

# On causes of GridFTP transfer throughput variance

---

Z. Liu, M. Veeraraghavan, J. Zhou  
University of Virginia

J. Hick  
NERSC

Y-T Li  
SLAC

IEEE/ACM NDM 2013 workshop  
In association with SC 2013  
Date: Nov. 17, 2013

Supported by NSF grants: OCI-1127340, OCI-1038058, OCI-1116081,  
DOE SC0007341 and ESnet Testbed grant DE-AC02-05CH11231

Thanks to Brian Tierney, Eric Pouyoul, and others on ESnet testbed team,  
Brent Draney, NERSC, and Wei Yang, SLAC, and Raj Kettimuthu, ANL



Questions? Contact Malathi Veeraraghavan, [mv5g@virginia.edu](mailto:mv5g@virginia.edu)

# Outline

---

- Background/contributions
- Transfer throughput variance observed
- Causes of variance
- **Throughput model** as a function of resource allocations (mem2mem)
- Impact of disk I/O contention
- Engineering solution (feedback?)
- Summary



# Related work

---

- OSCARS project: Use rate-guaranteed circuits across network to avoid packet losses
- Science DMZ project: Bypass firewall filters and campus LAN switches that could drop packets
- StorNet project: co-schedule network and storage resources
- Globus Online enables auto-selection of user-configurable parameters and thus avoids user mistakes (e.g., use -fast)



# Key contributions of this work

---

- Data Transfer Nodes (DTNs)
  - Dedicated servers that are used just for wide-area file transfers
  - Shared in **interactive** mode, not with a scheduler such as Portable Batch System (PBS)
- CPU cycles available at the DTNs are key determinants of mem2mem throughput
- Disk I/O is a major player for transfers involving disks
- Packet loss rate is a junior player on observed paths (research&education nets)



# Outline

---

- Background/contributions
- Transfer throughput variance observed
- Causes of variance
- Throughput model as a function of resource allocations (mem2mem)
- Impact of disk I/O contention
- Engineering solution (feedback?)
- Summary



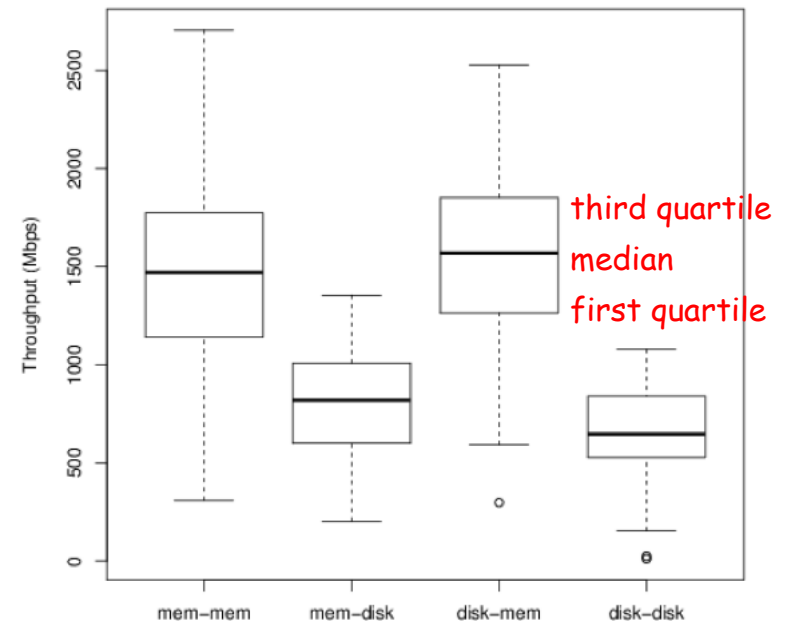
# Throughput variance

- 334 test transfers executed from ANL to NERSC production Data Transfer Nodes (DTNs)
- 4 types: mem-mem (84), mem-disk (78), disk-mem(87), disk-disk(85)

TABLE VI: Throughput of ANL-NERSC transfers (Mbps)

	mem-mem	mem-disk	disk-mem	disk-disk
Min	308.9	202.4	297.4	10.85
1st Qu.	1149	599.6	1265	527.3
Median	1472	819.0	1569	645.9
Mean	1463	789.6	1563	670
3rd Qu.	1772	1007	1851	841.3
Max	2706	1354	2529	1079
CV	35.69%	31.63%	30.80%	33.10%

Surprising result that mem2mem transfers had similar levels of variance as transfers involving disks



Z. Liu, M. Veeraraghavan, Z. Yan, C. Tracy, J. Tie, I. Foster, J. Dennis, J. Hick, Y. Li and W. Yang, "On using virtual circuits for GridFTP transfers," SC2012

# Problem Statement

---

- What are the causes of throughput variance?
- Find validated models for transfer throughput as a function of these factors
- Objective:
  - determine amount of concomitant resource allocations required at the two data transfer nodes (servers) and across the network to achieve a certain fixed throughput
  - schedule resources and achieve rate-guaranteed transfers



# Outline

---

- Background/contributions
- Transfer throughput variance observed
- **Causes of variance**
- Throughput model as a function of resource allocations (mem2mem)
- Impact of disk I/O contention
- Engineering solution (feedback?)
- Summary





# Three types of factors

---

- Intrinsic but controllable factors
- Intrinsic but less controllable factors
- Extrinsic factors



# Intrinsic but controllable factors

---

- **Application choices:**
  - GridFTP, scp, bbcp, FDT, RFTP
  - Run-time arguments such as -fast for data connection reuse for multiple file transfers, and -p for parallel connections in GridFTP
- **Transport-layer protocol choices:**
  - TCP or UDP based; If TCP: H-TCP, Cubic TCP?
  - Parameters: buffer size, window scaling, etc.
- **Solutions:**
  - User training: <http://fasterdata.es.net> makes recommendations
  - Globus Online



# Intrinsic but less controllable factors

---

- Round-trip propagation delay
  - Cannot be changed for transfers between specific **Data Transfer Nodes (DTNs): speed of light**
- Bottleneck link rate
  - Could upgrade links
- Maximum segment (frame) size
  - do all switches/routers support jumbo frames?
- File size
  - Small-size files will experience lower throughput
    - TCP Slow Start
    - Virtual circuit setup delay



# Extrinsic factors

---

- Competing tasks on DTNs impact resources available to a particular file transfer's processes
  - CPU cycles
  - Memory access rates
  - Disk I/O access rates
- Competing flows on shared network links
  - available capacity, packet losses, queueing delays (RTT impact)
- Focus of study



# Outline

---

- Background/contributions
- Transfer throughput variance observed
- Causes of variance
- Throughput model as a function of resource allocations (mem2mem)
- Impact of disk I/O contention
- Engineering solution (feedback?)
- Summary



# Methods used

---

- Developed scripts/tools by running experiments on testbeds
  - ESnet 100G testbed: CPU usage
  - U. Utah's Emulab: packet loss
- Used logins on **production DTNs** to execute instrumented transfers
  - Could observe actual contention for DTN resources
  - And for network resources



# Instrumented transfers between production DTNs

---

- Monitoring scripts: initiates *top* (CPU usage) and *tcpdump* (packet loss) before each transfer
- Scheduled hourly GridFTP transfers between NERSC and SLAC DTNs
- Stop data collection after transfers
- Analyze logs collected from GridFTP and monitoring tools
- Method: different from that used in previous papers in which independent sources of measurements were obtained (e.g., NWS for network)



# DTNs and network path

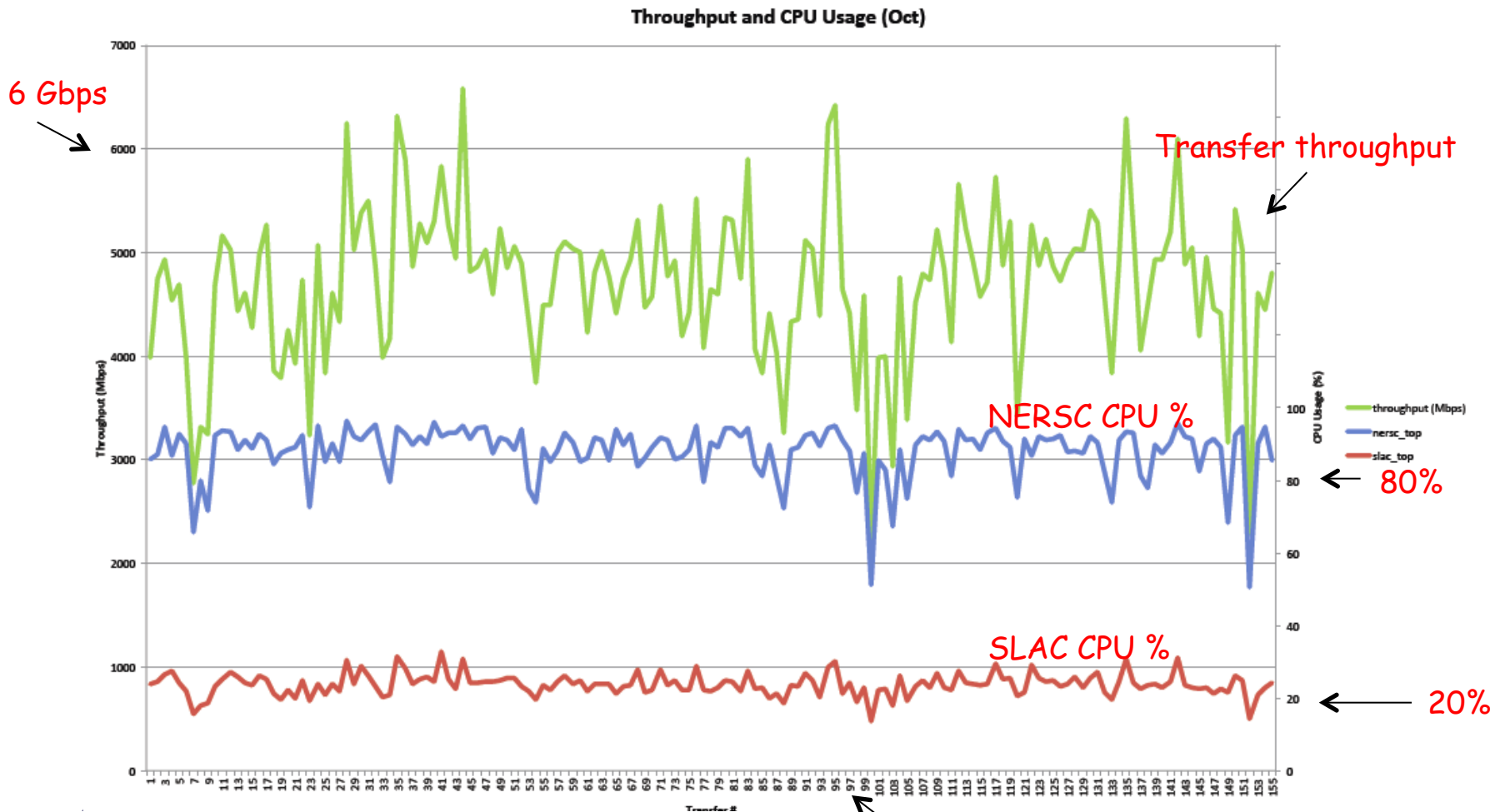
---

- NERSC Data Transfer Node:
  - 2 x AMD Opteron (dual-core)
  - 8 GB
  - CentOS 5.8 (2.6.18 x86\_64)
- SLAC Data Transfer Node:
  - 2 x Xeon (quad-core with HT enabled)
  - 48 GB
  - RHEL 5.9 (2.6.18 x86\_64)
- Bandwidth: 10 Gbps
- RTT: 2.47ms (10 routers on path)





# Transfer throughput dependence on CPU time



Consistent with ESnet background traffic load in the 2-3 Gbps range (ESnet4: 10 Gbps links)

Transfers (155 over a one-week period)

# Regression Model

---

- Dependent variable: throughput
- Independent variables:
  - NERSC CPU usage
  - packet loss rate
  - SLAC CPU usage?
  - SLAC CPU usage was highly correlated with NERSC CPU usage: linear model

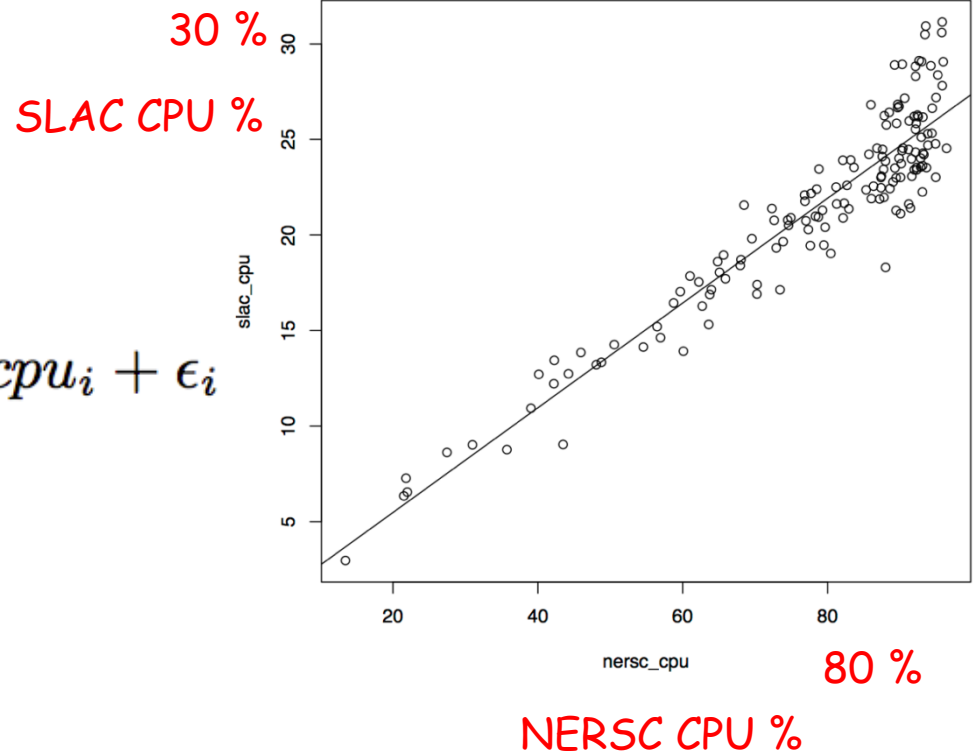


# Linear regression Model

- SLAC DTN had the more powerful CPUs; hence the y variable
- Represent SLAC CPU% consumed by our mem2mem transfer GridFTP process as a function of NERSC CPU% used by corresponding GridFTP process

$$SLACcpu_i = \beta_0 + \beta_1 NERSCcpu_i + \epsilon_i$$

Correlation coefficient: 0.94



# Non-linear regression Model

---

Throughput

Error term from  
previous regression

$$y_i = \beta_1' NERSCcpu_i + \beta_2' \epsilon_i + f(p_i) + e_i,$$

$$f(p_i) \approx \sum_{j=1}^{k+m} \alpha_j B_j(p_i)$$

$p_i$ : packet loss rate

recall Mathis et al. equation:  
 $y_i \sim 1/\sqrt{p_i}$

$k = 2, m = 3$   
chosen by Leave-One-Out  
Cross-Validation (LOOCV)

$f$  should be non-linear

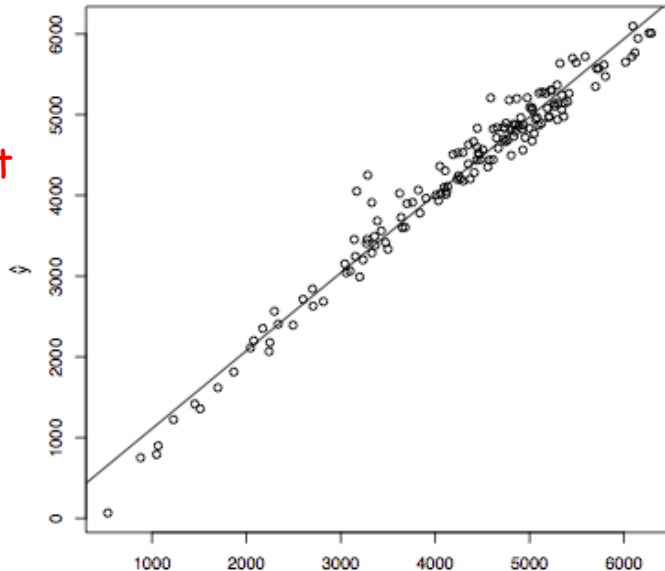
use B-Splines



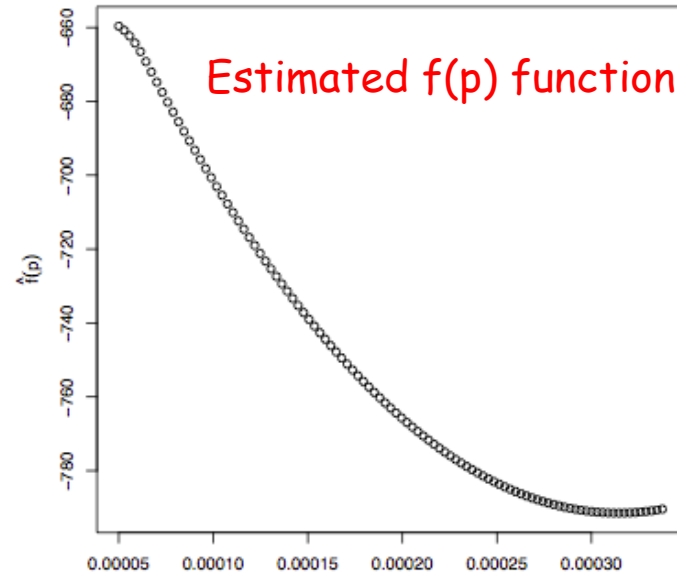
Acknowledgment: Prof. Jianhui Zhou, Stats Dept.

# Good fit: Adjusted R-squared = 0.997

Estimated throughput



Actual throughput



Packet loss rate

(a) Throughput vs. estimated throughput

(b) Estimated  $\hat{f}(p)$  function

Figure 8: Results of the regression model

- Range of packet loss rates:  $(5e-5, 3.5e-4)$
- Well-known boundary effects of smoothing methods



# Results

Variable	Mean value
Retransmission Rate ( $p_i$ )	4.1E-05
$f(p_i)$ in Mbps	-683.3
Throughput (Gbps)	4.226
NERSC CPU usage (%)	78.28
SLAC CPU usage (%)	21.45

$$y_i = \beta'_1 \text{NERSCcpu}_i + \beta'_2 \epsilon_i + f(p_i) + e_i,$$

Mean values of regression coefficients:  $\beta'_1 = 62.708$ ,  $\beta'_2 = 153.194$

$$\begin{aligned} \text{Mean value: } & 62.708 * 78.28 + 153.194 * 0 - 683.3 \\ & = 4.9 \text{ Gbps} - 0.683 \text{ Gbps} = 4225.48 \text{ Mbps} \end{aligned}$$

- CPU usage was the primary factor in determining throughput
- However packet loss rate, while small, contributes to throughput reduction ( $683.3/4225.48 = 16\%$ )



# Key finding

---

- To control variance, the number of concurrent processes on the DTNs have to be controlled
- Current approach:
  - DTNs are used in interactive mode
  - As users login to DTNs and initiate file transfer apps as needed, the amount of CPU and disk resources available to a particular transfer are not controlled



# Outline

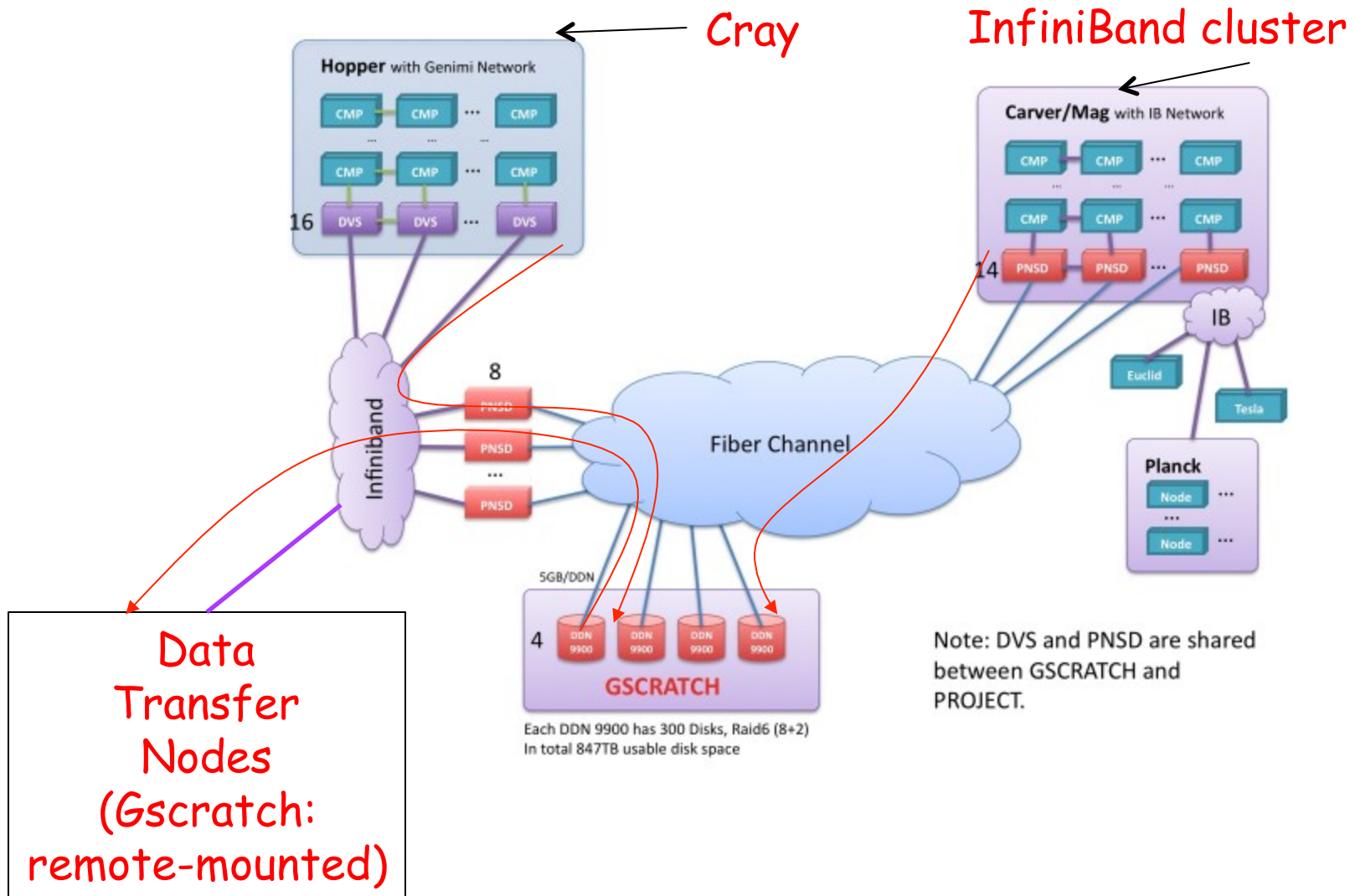
---

- Background/contributions
- Transfer throughput variance observed
- Causes of variance
- Throughput model as a function of resource allocations (mem2mem)
- **Impact of disk I/O contention**
- Engineering solution (feedback?)
- Summary





# NERSC systems (competition for disk access)



# Disk-mem experiment (NERSC-to-SLAC DTNs)

---

1. Invoke `dd` to write an 8 GB file (first file) to the global scratch file system with the `sync` system call to force the completion of pending disk writes, and record the time taken for the write operation.
2. Invoke `dd` to write another 8 GB file (second file) to the global scratch file system to ensure that the first file is no longer in the filesystem cache, which is required for the next step.
3. Invoke `dd` to read back the first file (which is now on disk, not cache) and record the time taken for the read operation.  
(Conveniently second file will be forced out of the cache).
4. Using `globus-url-copy`, transfer the second file to `/dev/null` on the SLAC DTN, and use `strace` to record system calls for further analysis. The only disk I/O operation required is the reading of the second file, which we know is not in the cache.



# Implemented scripts

---

- Characterize variability in disk read and disk write times for the shared file system global scratch at NERSC
- Parse strace output to find contribution of disk I/O access times to file transfer times
- Measured impact of strace: small



# Variability in disk I/O times

---

Table 5: NERSC Global Scratch file system operations

	Read (MB/s)	Write (MB/s)
Min	109.4	173.1
1st Qu.	247.6	530.2
Median	306.9	585.6
Mean	372.8	559.3
3rd Qu.	513.0	654.8
Max	746.8	747.4
CV	47.5%	23.5%

- 116 data points obtained
- Significant variance since file system is shared between computational systems and DTNs



# Results

- Disk access time matters, but is not the only resource that impacts throughput

This transfer's throughput was only 1.8 Gbps

Table 9: NERSC-to-SLAC disk-to-mem transfers

	Ratio of Disk IO time to total transfer time	Throughput (Gb/s)
Min	0.2647	0.768
1st Qu.	0.3415	1.687
Median	0.4010	1.994
Mean	0.4524	1.989
3rd Qu.	0.5523	2.32
Max	0.8191	3.228
CV	30.37%	26.42%

This transfer spent 81.24% of time in disk I/O

This transfer spent 35.08% of time in disk I/O



# Outline

---

- Background/contributions
- Transfer throughput variance observed
- Causes of variance
- Throughput model as a function of resource allocations (mem2mem)
- Impact of disk I/O contention
- Engineering solution (feedback?)
- Summary



# Engineering: proposed solution

---

- For low-variance, high-throughput transfers (part of workflow)
  1. run managed FT processes on DTNs (disable interactive; dedicated nodes for large dataset transfers)
  2. leverage Science DMZ to stage a local copy from shared filesystems
    - need a 2-phase cycle to alternate between unmanaged transfers from shared file systems and managed transfers across WAN
  3. calibrate required CPU times at two ends and VC rate (create nonlinear regression models - server dependent)
  4. schedule FT processes at two ends with PBS and schedule VC
- RoCE will reduce dependence on CPU time, but still need to schedule FT process on CPUs because of disk access competition
  - different nonlinear regression model for determining resource requirements



# Hobson's choice

---

- Choice between
  - low-variance, potentially increased waiting time for a high-throughput VC
  - high-variance, lower waiting time for a high-throughput IP-routed path
    - elephants stomp over mice
- if a DTN pair can sustain high throughput, request a high-rate VC, but wait time may be more for scheduler (OSCARS) to assign a circuit
  - to lower waiting time, lower utilization





# Summary

---

- **Science:**

- Theodore von Karman: "Scientists study the world as it is;  
Engineers create the world that never has been"
- Found answers to these questions by testing on operational DTNs
  - Are CPU resources in DTNs under contention: answer is yes for NERSC
  - Are disk I/O resources in contention: answer is yes for NERSC
  - Do packet losses occur on IP-routed paths: yes, low loss rate, but not insignificant impact on throughput
  - Can a model be found to express throughput as a function of resource allocations?  
Yes, but model required on a per-DTN pair basis

- **Engineering:** Given feasibility of developing a model, we proposed a solution using CPU schedulers, VCs, and file staging

- **Feedback:** Is it worthwhile developing this proposed solution?

- Backend transfers - administrative use? Reliable, fast, predictable transfers but possibly with delayed start times
- Did observe 12 TB transfers SLAC-BNL and 2.8 TB NCAR-NICS

