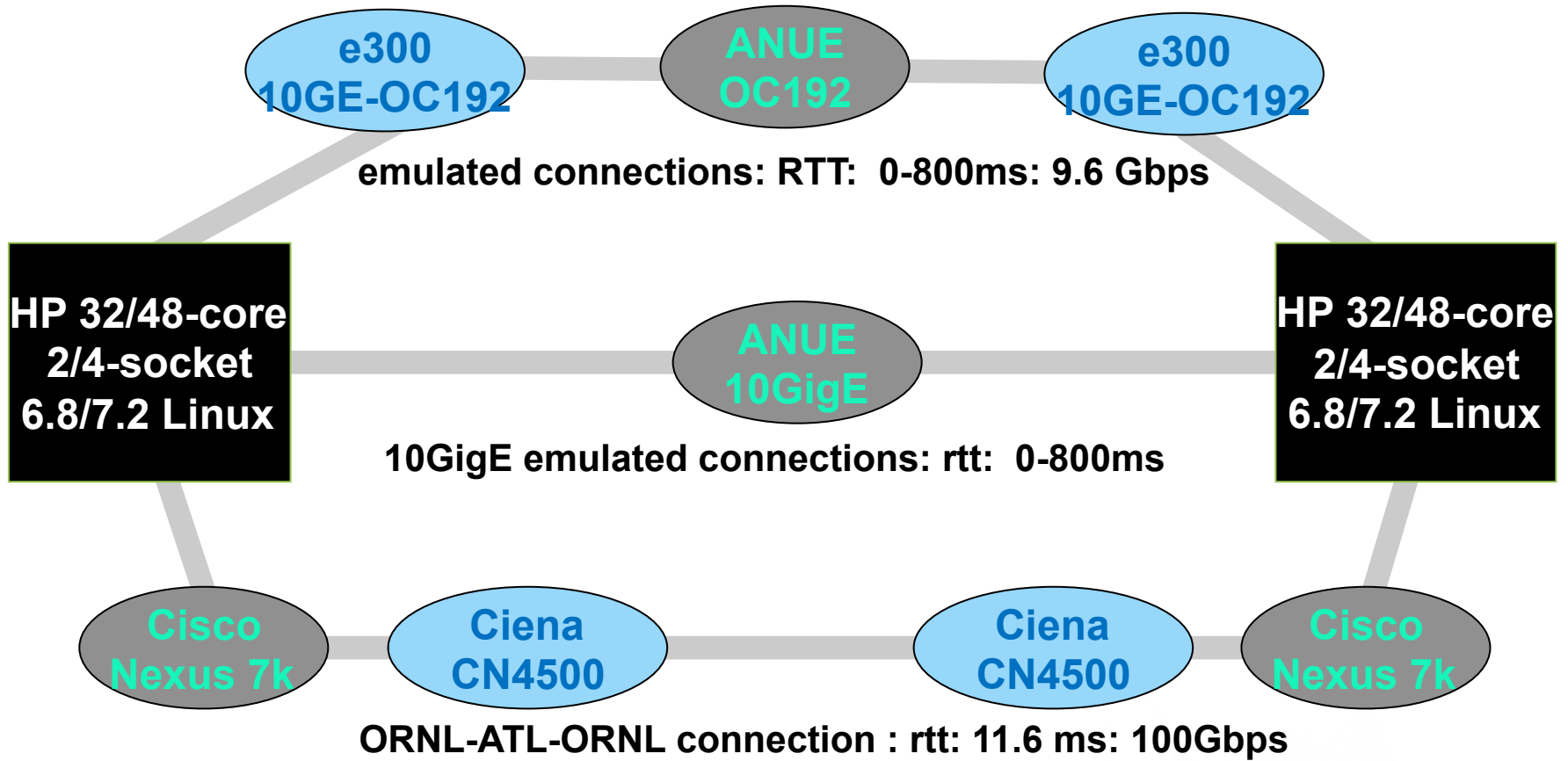
- CUBIC congestion control module- default under Linux
- TCP buffers tuned for 200ms rtt: 1-10 parallel streams

highly-desirable
concave region



RTT: cross-country (0-100ms), cross-continentals (100-200ms), across globe(366ms)

ORNL Network Testbed (since 2004)



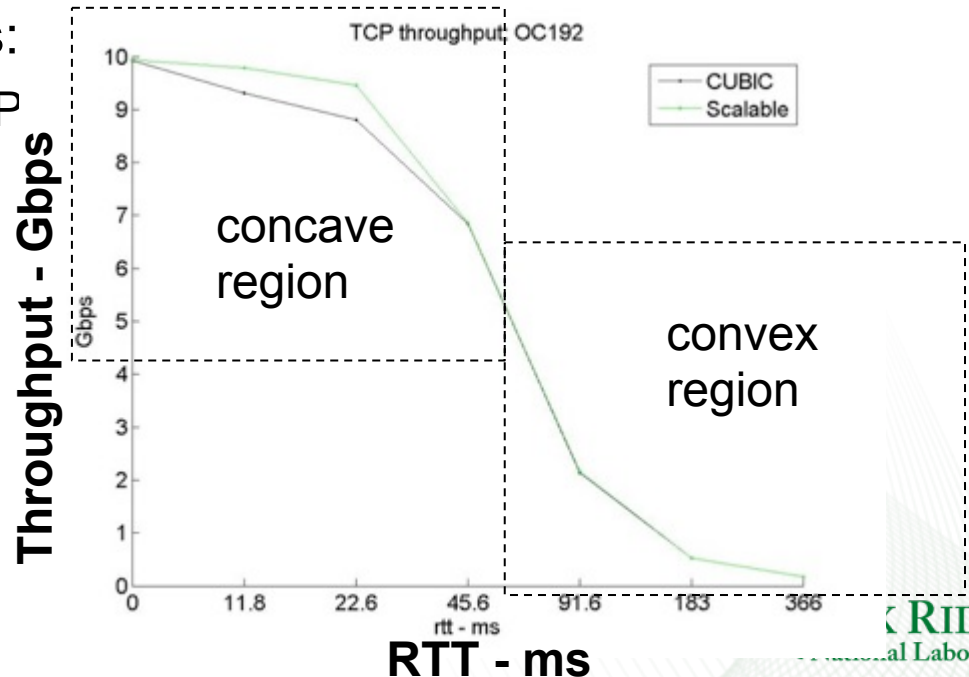
TCP Throughput Profiles

- Most common TCP throughput profile
 - convex function of rtt
 - example, Mathis et al (1997)

Throughput at rtt τ loss-rate p

$$T_M(\tau) = \frac{MSS * k}{\tau \sqrt{p}}$$

- Observed Dual-mode profiles:
 - CUBIC, STCP, HTCP, HS-TCP
- Smaller RTT
- Concave region
- Larger RTT
- Convex region



Desired Features of Concave Region

- Concave regions is very desirable
 - throughput does not decay as fast
 - rate of decrease slows down as rtt

Function $T(\tau), \tau \in I$ is concave iff derivative $\frac{dT}{d\tau}$ is non-increasing

not satisfied by Mathis model:

$$\frac{dT_M}{d\tau} = -\frac{C}{\tau^2 \sqrt{p}}$$

- Measurements: throughput profiles for rtt: 0-366ms

- Concavity: small rtt region

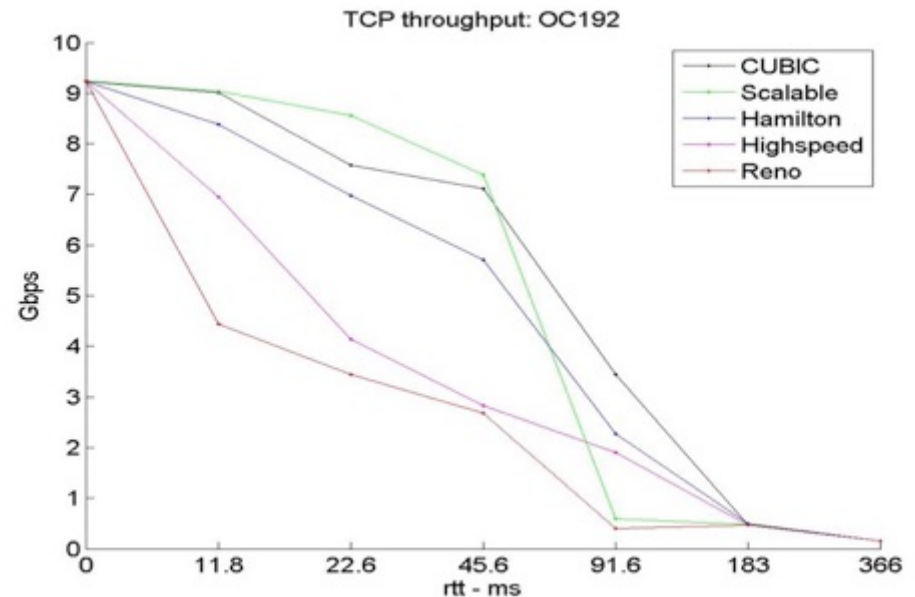
Only some TCP versions:

- CUBIC,
- Hamilton TCP
- Scalable TCP

Not for some TCP versions:

- Reno

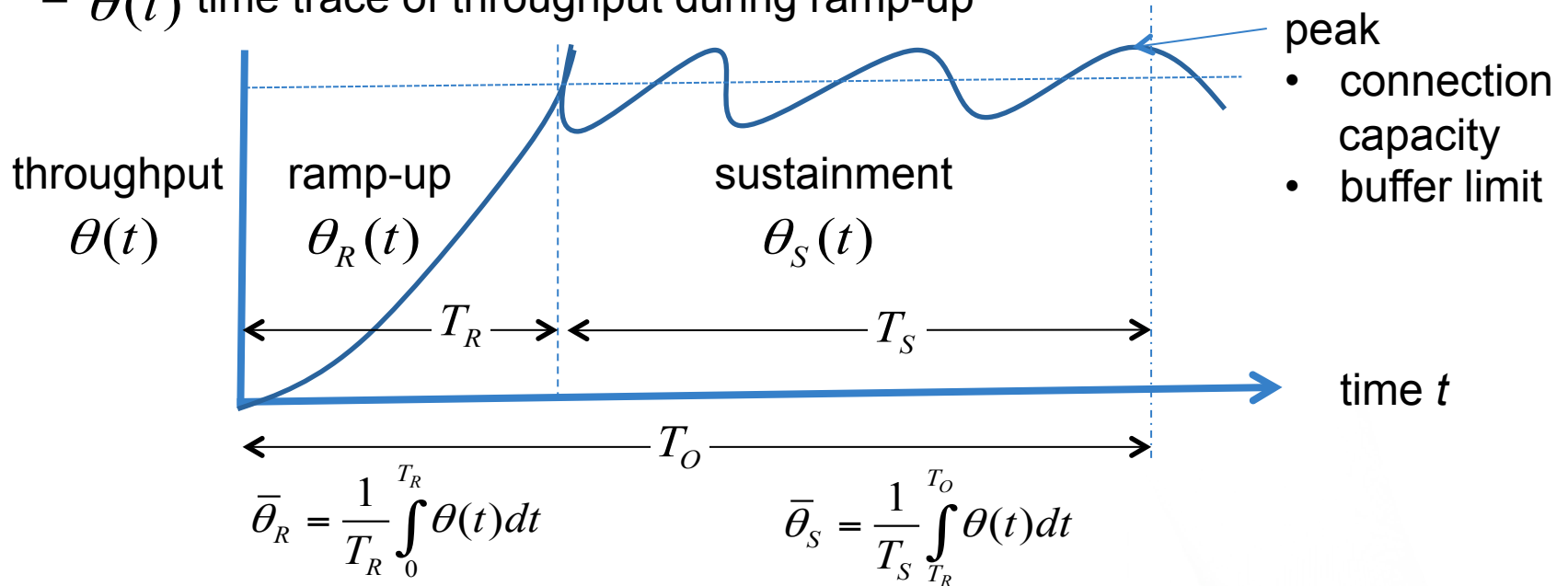
- These are loadable linux modules



Generic Transport Model

- Ramp-up Phase

- throughput increases from initial value to around capacity
e.g. TCP slow start, UDT ramp-up
- $\theta(t)$ time trace of throughput during ramp-up

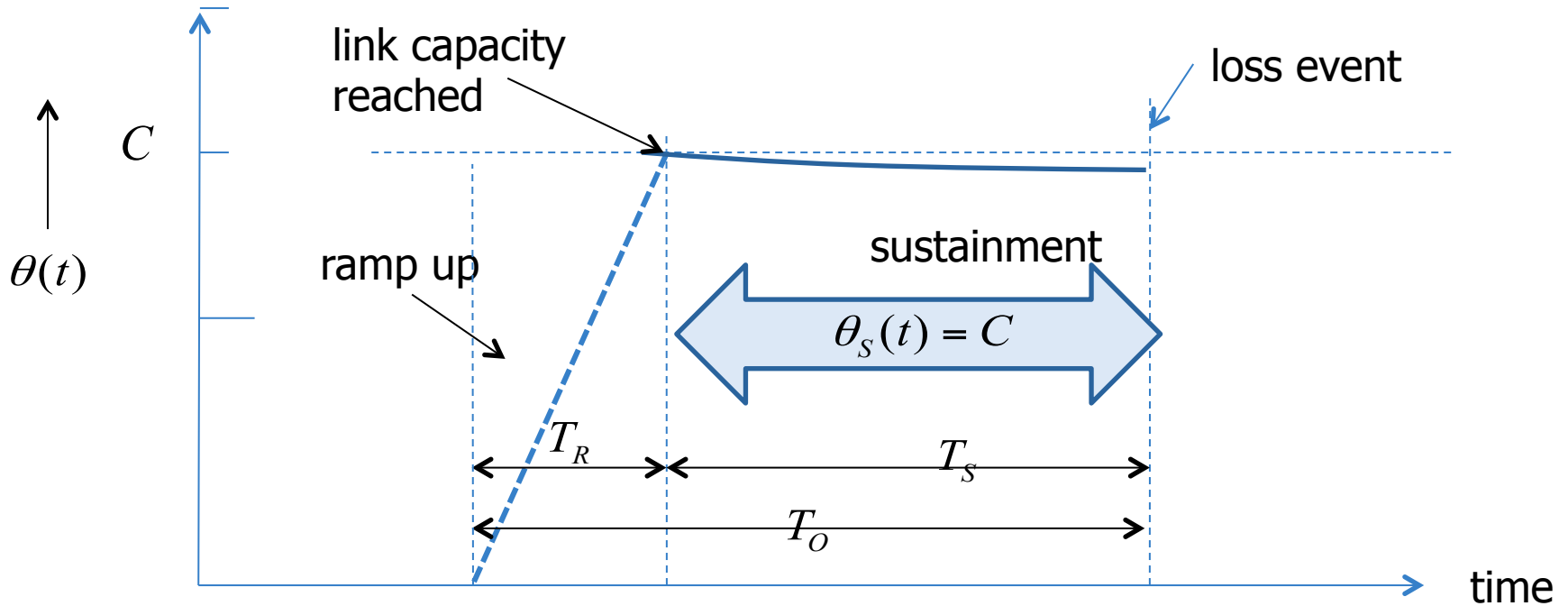


- Sustainment Phase

- Throughput is maintained around a peak value
 - TCP congestion avoidance, UDT stable peak
- $\theta(t)$ time trace of throughput during sustainment

$$\Theta_O(\tau) = \frac{1}{T_O} \int_0^{T_O} \theta(\tau, t) dt$$

Boundary Case



Average Throughput: Monotonicity

$$\Theta_O = \bar{\theta}_R \frac{T_R}{T_O} + \bar{\theta}_S \left[\frac{T_O - T_R}{T_O} \right]$$

$$= \bar{\theta}_R f_R + \bar{\theta}_S [1 - f_R]$$

$$= \bar{\theta}_S - f_R \left[\underbrace{\bar{\theta}_S - \bar{\theta}_R}_{+} \right]$$

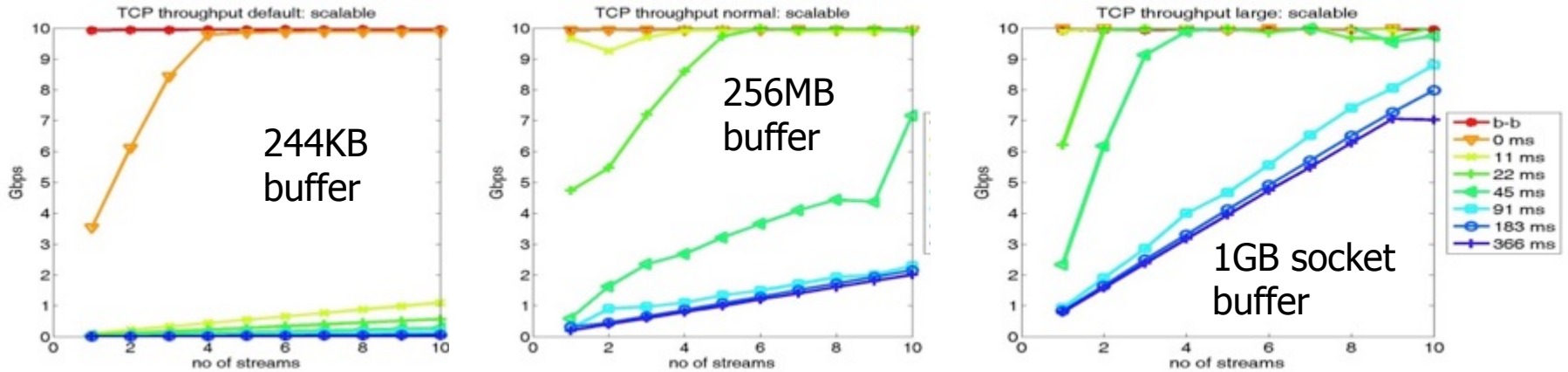
increases with τ

decreases with τ

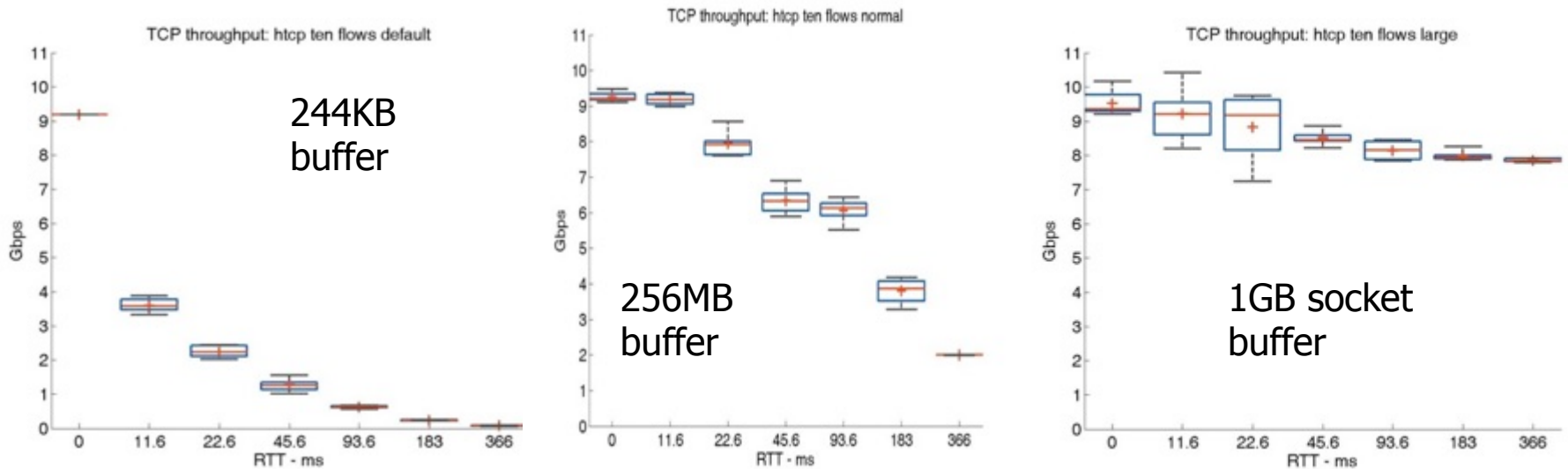
Multiple TCP Streams and Large Buffers

Both provide higher throughput and expanded concave region

- Increase in average throughput: STCP over 10GigE



- Expanded Concave Region: Hamilton TCP: 10 flows over OC192



Concavity: Faster than Slow Start – Multiple TCP flows

Faster than Slow Start:

More increases than slow start: $n_k = \tau^{\epsilon_\tau} \log C$ $\epsilon_\tau > 0, \tau > 1$

$$T_R = \tau n_k = \tau^{1+\epsilon_\tau} \log C$$

$$\text{data sent: } 1 + 2 + L + 2^{n_k} = 2^{n_k+1} - 1 = 2^{1+\tau^{\epsilon_\tau}} C - 1$$

$$\bar{\theta}_R \approx \frac{2^{1+\tau^{\epsilon_\tau}} C}{\tau^{1+\epsilon_\tau} \log C}$$

Average Throughput:

$$\Theta_o(\tau) = \frac{2^{1+\tau^{\epsilon_\tau}} C}{T_o} + C \left[\frac{T_o - \tau^{1+\epsilon_\tau} \log C}{T_o} \right]$$

$$\frac{d\Theta_o}{d\tau} = - \underbrace{\frac{(1+\epsilon_\tau) \tau^{\epsilon_\tau} C \log C}{T_o}}_{\text{decreasing function of } \tau}$$

decreasing function of τ

implies concavity of $\Theta_o(\tau)$

Monotonicity Conditions: Not always decreasing in τ

Average Throughput:

$$\begin{aligned}\Theta_O(\tau) &= \bar{\theta}_R(\tau) \frac{T_R(\tau)}{T_O} + \bar{\theta}_S(\tau) \left[\frac{T_O - T_R(\tau)}{T_O} \right] \\ &= \bar{\theta}_S(\tau) + f_R(\tau) [\bar{\theta}_R(\tau) - \bar{\theta}_S(\tau)]\end{aligned}$$

Effective sustained phase: $\bar{\theta}_S(\tau) > \bar{\theta}_R(\tau)$

$\Theta_O(\tau)$ monotonically decreases in τ

$$\bar{\theta}_S(\tau) \downarrow \qquad f_R(\tau) \uparrow$$

Ineffective sustained phase: $\bar{\theta}_S(\tau) < \bar{\theta}_R(\tau)$ ←

$\Theta_O(\tau)$ may increase in τ

may lead to:

- lower throughput
- convex region

If $f_R(\tau)$ decreases “faster” than $\bar{\theta}_S(\tau)$

Some UDT measurements show this behavior

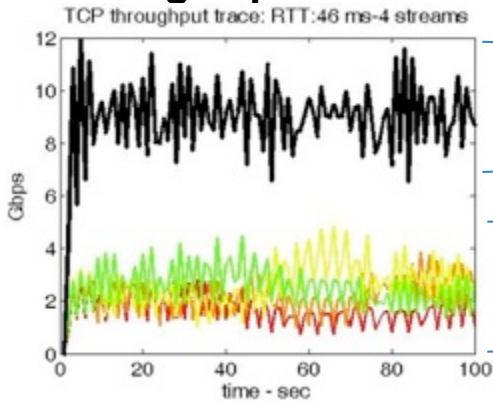
Poincare Map

Well-Known tool for analyzing time series – used in chaos theory

- Poincare map $M : \mathcal{R}^d \rightarrow \mathcal{R}^d$
 - Time series: $X_0, X_1, \dots, X_i, X_{i+1}, \dots$
 - generated as $X_{i+1} = M(X_i)$ $X_i = M^i(X_0)$

- Effect of Poincare map:

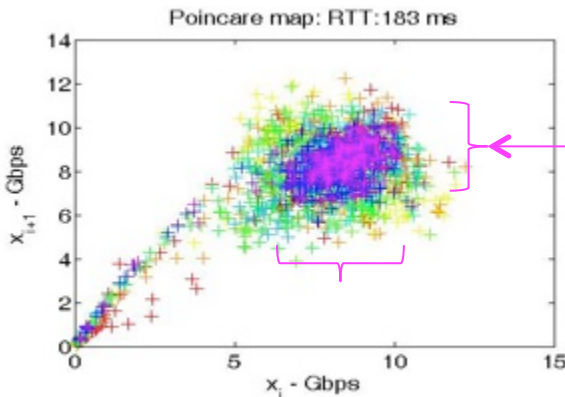
- range specifies achievable throughput



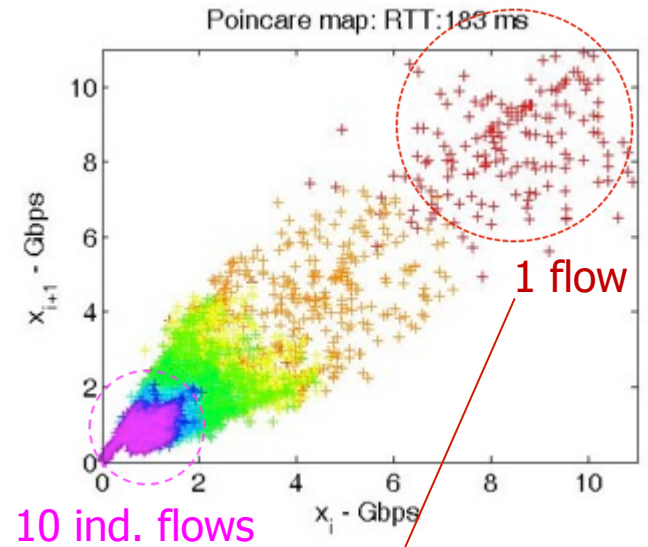
combined throughput

4 individual flows
CUBIC 46ms

- complexity indicates rich dynamics – lower throughput and narrow concave

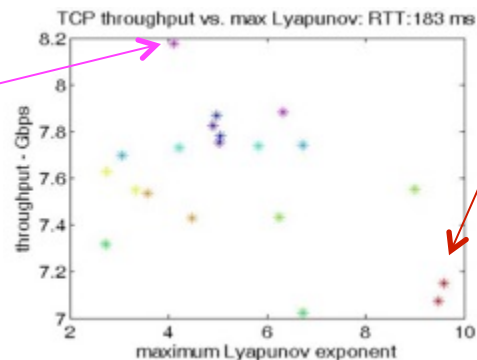


sum of 10 flows



10 ind. flows

1 flow



large Lyapunov exponent
• low throughput
next slide

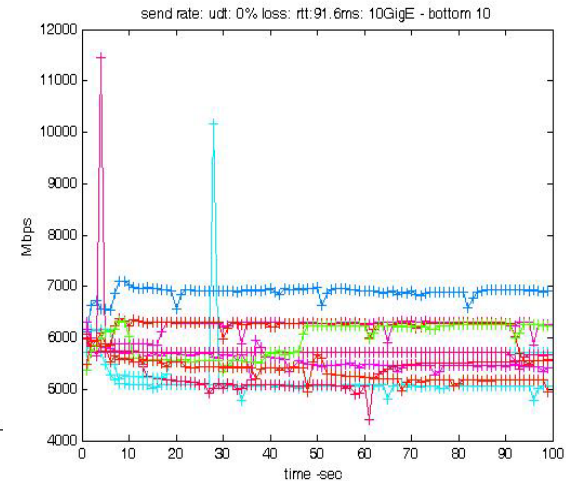
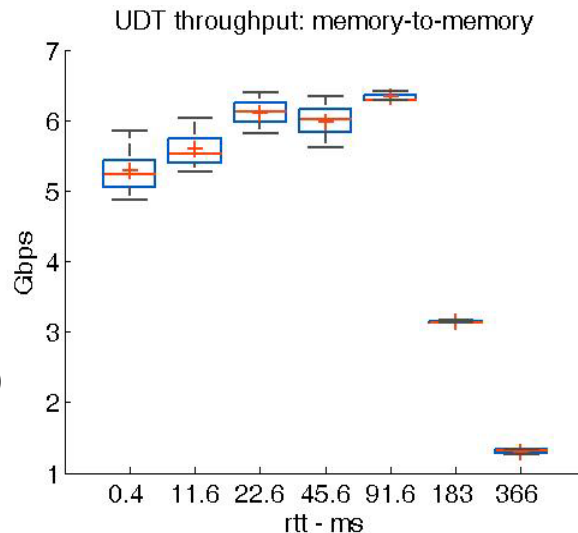
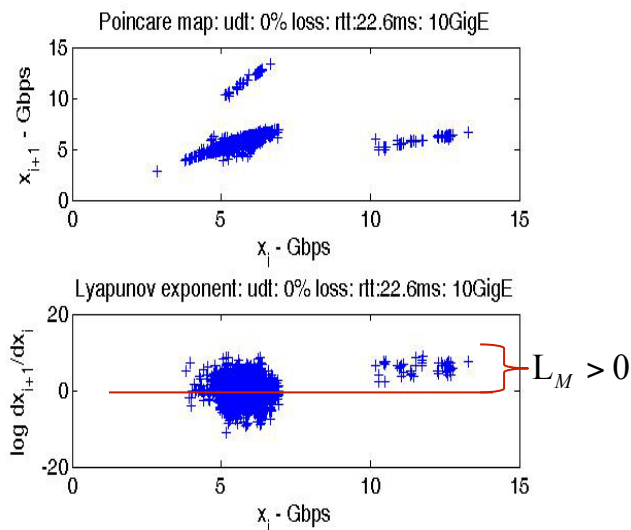
Lyapunov Exponent: Stability and Concavity

- Log derivative of Poincare map

$$L_M = \ln \left| \frac{dM}{dX} \right|$$

- Provides critical insights into dynamics

- Stable trajectories: $L_M < 0$
- Chaotic trajectories: $L_M > 0$
 - indicate exponentially diverging trajectories with small state variations
 - larger exponents indicate large deviations
- protocols are operating at peak at rtt
 - stability implies average close to peak - implies concavity
 - positive exponents imply lowered throughput – trajectories can only go down
 - » then, weak sustainment implies convexity



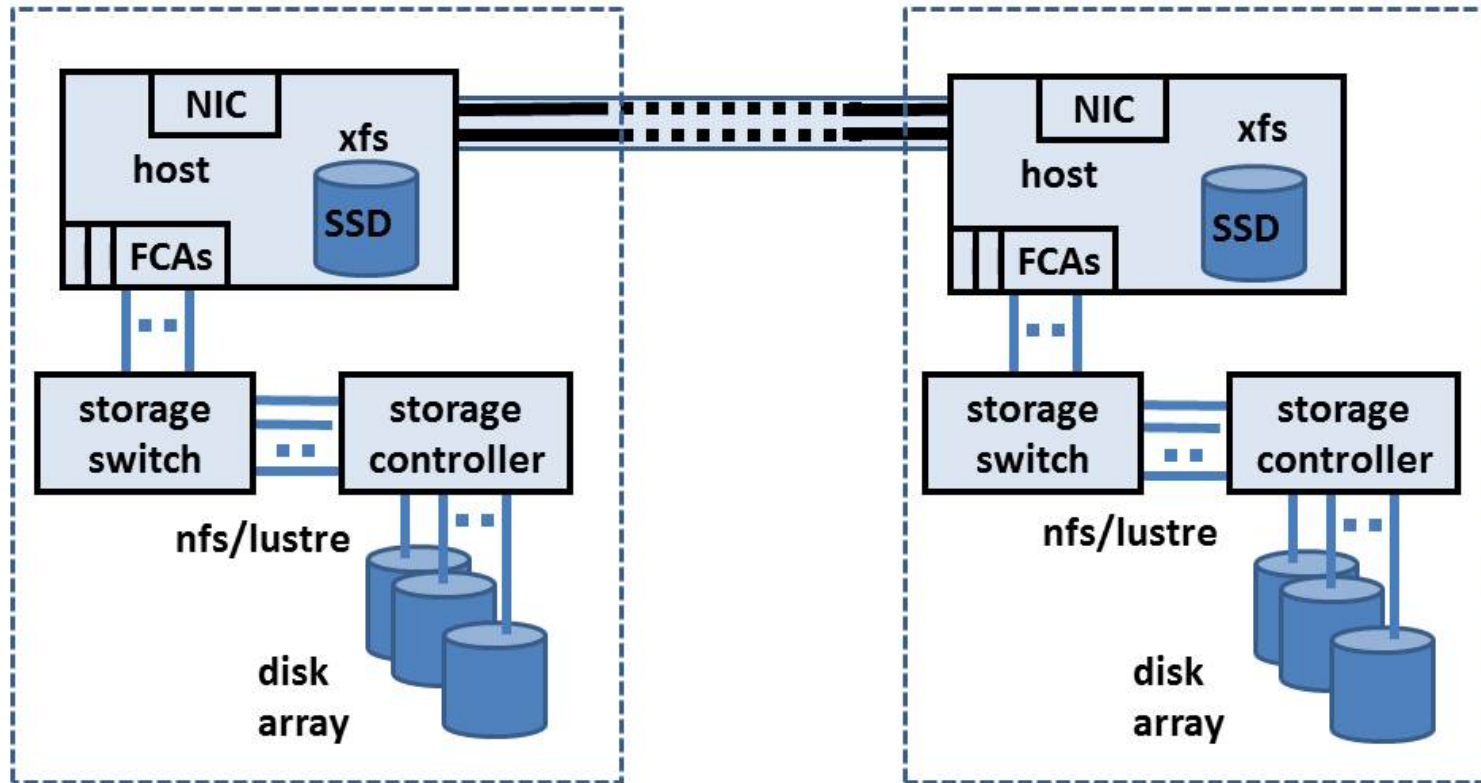
XDD: host-to-host file transfer tool

- XDD started as a file and storage benchmarking toolkit
 - storage tuning options
- Added a network implementation, python frontend, multi-host coordination, and NUMA tuning options
 - multiple NIC's from a single process for “better” storage access
- xddprof: sweep relevant tuning parameters
 - identify storage parameters to align with network performance profiles

xddmcp: Composing host-to-host flows

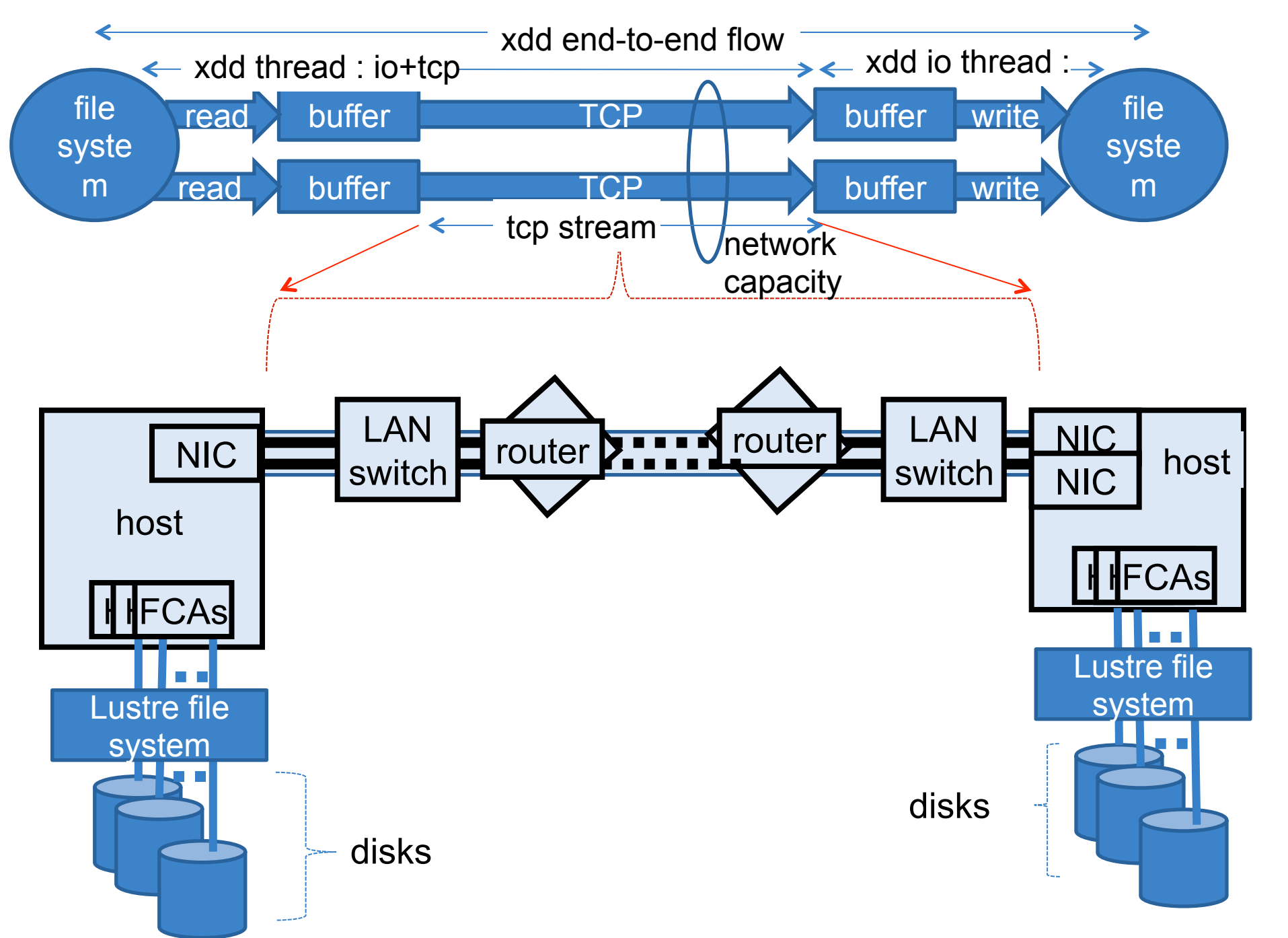
- Network and storage tuning optimizations aren't always complementary
 - network may prefer high number of streams
 - storage may prefer lower thread counts
- Leverage profiling information to understand performance
- Identify *compatible* network and storage parameters

ORNL Testbed: nfs, xfs and lustre file systems



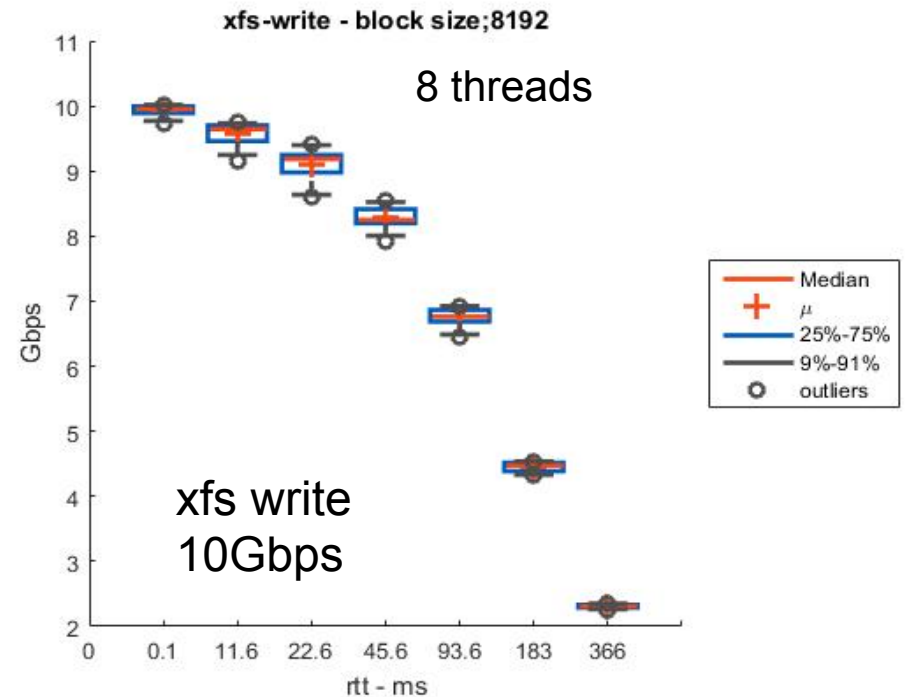
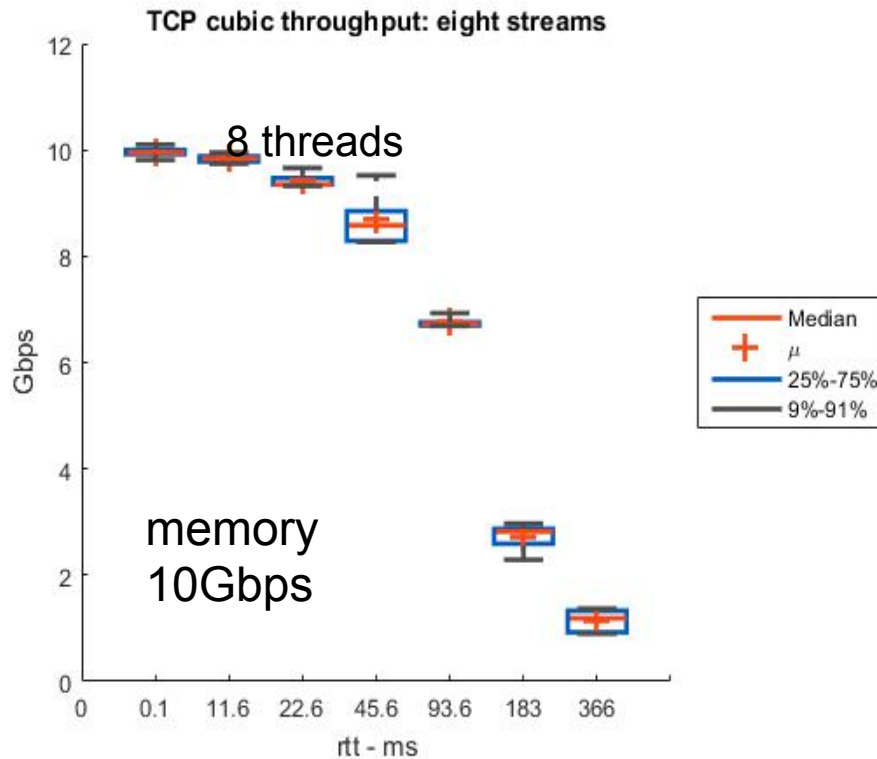
Peak IO rates: xddprof on hosts
nfs: ~2Gbps
xfs: ~40 Gbps
lustre: ~32 Gbps

Peak n/w throughput: iperf
TCP: > 9Gbps
UDP/T: > 8Gbps
for 0ms rtt



TCP CUBIC and xfs

- xddmcp host-to-host file transfers: peak: 10Gbps

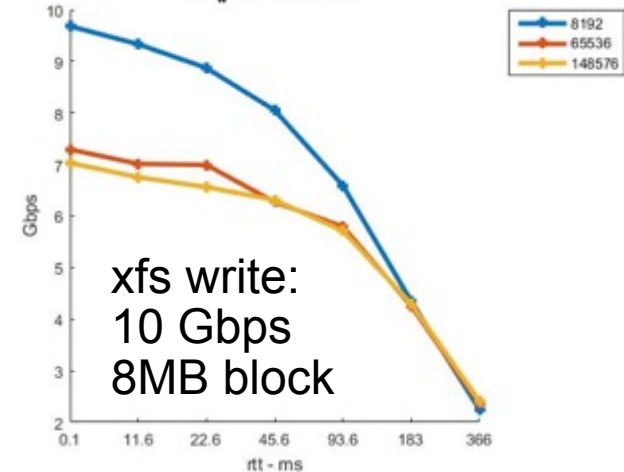
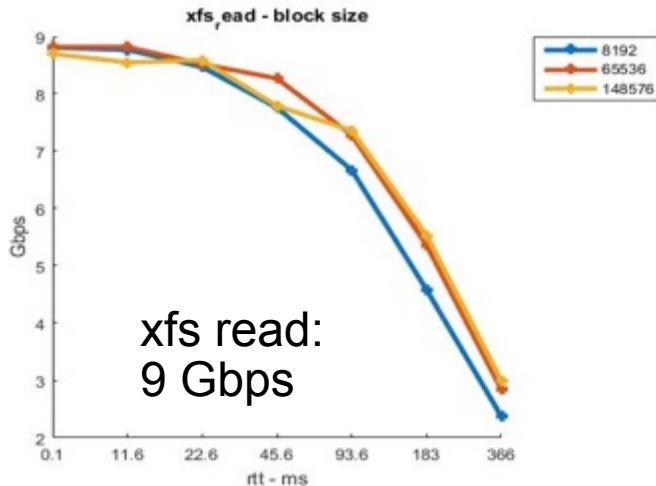
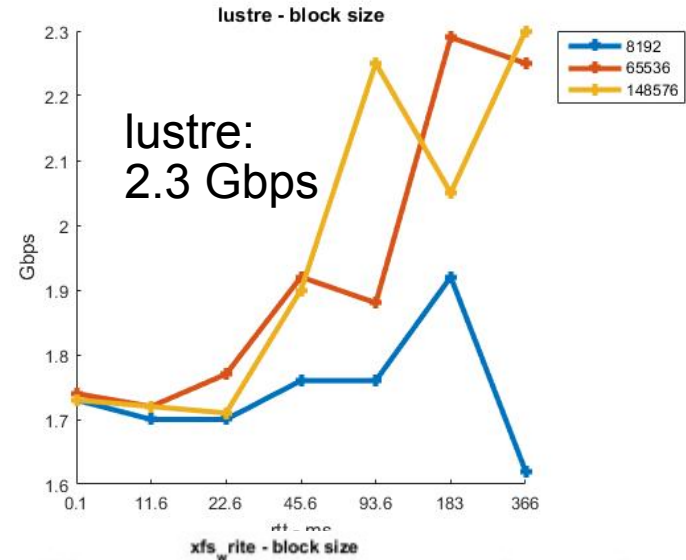
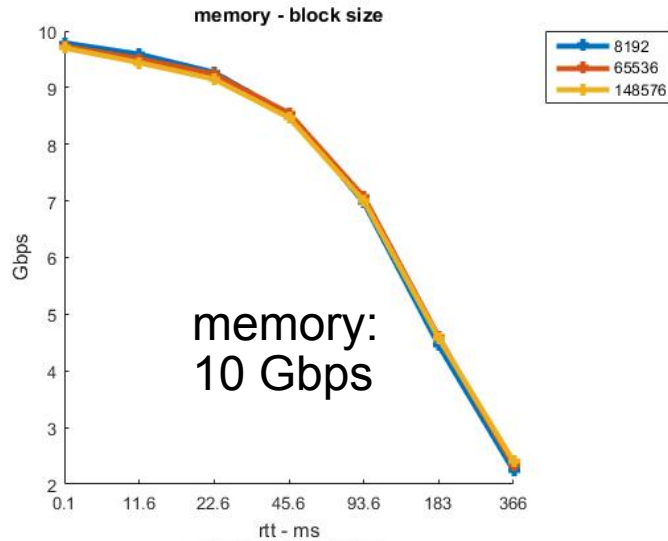


xdd file IO throughput is close to TCP throughput

- 8 IO threads and 8 TCP parallel streams
- Impedance mismatch is quite small

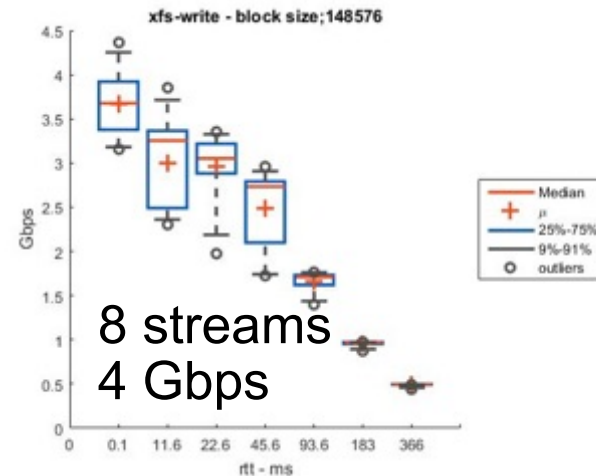
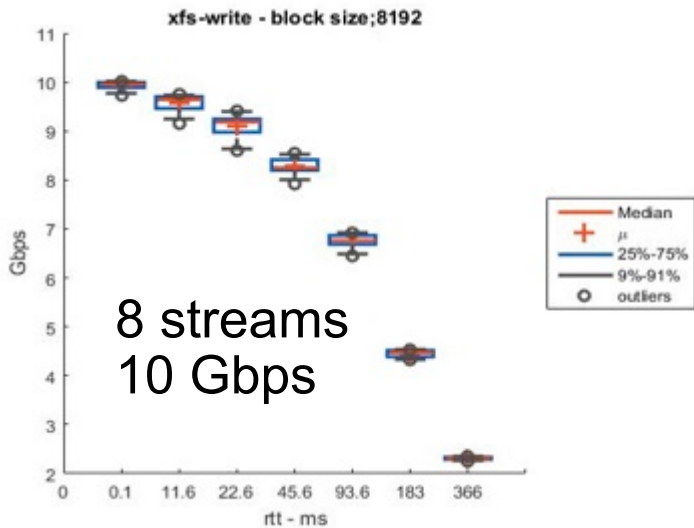
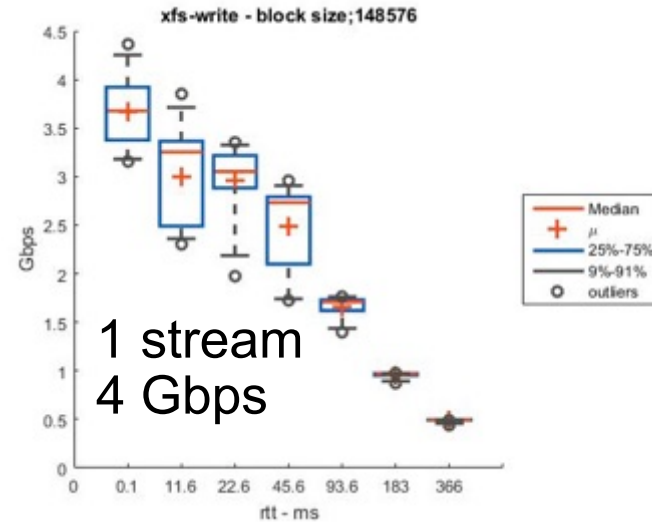
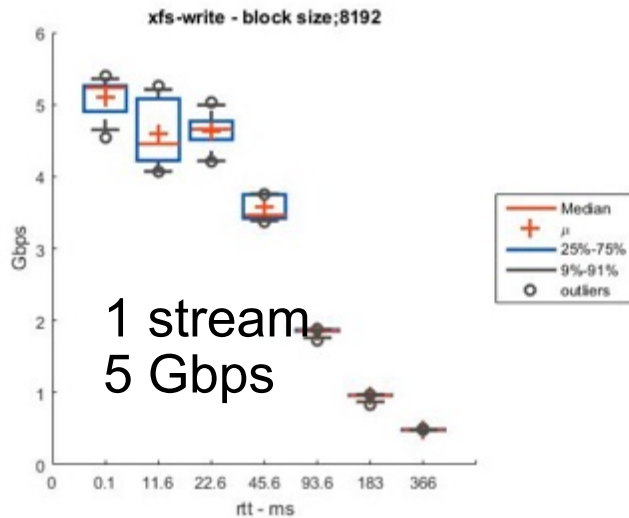
Average Throughput: lustre, xfs

8 streams: lustre throughput is lower compared to 1 stream



xfs: write

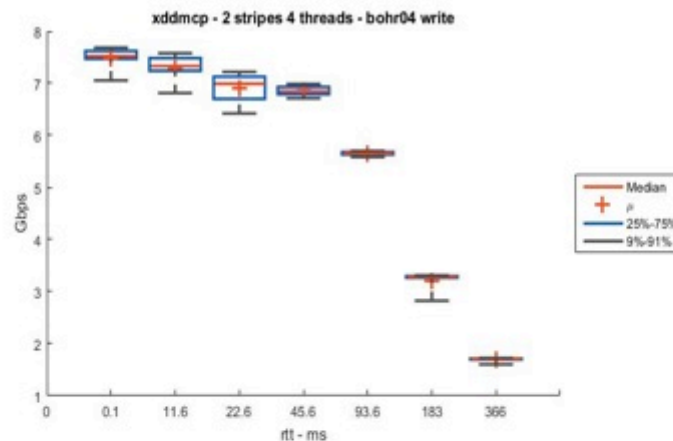
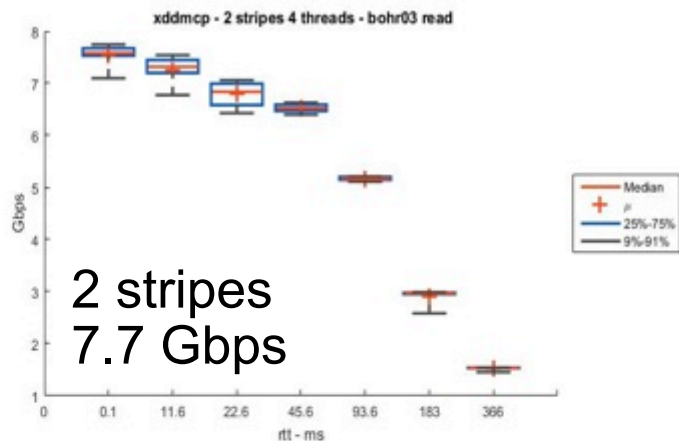
- 8 streams, 8MB blocks: 10Gbps



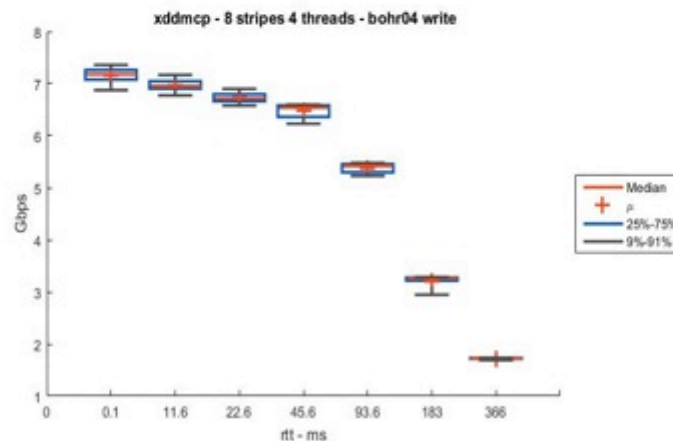
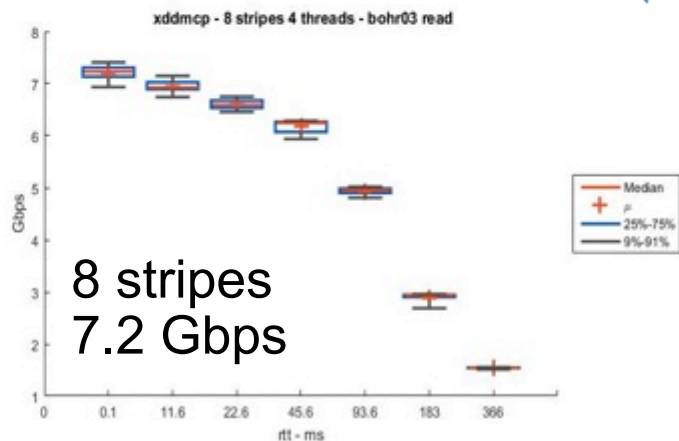
Smaller 8MB blocks for higher throughput

Best Case: Lustre default IO

4 streams – 2 stripes: 7.5 Gbps



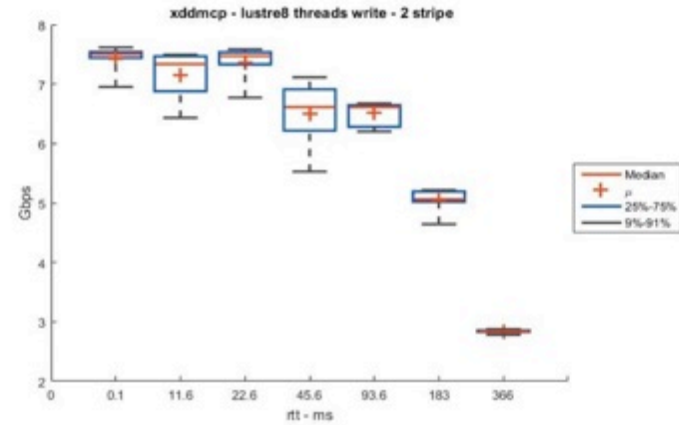
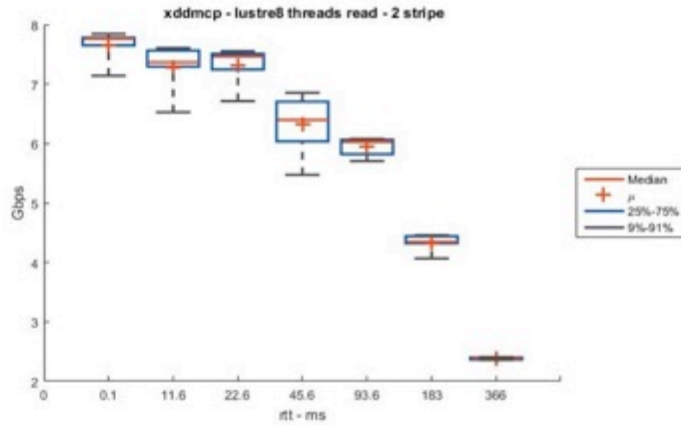
read ↔ write



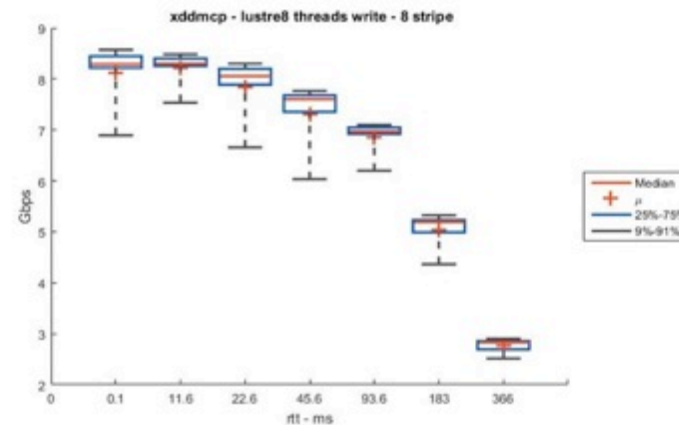
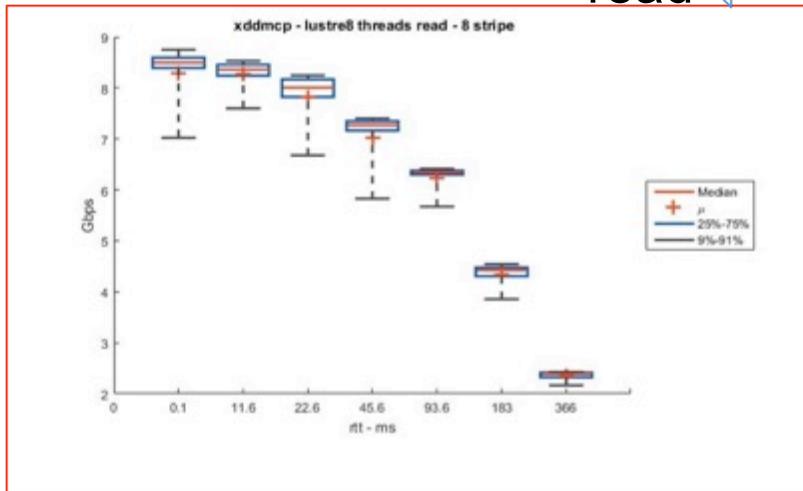
- 2 stripes provide higher throughput

Best Case with direct IO

- 8 threads – 8 stripes : 8.5Gbps



read ↔ write



8 stripes provide higher throughput by 1Gbps over default IO best case

d-w Method

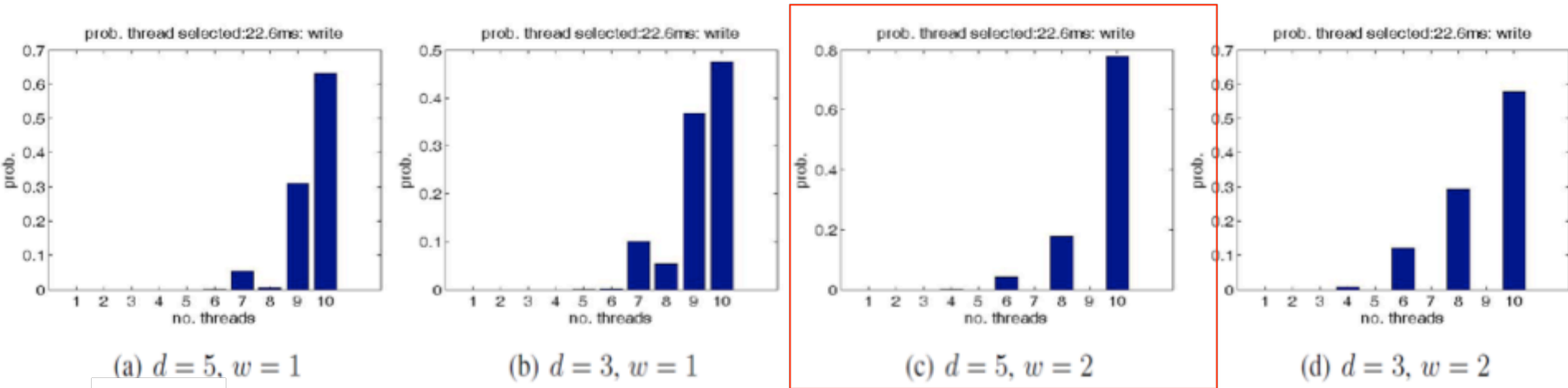
- Joint file I/O and network transport parameters for peak rate
 - Deriving them from full-profile takes months of measurements
 - Statistical variations are significant – gradient estimate is noisy
- Developed a *Depth-Width* or - method
 - exploits overall unimodality: throughput vs. number of streams
 - stochastic gradient search method: using repeated measurements over window
- Result: Identified peak configurations:
 - 97% of peak transfer rate for XFS and Lustre
 - probing 12% of parameter space – days vs. months for full profile

d-w algorithm

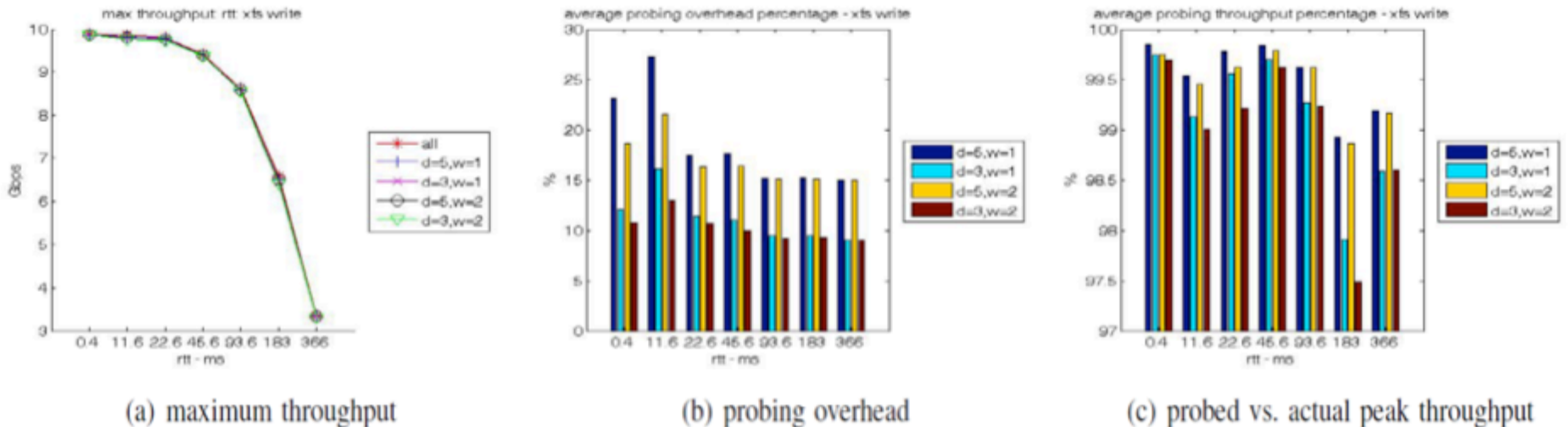
- initialization: start with largest number of flows
- repeat until halting criterion:
 - jump over a W -sized window for new probing configuration (different number of flows)
 - compute maximum throughput of d collected measurements at current configuration
- halting criterion:
 - maximum throughput decreases for two consecutive iterations

d-w Probing XFS Write

Confidence Estimate: fraction of 700 configurations conducive to d-w method



Confidence estimates of each configuration being selected with XFS file write transfer for $RTT = 22.6$ ms



Profile of d-w probing with XFS file write transfer

Conclusions

Contributions

- **Collected extensive transport measurements over dedicated connections**
 - **Multiple TCP variants and UDT: new insights**
 - concave region of throughput profile
 - rich dynamics: complex Poincare map and positive Lyapunov exponents
- **Developed throughput model**
 - simple enough not to require detailed protocol parameters
 - still explains the basic qualitative
 - concavity analysis
 - Poincare Maps and Lyapunov exponents: link dynamics to profiles
- **Applied for fine tuning XDD file transfers**

Future Directions

- **Detailed analytical models to explain concavity**
- **Automatic parameter optimization methods**
- **Differential methods:**
 - align analytical models with measurements
 - capture difference and apply corrections

Thank you Questions?

