

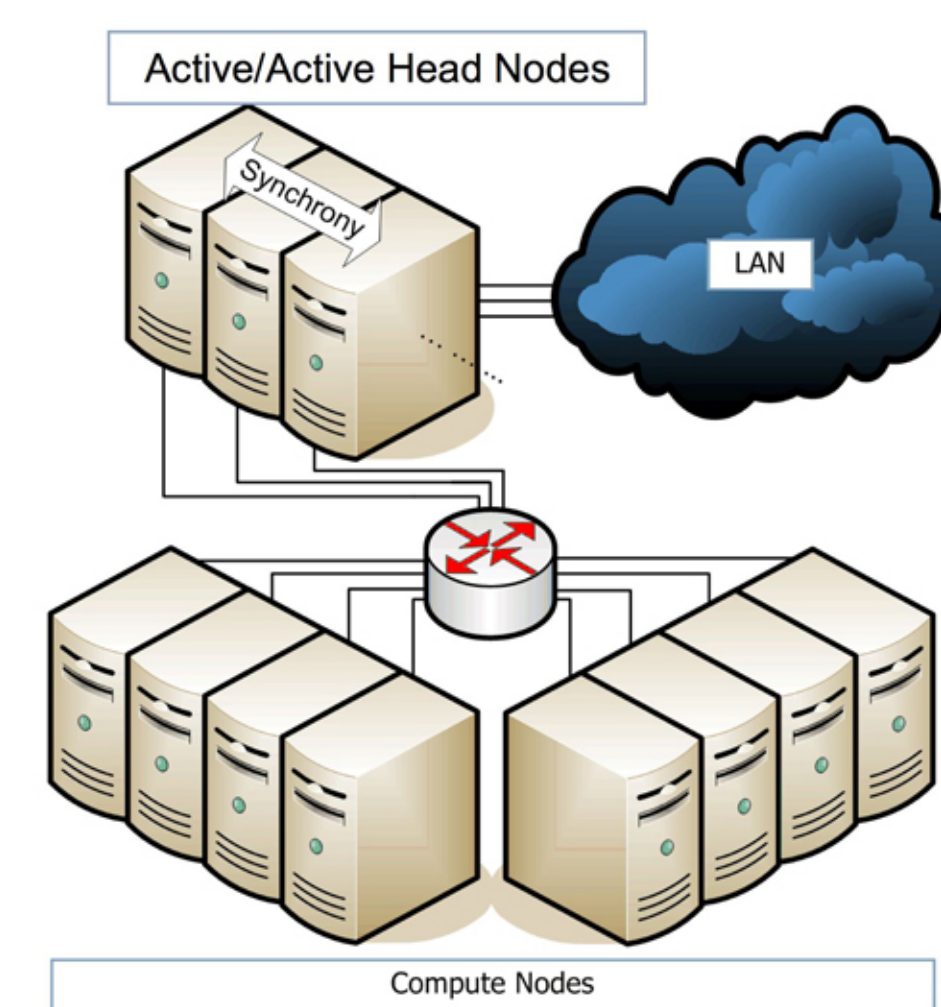
# Resiliency for High-Performance Computing Systems

Stephen L. Scott, Christian Engelmann, Hong H. Ong, Geoffroy R. Vallée, Thomas Naughton, Anand Tikotekar, George Ostrouchov (Oak Ridge National Laboratory)  
 Chokchai Leangsuksun, Nichamon Naksinehaboon, Raja Nassar, Mihaela Paun (Louisiana Tech University)  
 Frank Mueller, Chao Wang, Arun Nagarajan, Jyothish Varma (North Carolina State University)  
 Xubin He, Li Ou, Xin Chen (Tennessee Technological University)

## Research and development goals

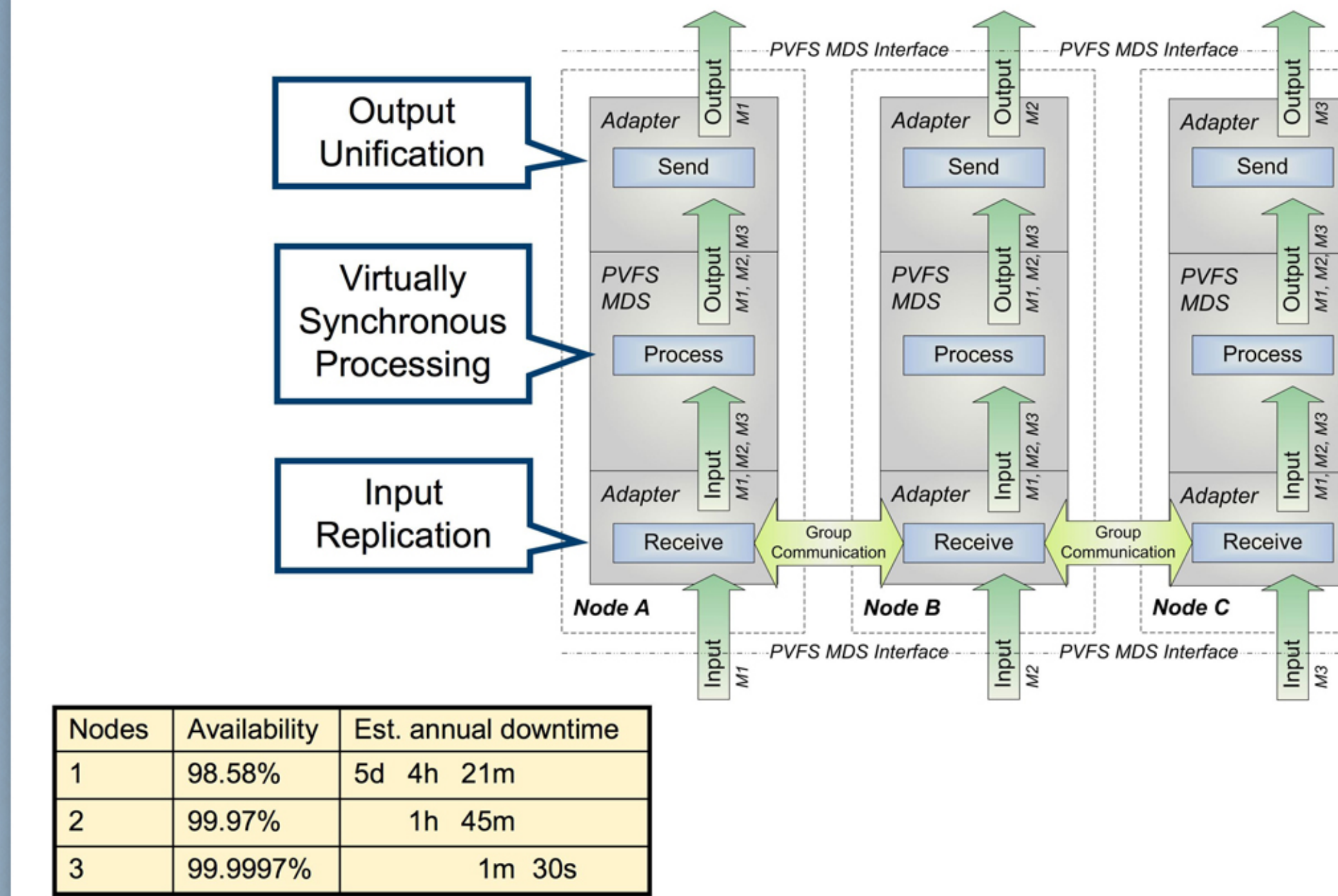
- Efficient redundancy strategies for head and service nodes in HPC systems to provide high availability as well as high performance of critical infrastructure services
- Reactive fault tolerance for HPC compute nodes utilizing the job pause approach as well as checkpoint interval and placement adaptation to actual and predicted system health threats
- Proactive fault tolerance using system-level virtualization in HPC environments for preemptive migration of computation away from compute nodes that are about to fail
- Reliability analysis for identifying pre-fault indicators, predicting failures, and modeling and monitoring of individual component and overall HPC system reliability
- Holistic fault tolerance technology through combination of adaptive proactive and reactive fault tolerance mechanisms in conjunction with system health monitoring and reliability analysis

## Symmetric active/active redundancy for head and service nodes



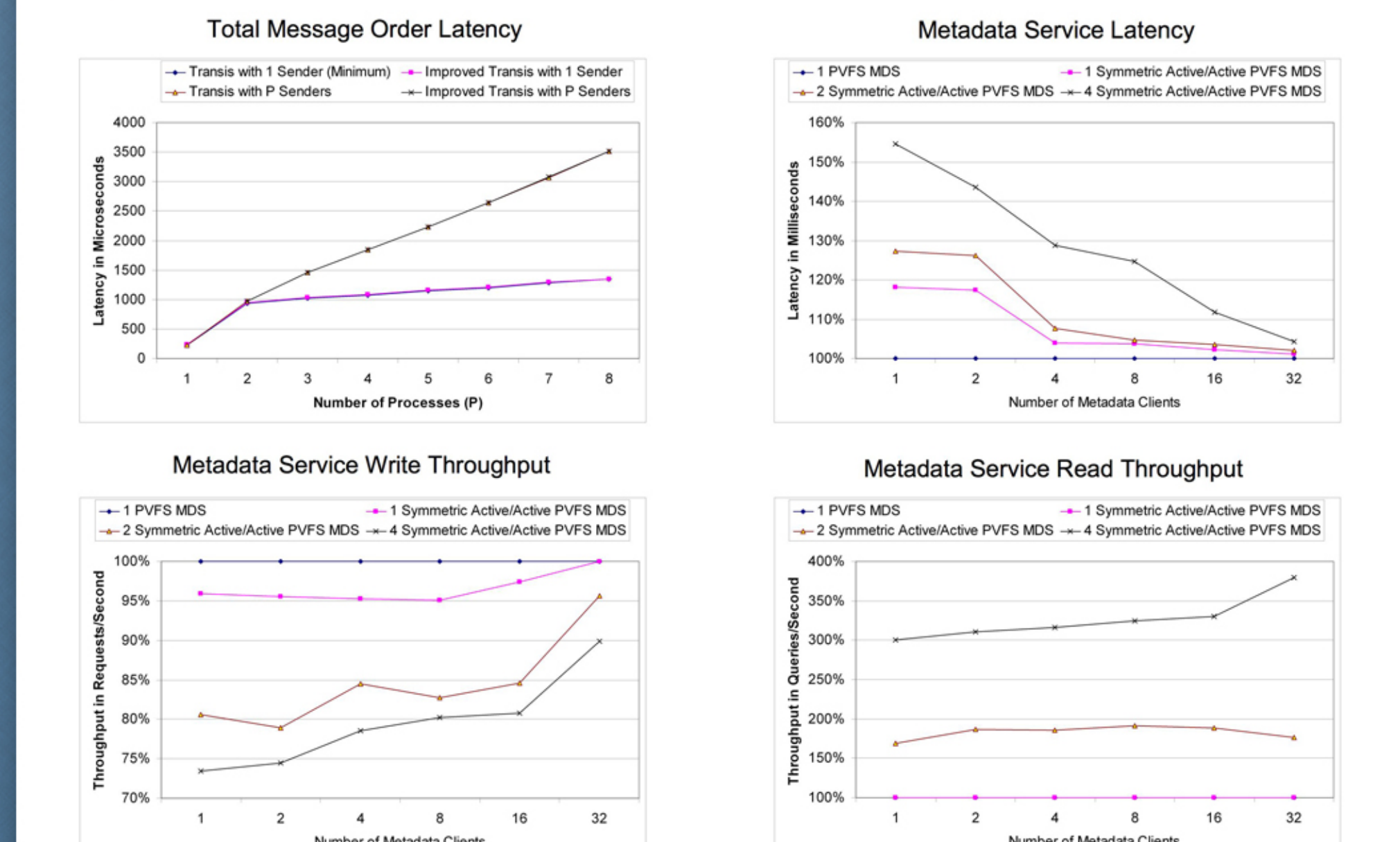
- Many active head nodes
  - Work load distribution
  - Symmetric replication between head nodes
  - Continuous service
  - Always up to date
  - No fail-over necessary
  - No restore-over necessary
  - Virtual synchrony model
  - Complex algorithms
- Prototypes for PBS Torque and Parallel Virtual File System metadata server

## Symmetric active/active replication



Nodes	Availability	Est. annual downtime
1	98.58%	5d 4h 21m
2	99.97%	1h 45m
3	99.9997%	1m 30s

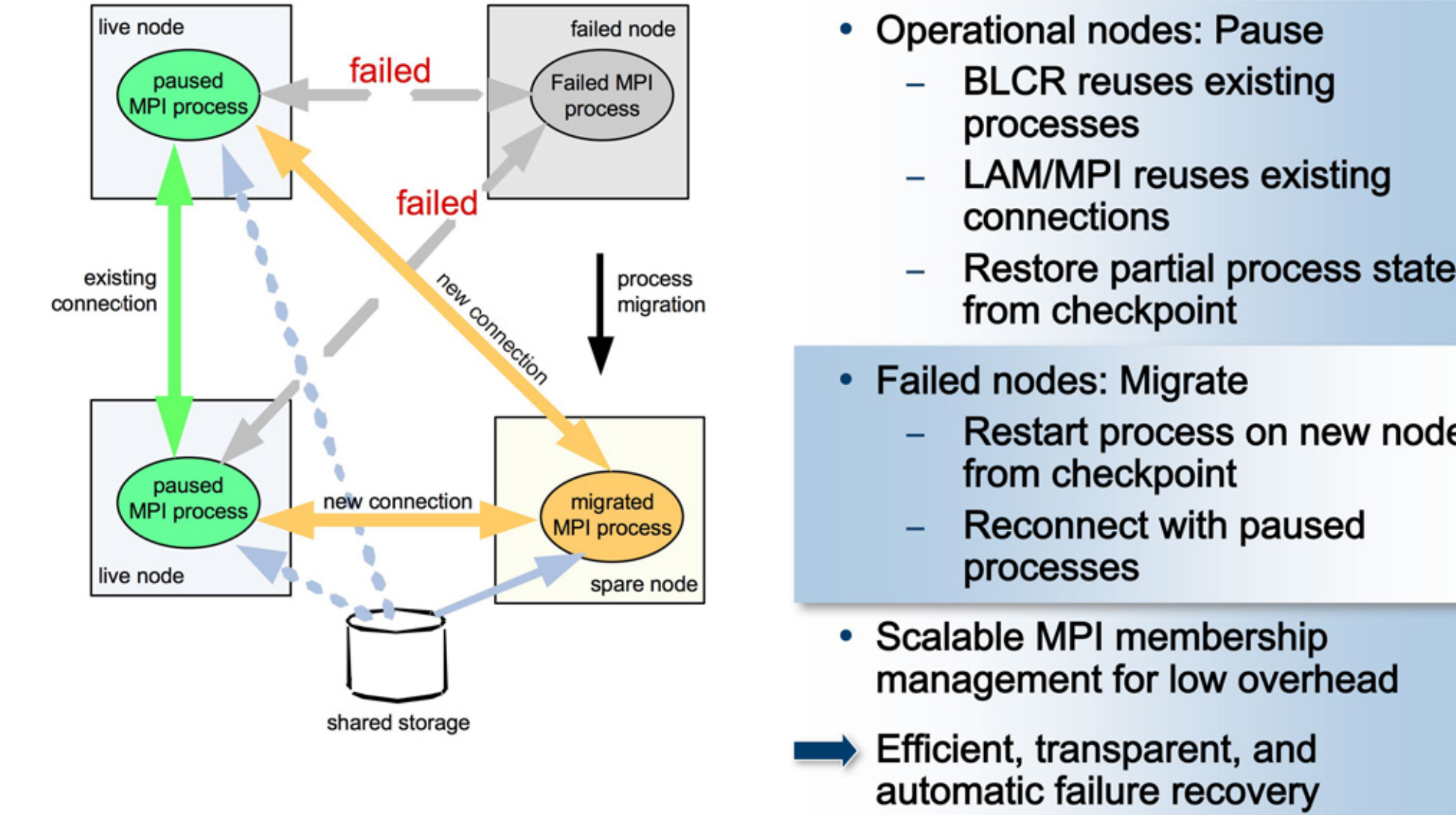
## Symmetric active/active Parallel Virtual File System metadata service



## Reactive vs. proactive fault tolerance for compute nodes

- Reactive fault tolerance:
  - State saving during failure-free operation
  - State recovery after failure
  - Assured quality of service, but limited scalability
- Proactive fault tolerance:
  - System health monitoring and online reliability modeling
  - Failure anticipation and prevention through prediction and reconfiguration before failure
  - Highly scalable, but not all failures can be anticipated
- Ideal solution: Matching combination of both

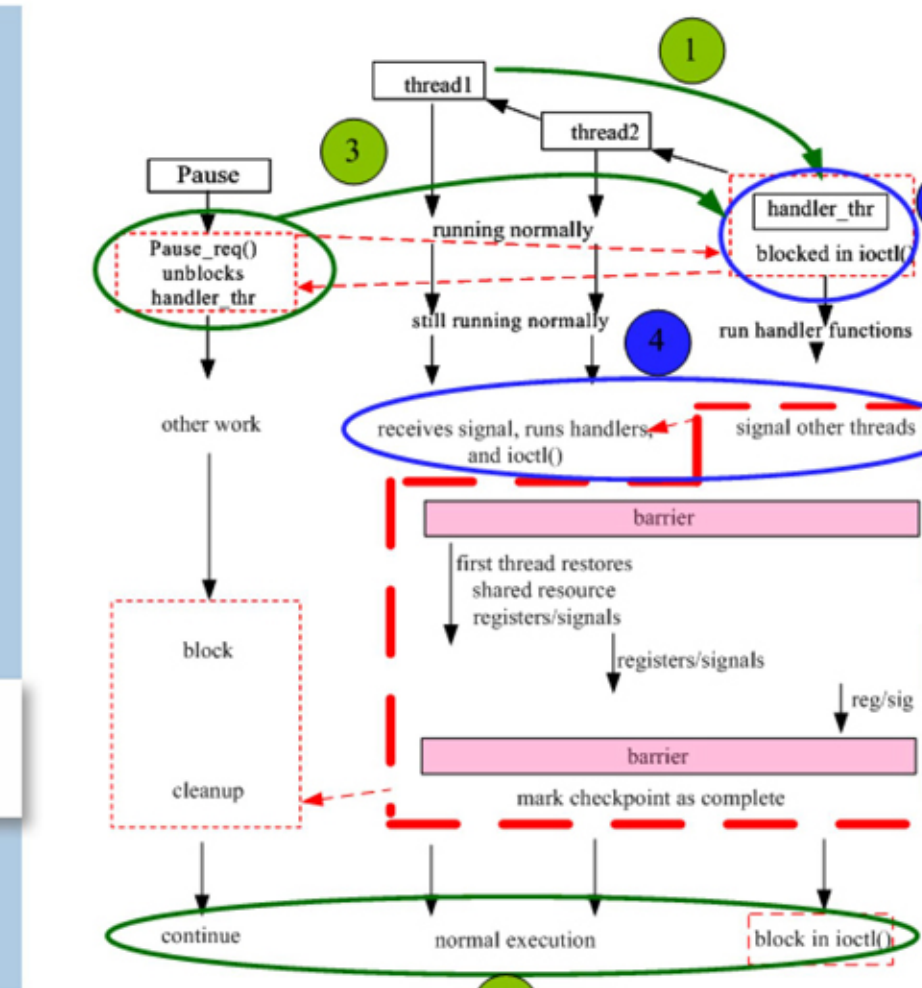
## Enhanced reactive fault tolerance with LAM/MPI+BLCR job pause mechanism



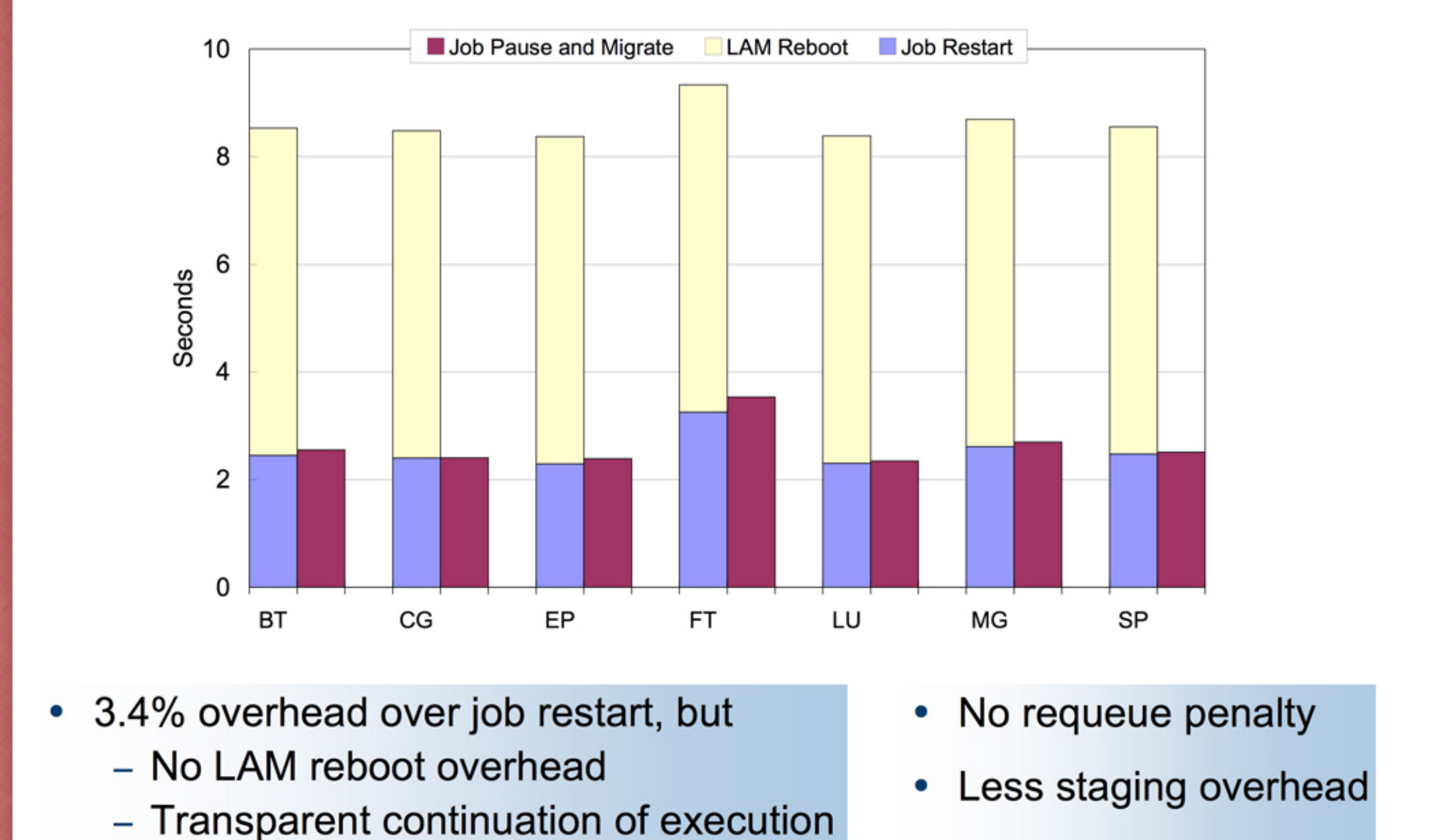
- Operational nodes: Pause
  - BLCR reuses existing processes
  - LAM/MPI reuses existing connections
  - Restore partial process state from checkpoint
- Failed nodes: Migrate
  - Restart process on new node from checkpoint
  - Reconnect with paused processes
- Scalable MPI membership management for low overhead
- Efficient, transparent, and automatic failure recovery

## New job pause mechanism in BLCR

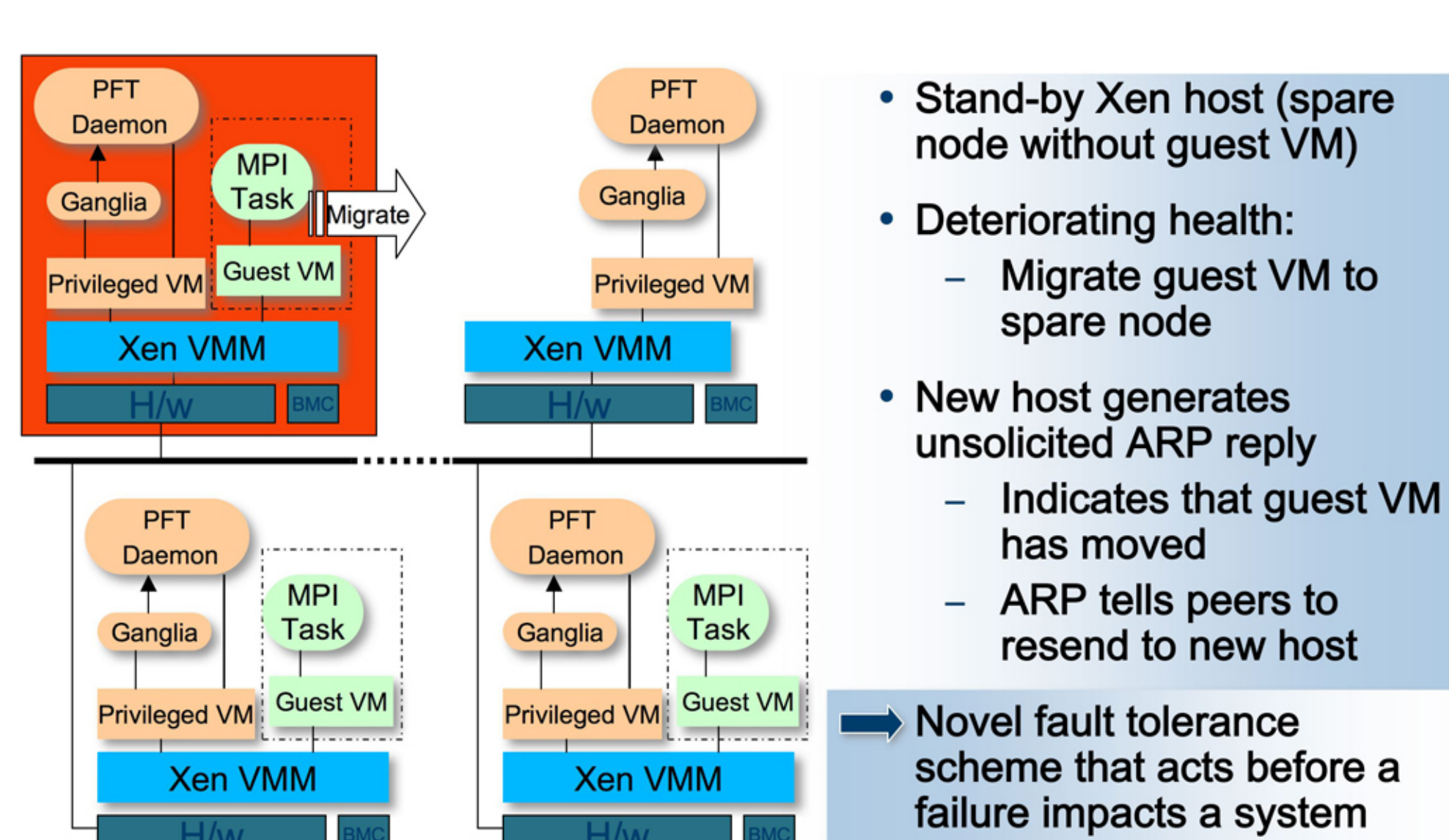
- Application registers threaded callback → spawns callback thread
- Thread blocks in kernel
- Pause utility calls ioctl(), unblocks callback thread
- All threads complete callbacks and enter kernel
- New: All threads restore part of their states
- Run regular application code from restored state



## LAM/MPI+BLCR job pause performance

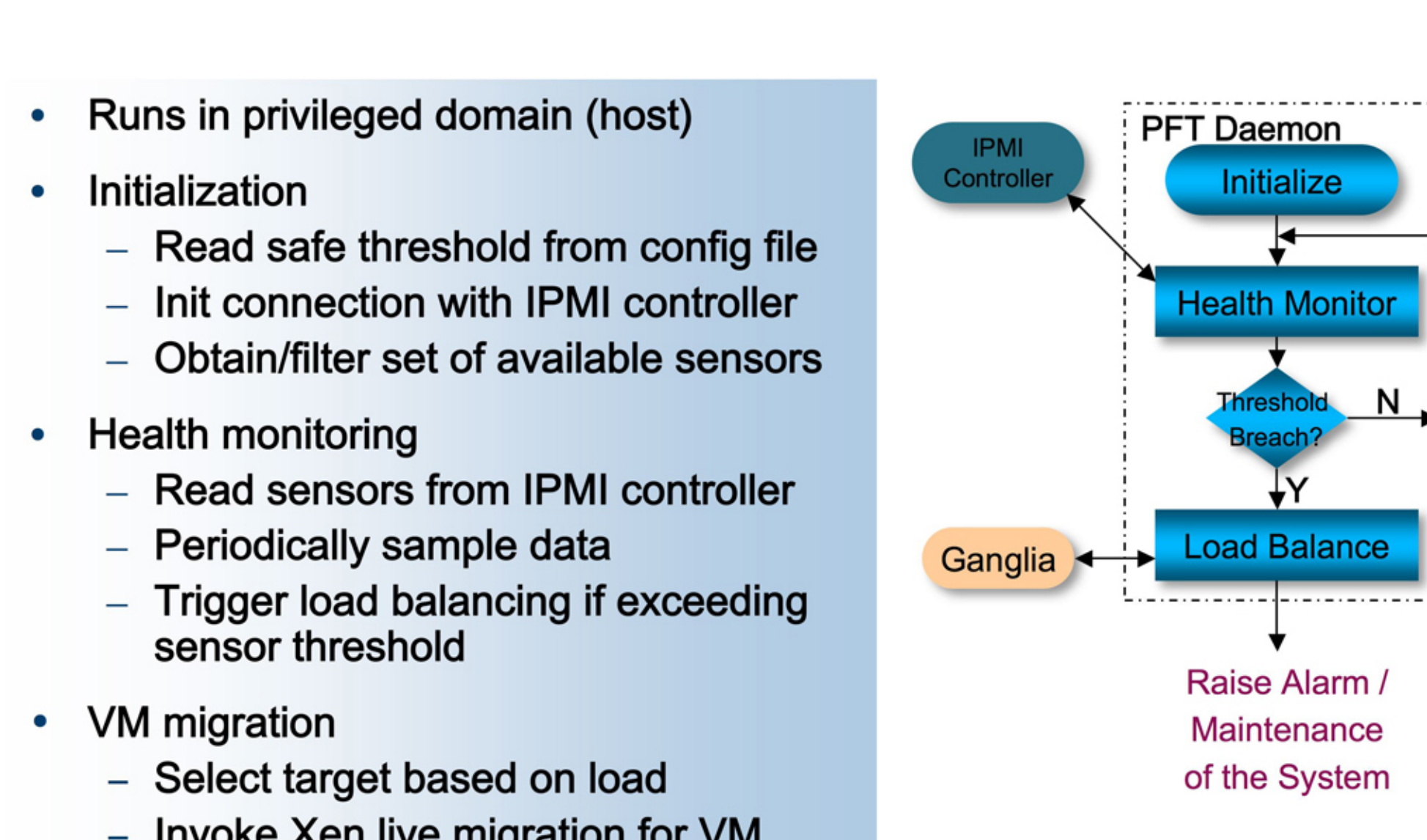


## Proactive fault tolerance using Xen virtualization



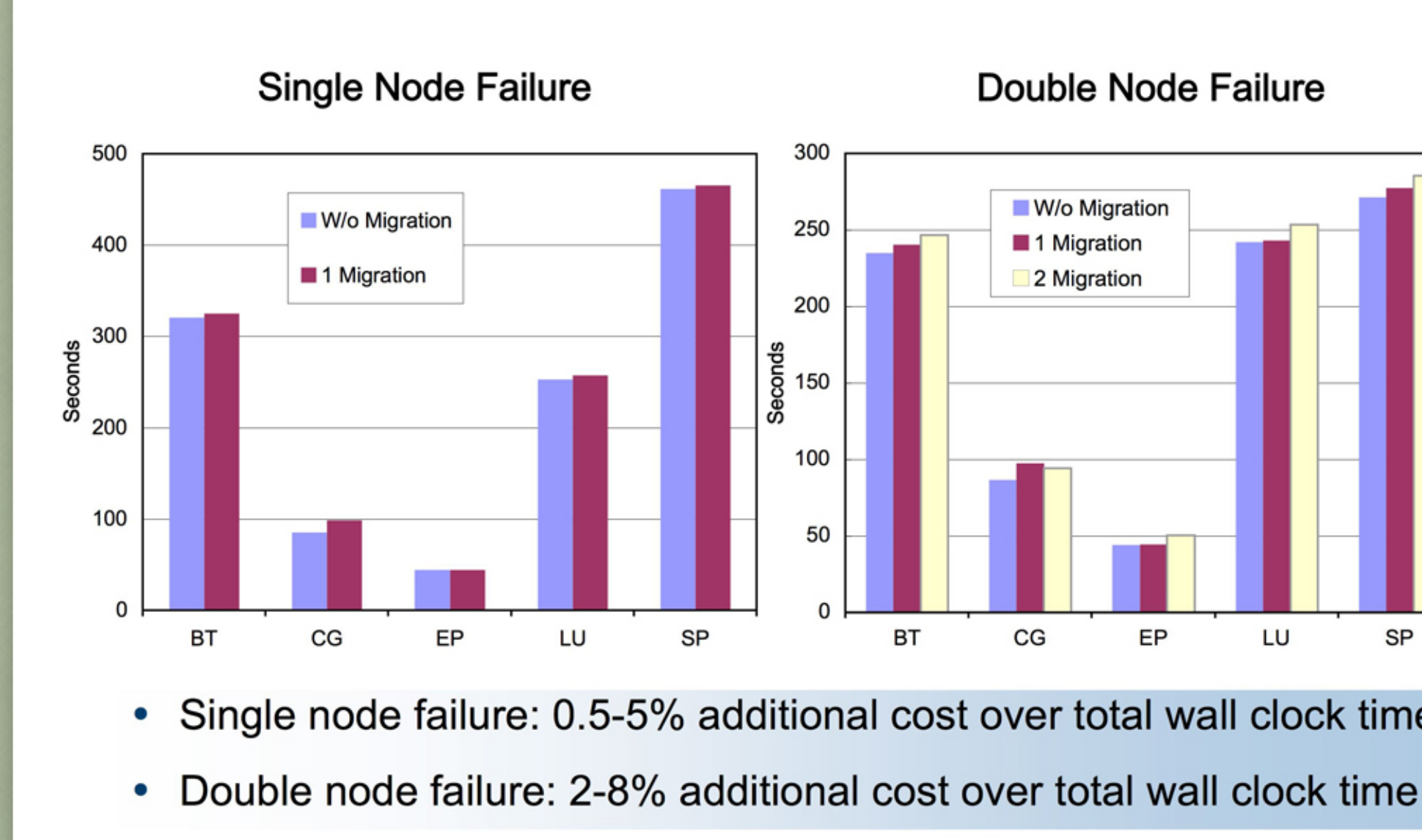
- Stand-by Xen host (spare node without guest VM)
- Deteriorating health:
  - Migrate guest VM to spare node
- New host generates unsolicited ARP reply
  - Indicates that guest VM has moved
  - ARP tells peers to resend to new host
- Novel fault tolerance scheme that acts before a failure impacts a system

## Proactive fault tolerance daemon

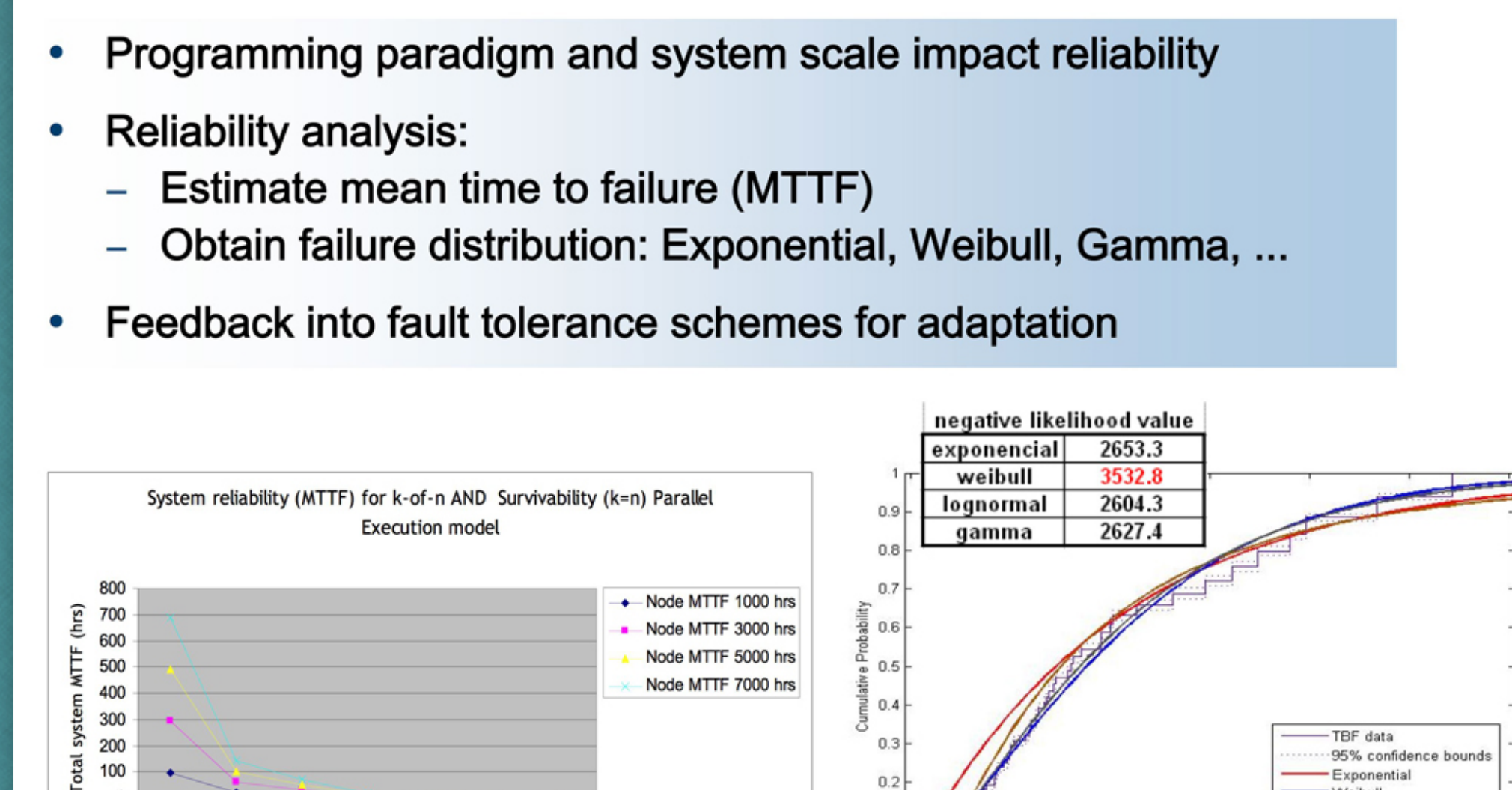


- Runs in privileged domain (host)
- Initialization
  - Read safe threshold from config file
  - Init connection with IPMI controller
  - Obtain/filter set of available sensors
- Health monitoring
  - Read sensors from IPMI controller
  - Periodically sample data
  - Trigger load balancing if exceeding sensor threshold
- VM migration
  - Select target based on load
  - Invoke Xen live migration for VM

## VM migration performance impact

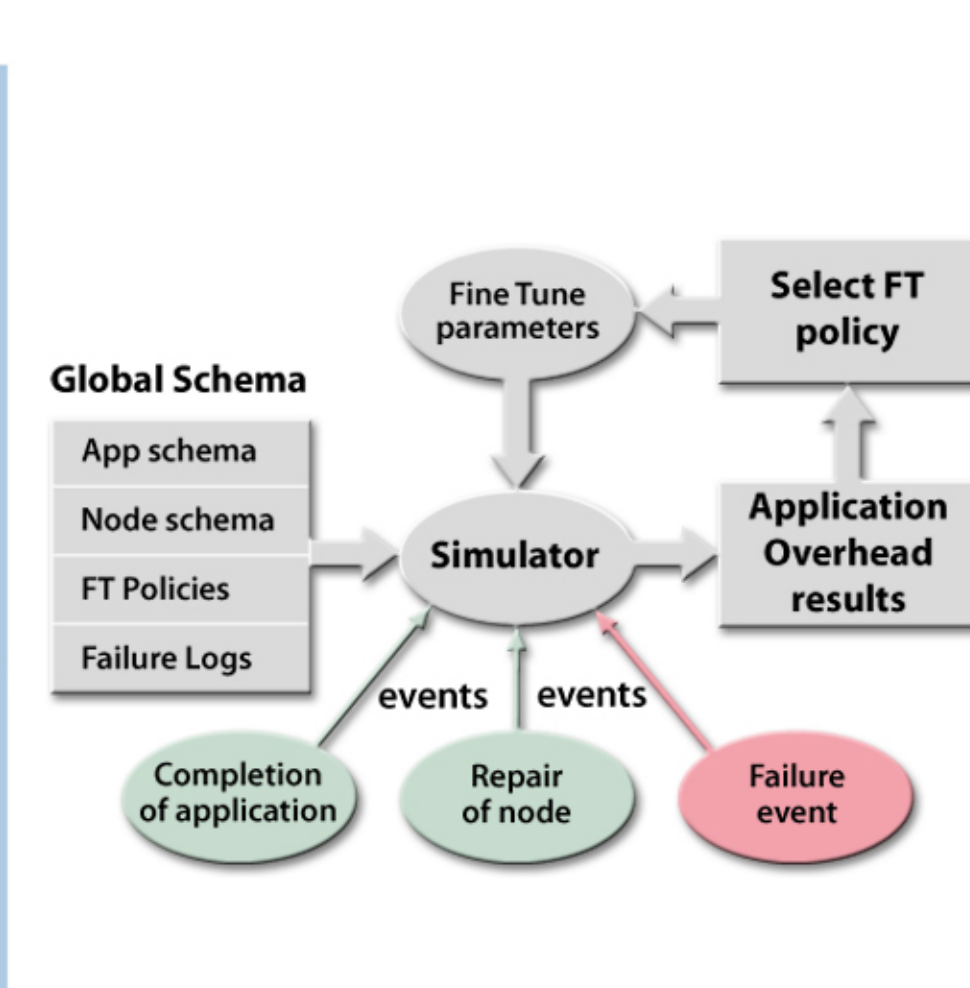


## HPC reliability analysis and modeling for prediction and anticipation

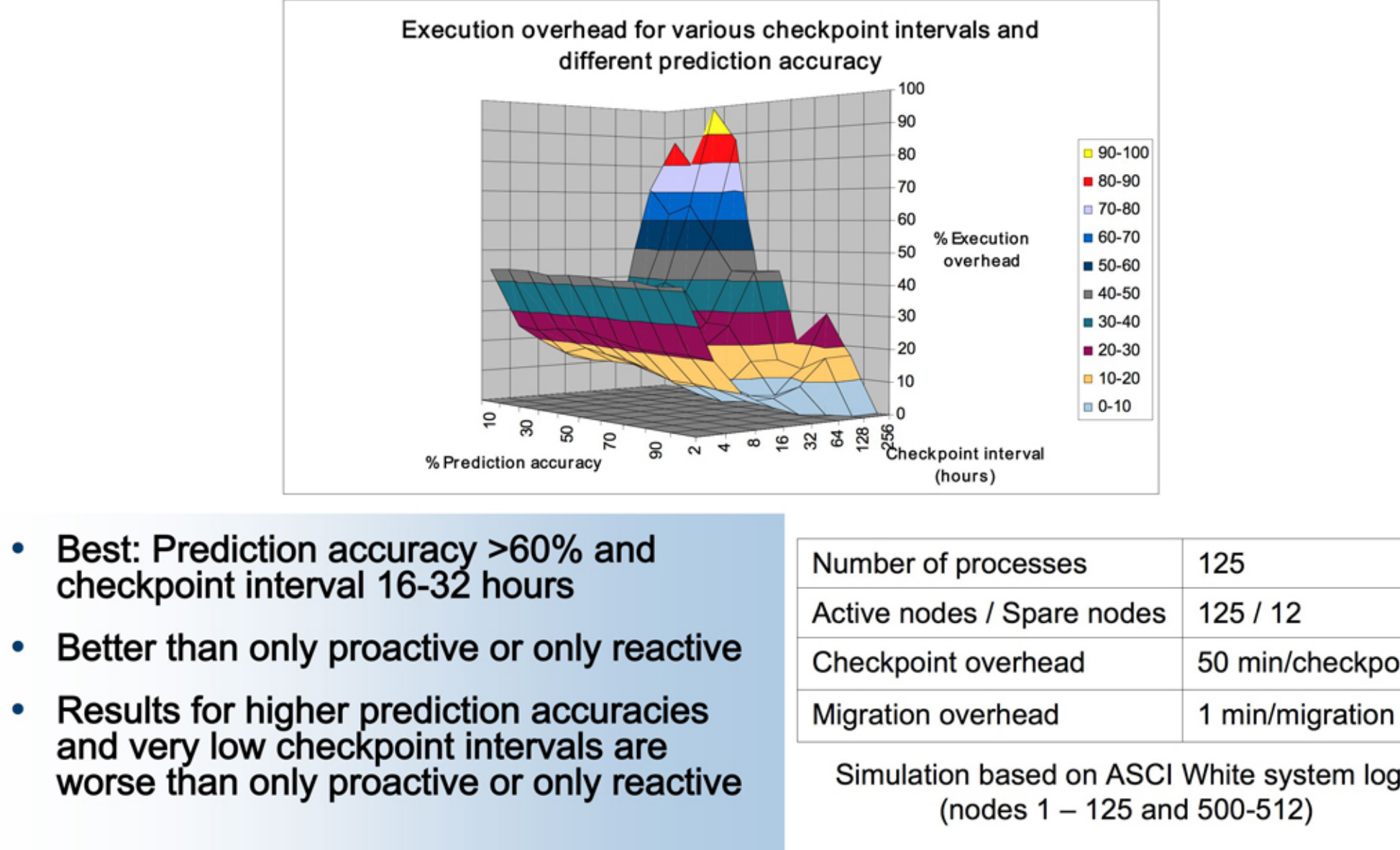


## Simulation framework for HPC fault tolerance policies

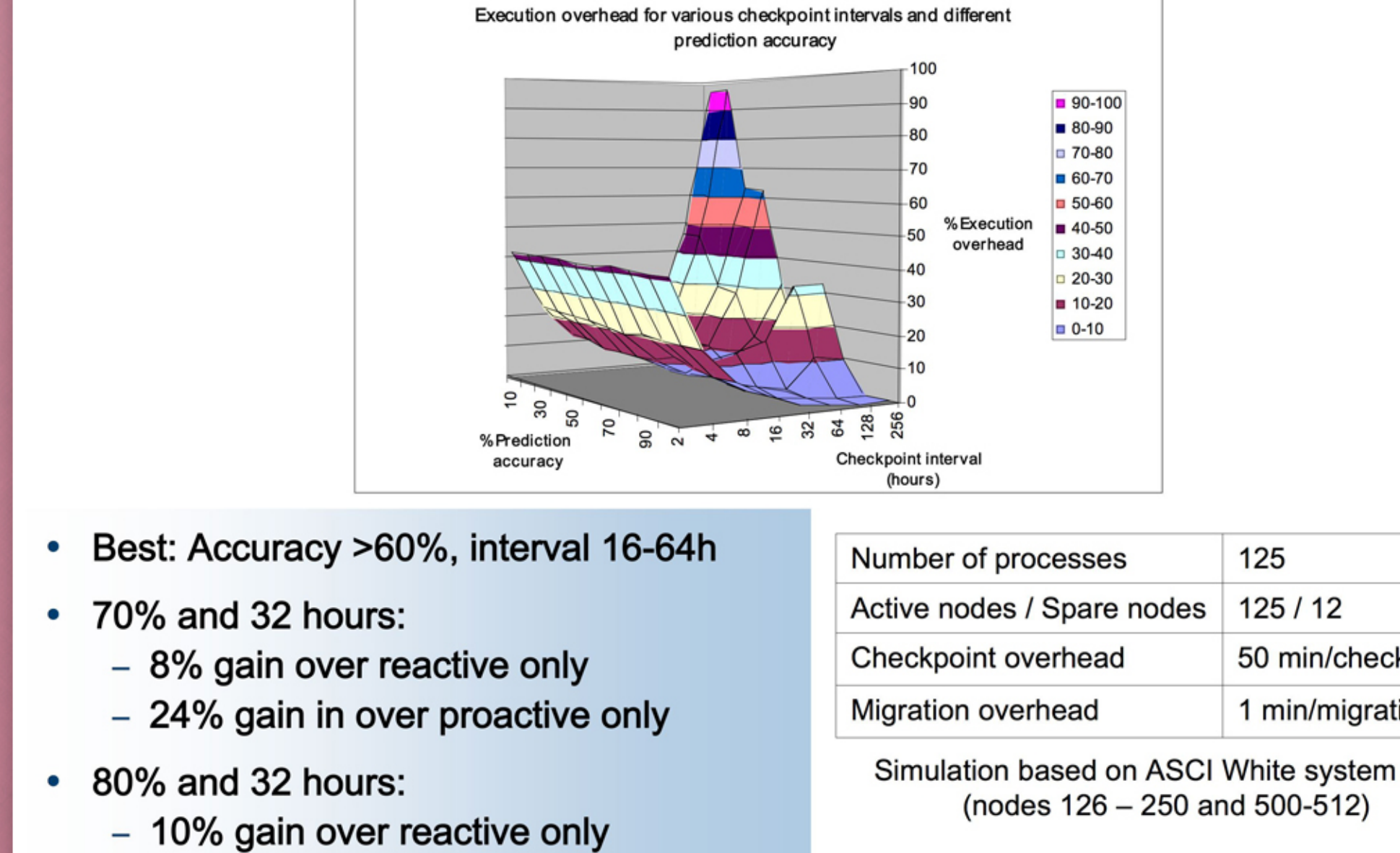
- Evaluation of fault tolerance policies
  - Reactive only
  - Proactive only
  - Reactive/proactive combination
- Evaluation of fault tolerance parameters
  - Checkpoint interval
  - Prediction accuracy
- Event-based simulation framework using actual HPC system logs
- Customizable simulated environment
  - Number of active and spare nodes
  - Checkpoint and migration overheads



## Combination of proactive and reactive fault tolerance: Simulation example 1



## Combination of proactive and reactive fault tolerance: Simulation example 2



## A holistic resiliency framework for high-performance computing

