

*Exceptional service in the national interest*



SAND2014-15086PE

Sandia  
National  
Laboratories



# An Evaluation of BitTorrent's Performance In HPC Environments

Matthew G. F. Dosanjh, Patrick G. Bridges,  
Suzanne M. Kelly, James H. Laros III,  
Courtenay T. Vaughan



Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

# Motivation

- Need scalable system services in large HPC systems
  - Example services: file systems, job launch, system monitoring
  - Communication bottlenecks; performance degradation
- Can we use recent distributed systems techniques here?
  - E.G. Peer-to-peer services developed for Internet systems
  - Must scale to large systems with complex topologies and bottlenecks
  - Not designed for HPC use-cases, networks, etc.

# An Example Problem

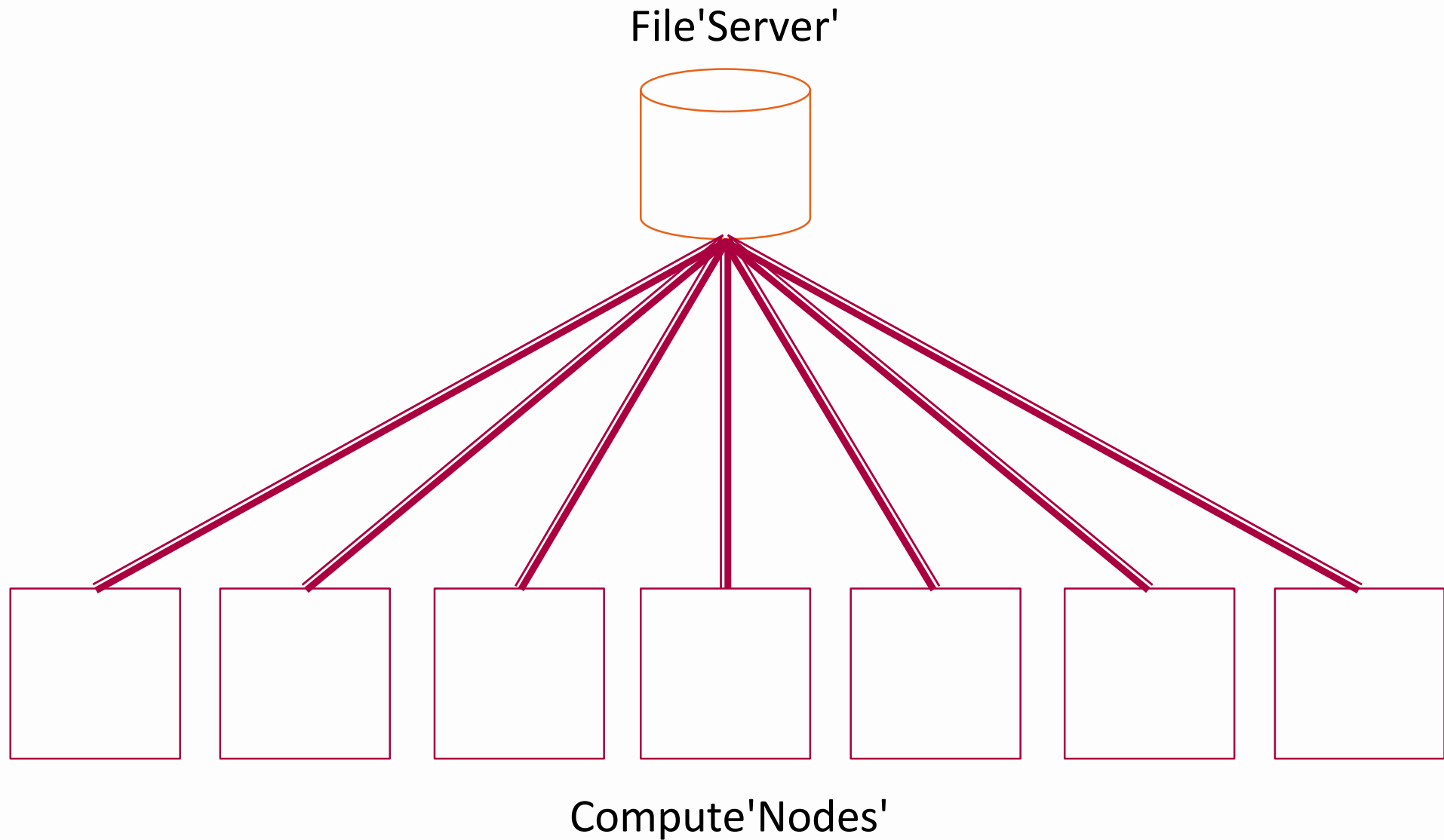
- Dynamic Shared Libraries distribution on HPC platforms
  - Read-only data-transfer
  - Many simultaneous downloaders
  - Looks a lot like a peer-to-peer problem
- Can we use BitTorrent?

# DYNAMIC SHARED LIBRARIES

# The Problem

- Demand for Dynamic Shared Library support is increasing
- Two methods of use:
  - Linked at program launch
  - `dlopen()` during runtime

# The Naïve Approach



# Current Solutions

- There are a number of current solutions to this problem
  - Cache nodes
  - Bulk pre-distribution
  - Tree-based overlay networks
- These solutions look like peer-to-peer networks
  - Can we use an off-the-shelf peer-to-peer technology?

# BITTORRENT

# Why BitTorrent?

- One of the more popular peer-to-peer protocols
- Open-source implementations
- Private networks
- Extensible protocol
  - Peer Exchange
  - Magnet Links

# Preparation

- Tracker
  - Daemon listens for incoming peers
  - Maintains a list of peers
- Description File
  - User specified file or directory
  - Splits file(s) into smaller pieces
  - Lists tracker information
- Initial Seeder
  - Starts with a full set of data
  - Registers with tracker

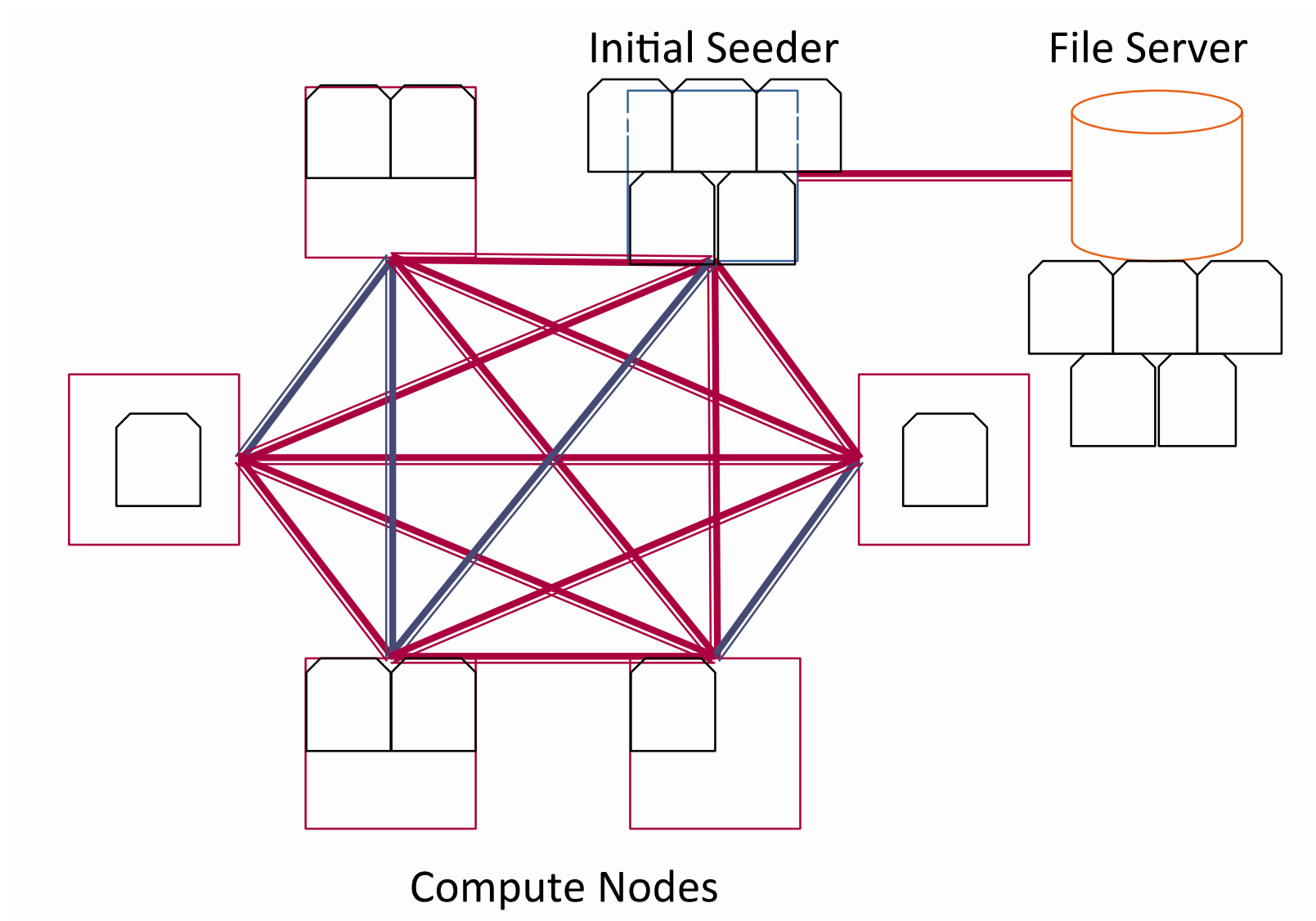
# Starting a Client

- Reads description file
- Connects to tracker
  - Registers as a peer
  - Gets a list of peers
- Connects to other peers

# Data Distribution

- Leechers
  - Keep a list of piece availability
  - Request pieces they don't have
- Seeders and Leechers
  - Respond to requests
  - Distributes Data according to fairness algorithm
- Fairness algorithm
  - Leechers use Tit-for-tat
  - Seeders try to push complete copies

# Data Distribution

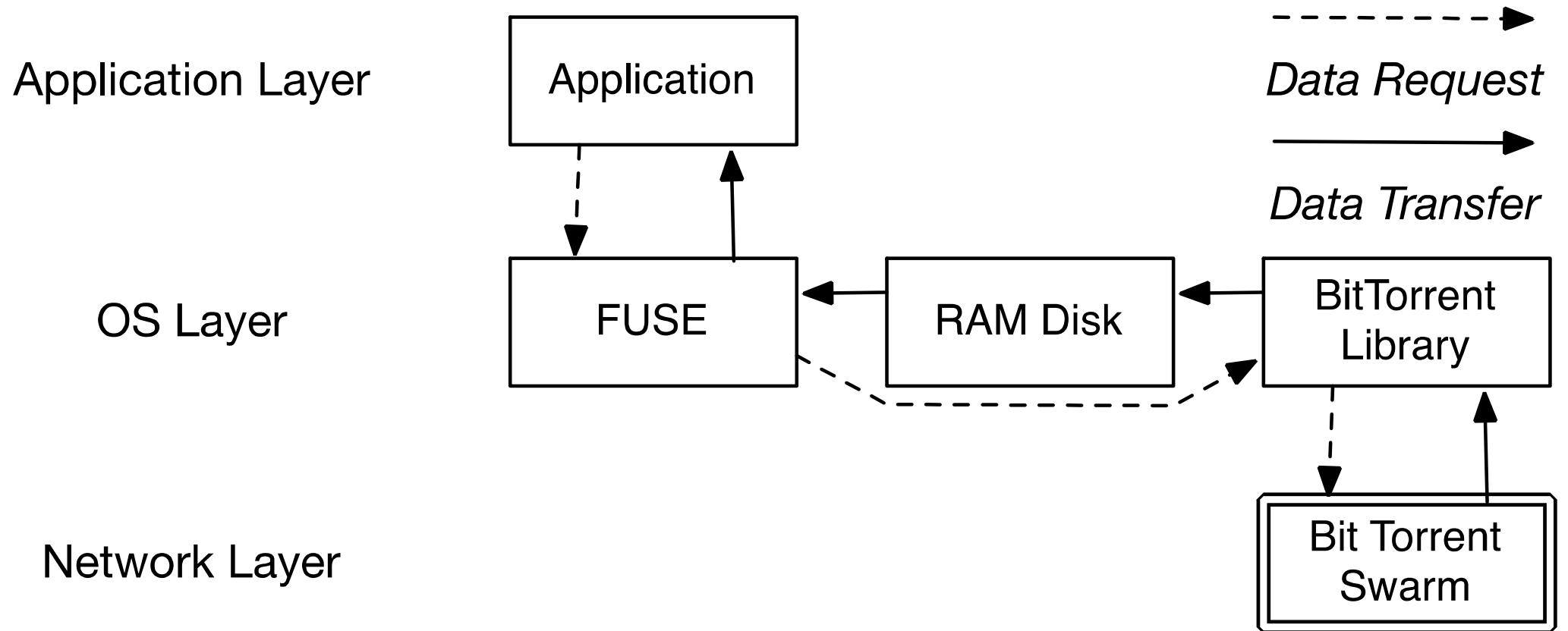


# ARCHITECTURE

# Added Actions

- Before Job Launch
  - Create the description file
  - Start Tracker & Initial Seeder
- During Job Launch
  - Distribute the description file
  - On each node
    - Start FUSE client
    - Add directory to LD\_LIBRARY\_PATH
- During Runtime
  - DL requests first search the FUSE directory
  - If a file isn't downloaded yet, request it from the BitTorrent swarm
  - If a file is downloaded, pass through to a local RAMDisk

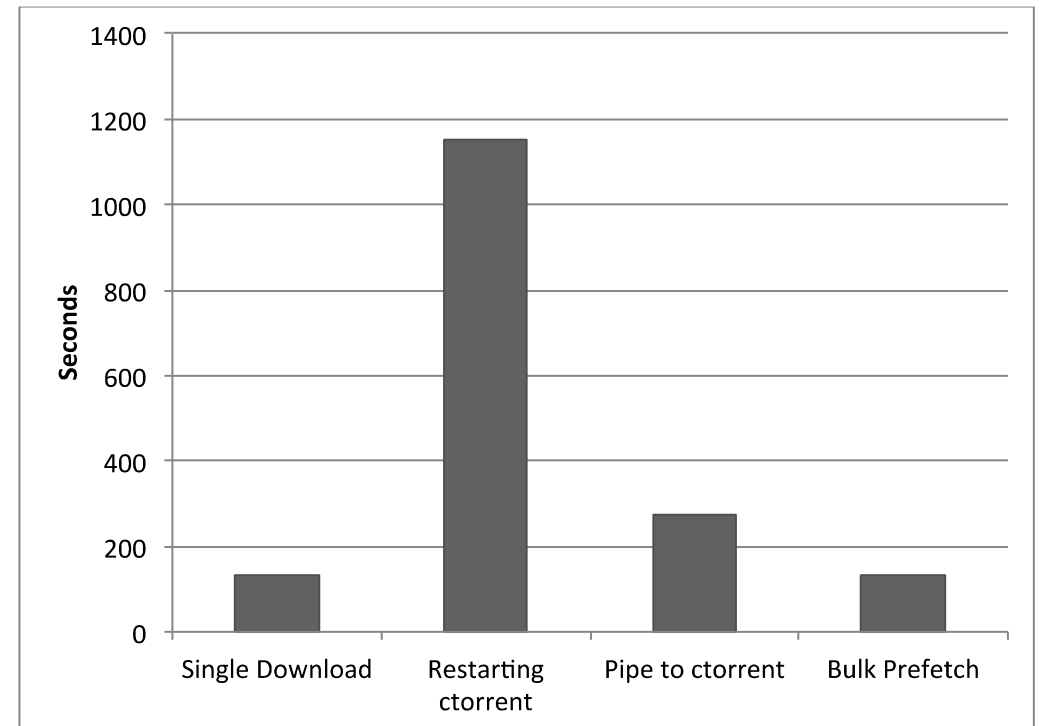
# Data Request Path



# MAKING BITTORRENT WORK FOR HPC

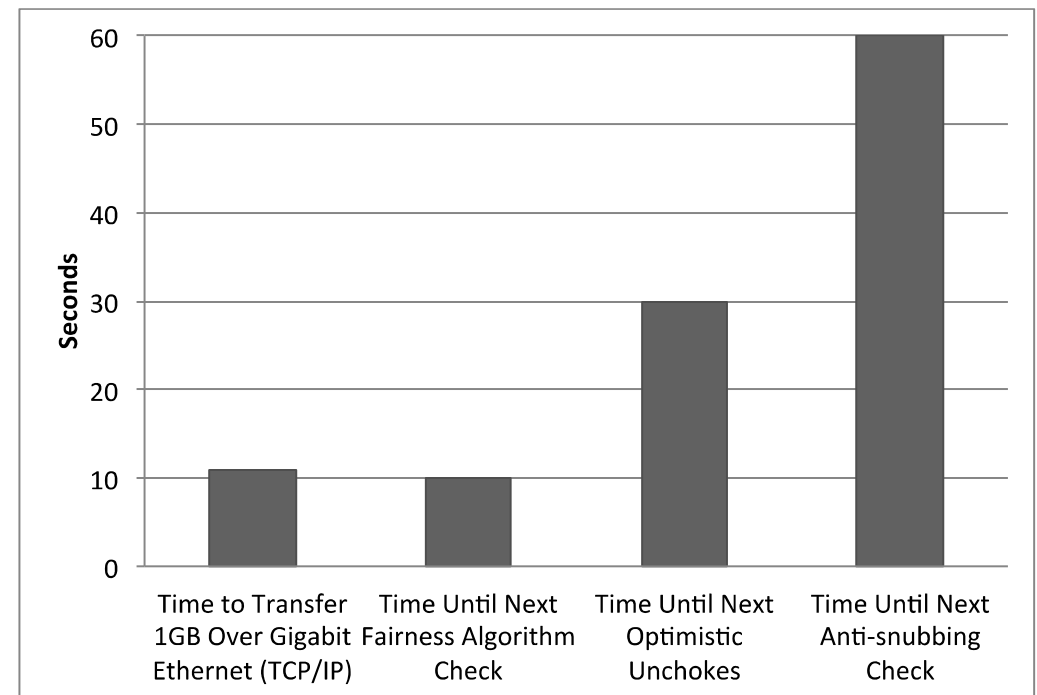
# BitTorrent Assumptions Get in the Way

- Basic approach:  
Downloads individual files on-demand
- Problem: Violates basic BitTorrent usage assumption



# Fairness Algorithm Problems

- BitTorrent tries to tolerate adversarial peers
- Fairness algorithm tuned for Internet speeds and timescales
- Sequential requests interfere with tit-for-tat choking algorithm



# Working Around Fairness

- Solution: Bulk prefetching of launch loaded libraries
- Libraries can be loaded two ways
  - On launch
  - On demand (dlopen)
- List of launch-loaded libraries provided by ldd
- Downloading these together, we match BitTorrent's use case
- We still support on-demand requests

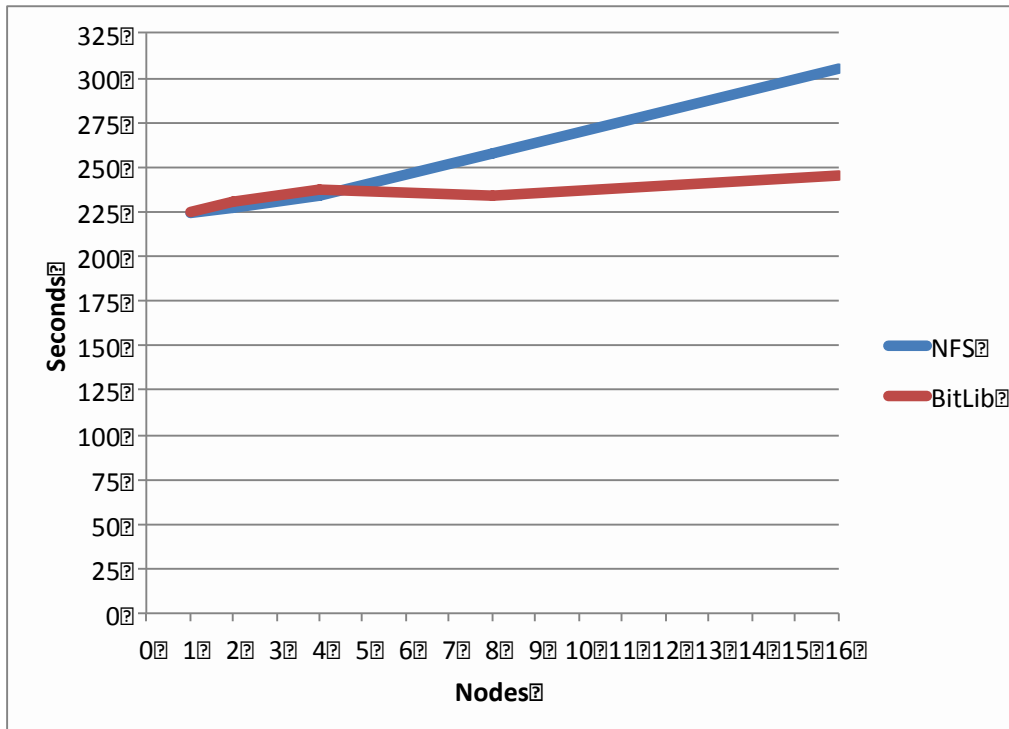
# RESULTS

# Experimental Setup

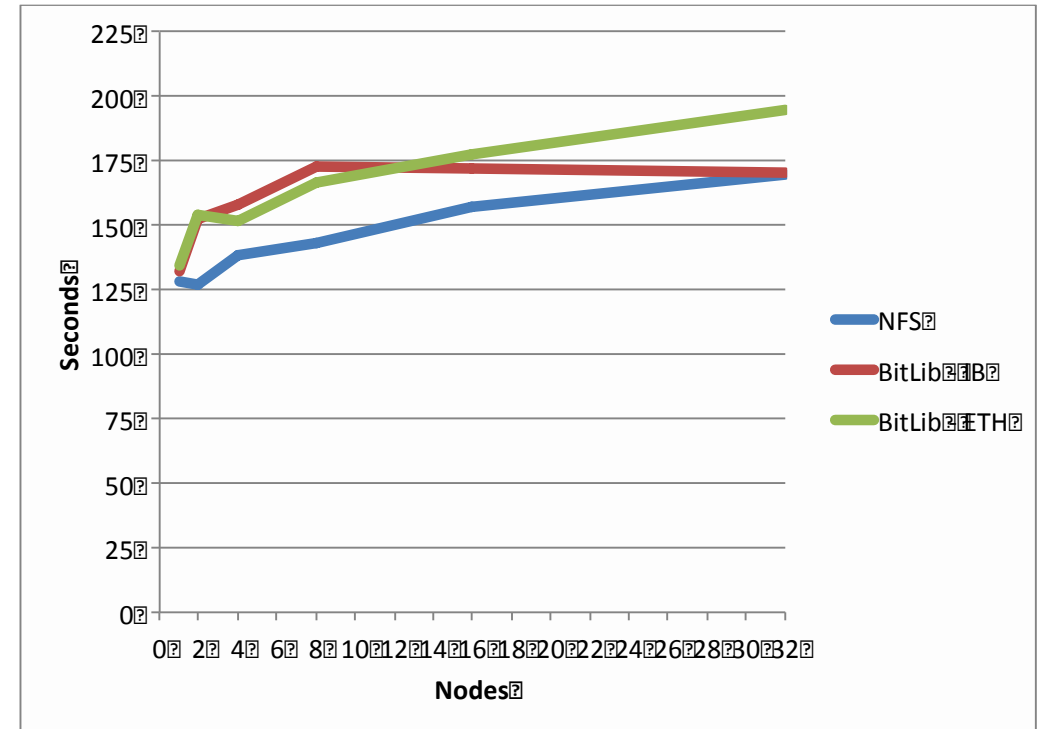
- The Pynamic benchmark
  - Test case has 495 libraries totaling ~1.1GB
  - Test case based on an LLNL scientific code
- Three Machines
  - Teller – 100 node, AMD Fusion with Infiniband
  - Muzia – 20 node, Cray XE6
  - Cielo - 8,944 Node, Cray XE6

# Small Scale

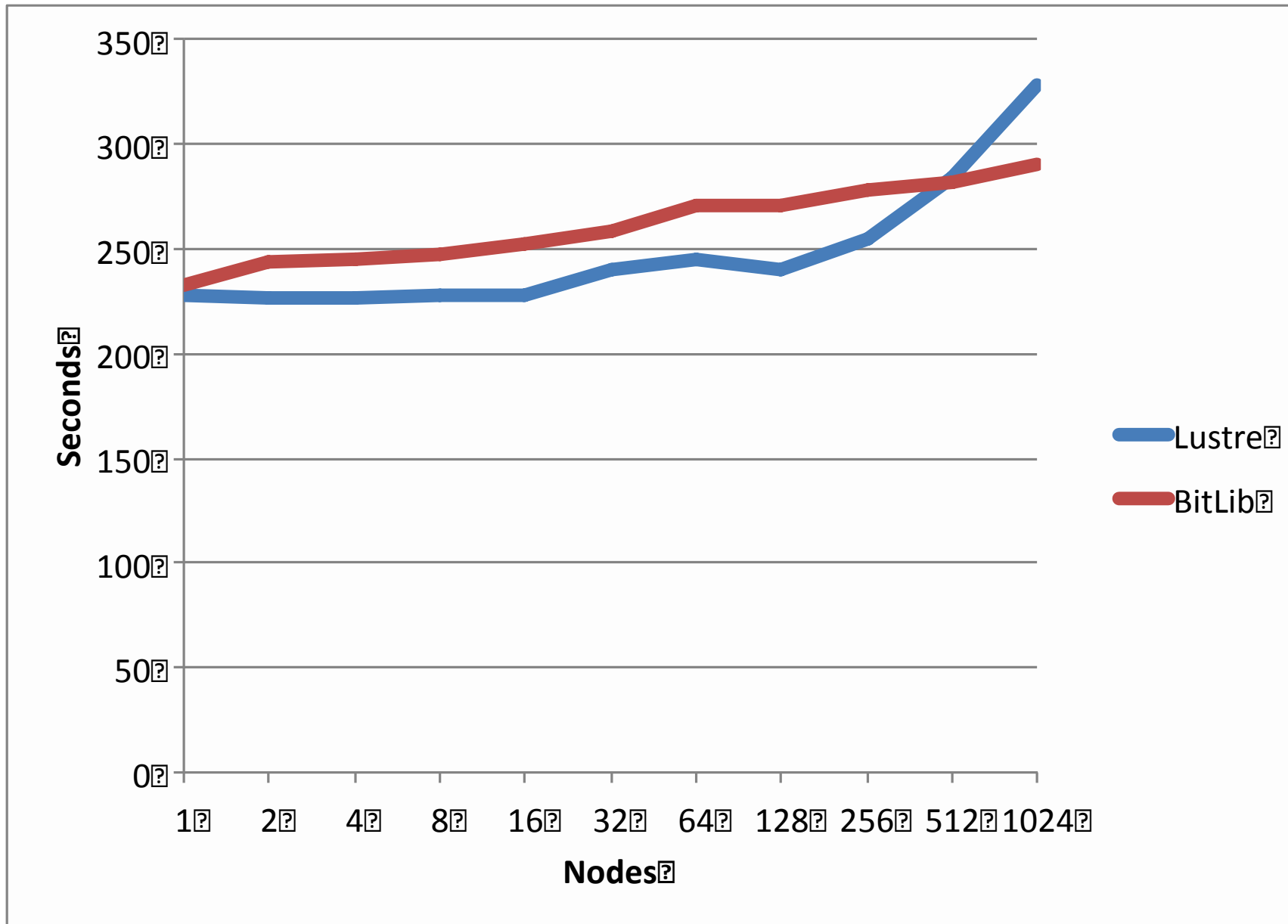
## Muzia: a Gemini system



## Teller: an Infiniband system



# Cielo



# ISSUES AND FUTURE WORK

# Network Features

- Can we use HPC network features to enhance and optimize BitTorrent?
- Multicast
- Topology-aware
- Multiple-paths

# Acknowledgements

This work was funded by NNSA's Advanced Simulation and Computing (ASC) Program. Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under contract DE- AC04-94AL85000.

**QUESTIONS?**