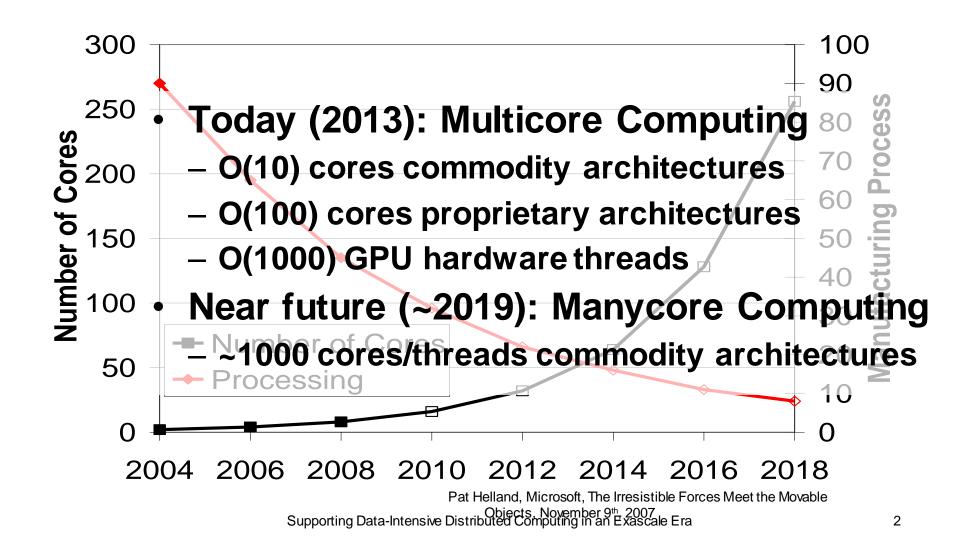# Supporting Data-Intensive Distributed Computing in an Exascale Era
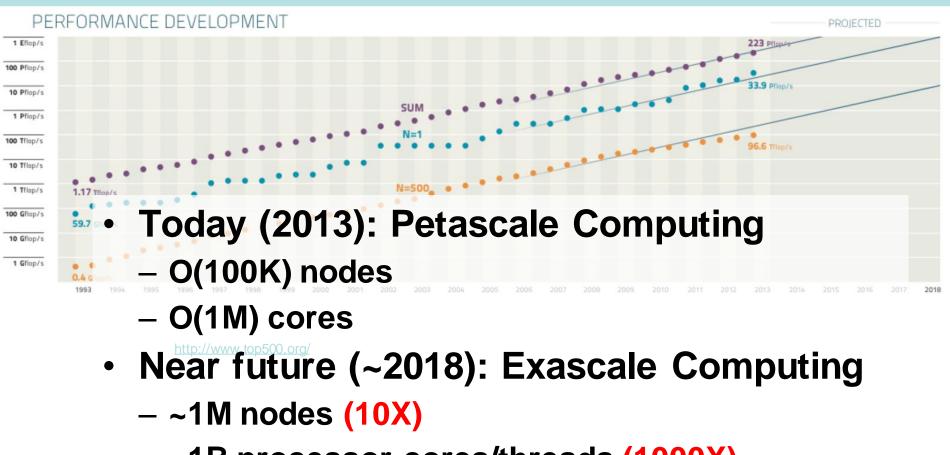
**Ioan Raicu**

Computer Science Department, Illinois Institute of Technology
Math and Computer Science Division, Argonne National Laboratory

August 7th, 2013
MAGIC Meeting: 2020-2025 Scientific Computing Environments

digitalblasphemy

# Manycore Computing

- **Today (2013): Multicore Computing**
  - **O(10) cores commodity architectures**
  - **O(100) cores proprietary architectures**
  - **O(1000) GPU hardware threads**
- **Near future (~2019): Manycore Computing**
  - **~1000 cores/threads commodity architectures**

Chart axes: Number of Cores (left, 0–300), Manufacturing Process (right, 0–100), years 2004–2018. Legend: Number of Cores, Processing.

# Exascale Computing

PERFORMANCE DEVELOPMENT

PROJECTED

| | | |
|---|---|---|
| 1 Eflop/s | | |
| 100 Pflop/s | | 223 Pflop/s |
| 10 Pflop/s | | 33.9 Pflop/s |
| 1 Pflop/s | SUM | |
| 100 Tflop/s | N=1 | 96.6 Tflop/s |
| 10 Tflop/s | | |
| 1 Tflop/s | 1.17 Tflop/s | N=500 |
| 100 Gflop/s | 59.7 G... | |
| 10 Gflop/s | | |
| 1 Gflop/s | 0.4 G... | |

1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018

- **Today (2013): Petascale Computing**
  - **O(100K) nodes**
  - **O(1M) cores**

http://www.top500.org/

- **Near future (~2018): Exascale Computing**
  - **~1M nodes (10X)**
  - **~1B processor-cores/threads (1000X)**

http://s.top500.org/static/lists/2013/06/TOP500_201306_Poster.png
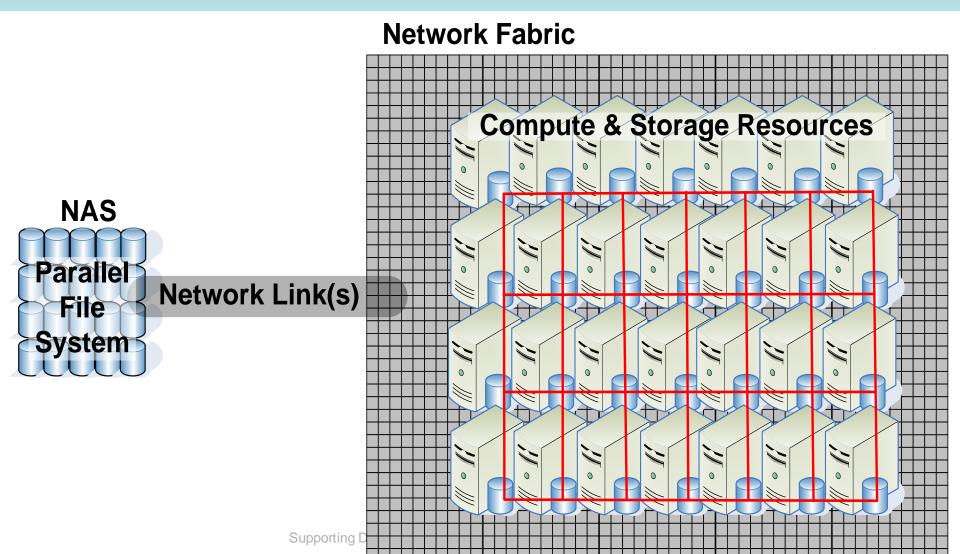
# Exascale Computing Architecture

- Compute
  - 1M nodes, with ~1K threads/cores per node

- Networking
  - N-dimensional torus
  - Meshes

- Storage
  - SANs with spinning disks will replace today's tape
  - SANs with SSDs might exist, replacing today's spinning disk SANs
  - SSDs might exist at every node

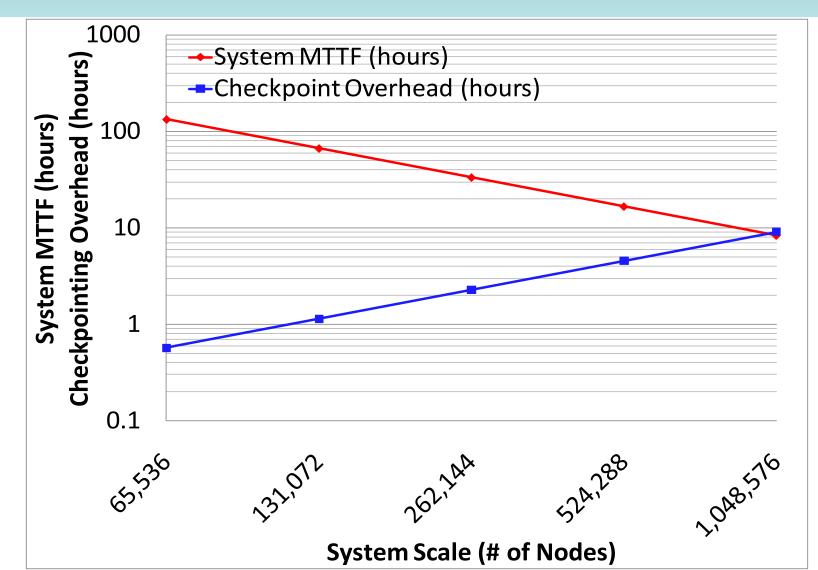# State-of-the-Art Storage Systems in HEC Parallel File Systems

- Segregated storage and compute
  - NFS, GPFS, PVFS, Lustre, Panasas
  - Batch-scheduled Supercomputers
  - Programming paradigms
- Co-located storage
  - BFS, GFS
- Data centers at scale
  - Programming paradigms
  - Others from academia

**NAS**

**Network Link(s)**

**Network Fabric**

**Compute Resources**

# Future Storage System Architecture for Extreme Scale HEC

**Network Fabric**

**Compute & Storage Resources**

**NAS**

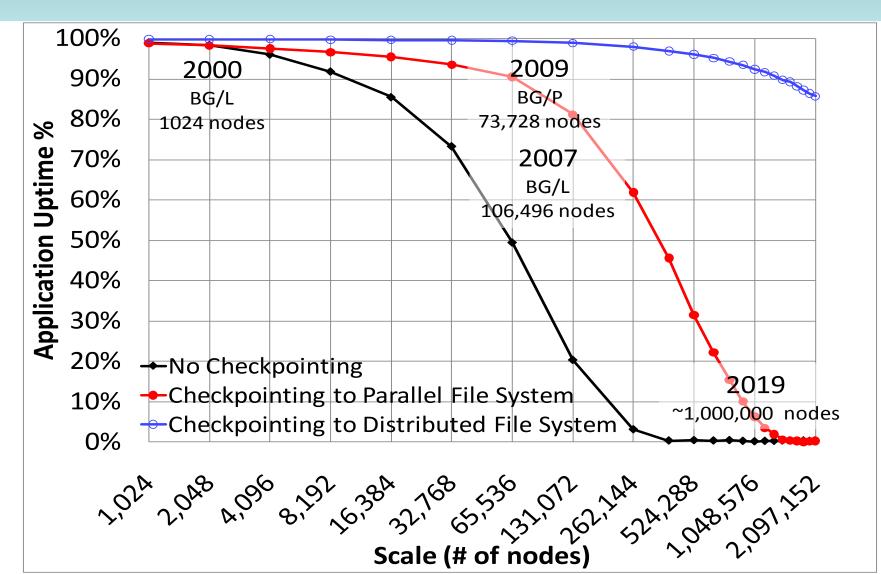**Parallel File System**

**Network Link(s)**

# Some Challenges to Overcome at Exascale Computing

- Programming paradigms
    - HPC is dominated by MPI today
    - Will MPI scale another 3 orders of magnitude?
    - Other paradigms (including loosely coupled ones) might emerge to be more flexible, resilient, and scalable
- Storage systems will need to become more distributed to scale ➜ Critical for resilience of HPC
- Network topology must be used in job management, data management, compilers, etc
- Power efficient compilers and run-time systems

# Expected checkpointing cost and MTTF towards exascale

# Simulation application uptime towards exascale

Application Uptime %

100%
90%
80%
70%
60%
50%
40%
30%
20%
10%
0%

2000
BG/L
1024 nodes

2009
BG/P
73,728 nodes

2007
BG/L
106,496 nodes

2019
~1,000,000 nodes

- No Checkpointing
- Checkpointing to Parallel File System
- Checkpointing to Distributed File System

Scale (# of nodes)

1,024  2,048  4,096  8,192  16,384  32,768  65,536  131,072  262,144  524,288  1,048,576  2,097,152

# Main Message

- ***Decentralization is critical***
  - Computational resource management (e.g. LRMs)
  - Storage systems (e.g. parallel file systems)

- ***Preserving locality is critical!***
  - POSIX I/O on shared/parallel file systems ignore locality
  - Data-aware scheduling coupled with distributed file systems that expose locality is the key to scalability over the next decade

- *Co-locating storage and compute is **GOOD***
  - Leverage the abundance of processing power, bisection bandwidth, and local I/O

Supporting Data-Intensive Distributed Computing in an Exascale Era

# Critical Technologies Needed to achieve Extreme Scales

- Fundamental Building Blocks (with a variety of resilience and consistency models)
  - Distributed hash tables (aka NoSQL data stores)
  - Distributed Message Queues
- Deliver future generation distributed systems
  - Global File Systems, Metadata, and Storage
  - Job Management Systems
  - Workflow Systems
  - Monitoring Systems
  - Provenance Systems
  - Data Indexing

# More Information

- More information:
  - http://www.cs.iit.edu/~iraicu/
  - http://datasys.cs.iit.edu/
- Contact:
  - iraicu@cs.iit.edu
- Questions?