



STORAGE DEVELOPER CONFERENCE

SNIA ■ SANTA CLARA, 2015

**System Verification At Scale:**

**Thousands of Users**

**Do you need to test with them or not?**

**Steven Buller**

**Julian Cachua**

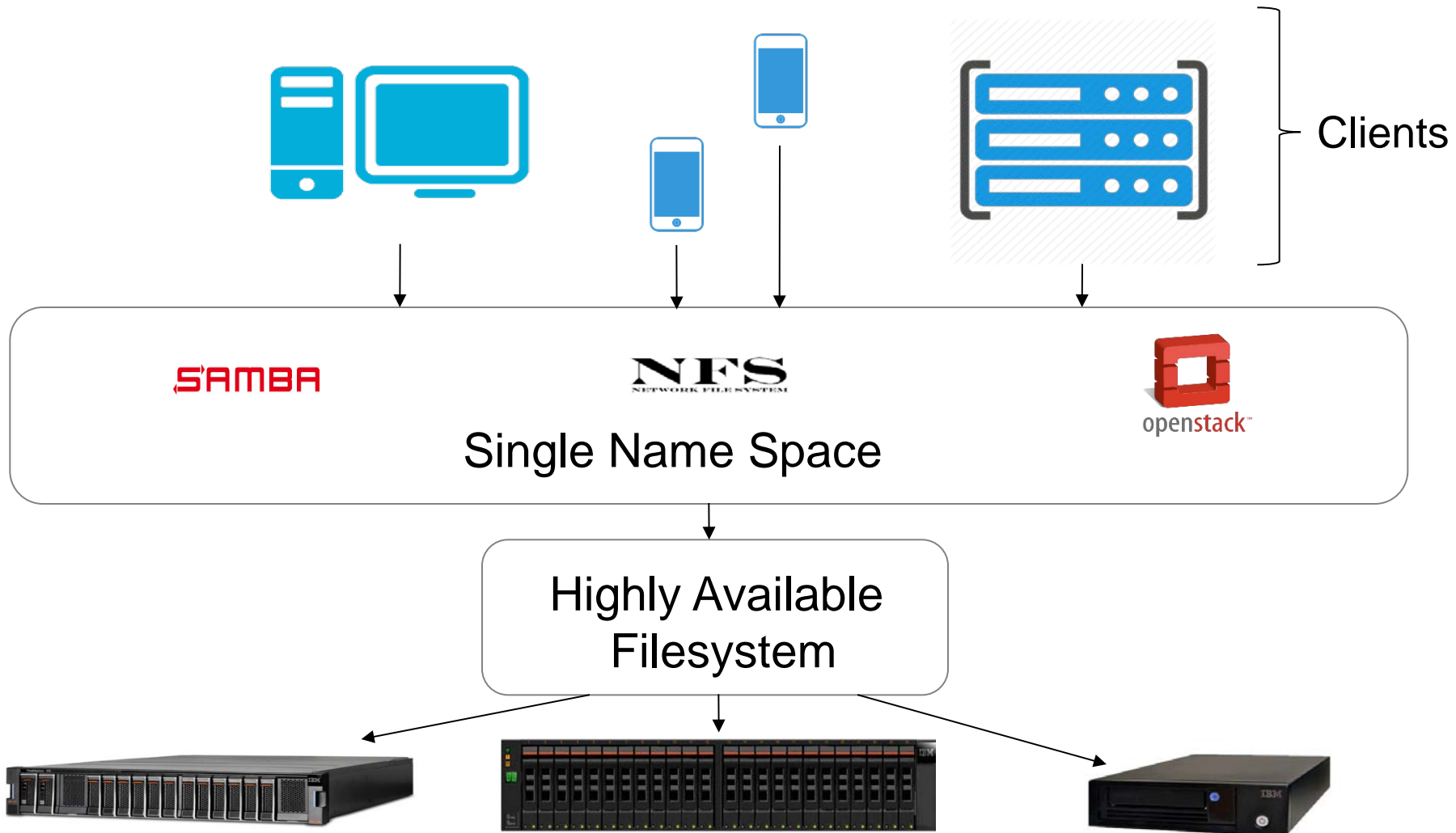
**Christina Lara**

**IBM**

# AGENDA

- ❑ WHY thousands of users
- ❑ HOW to test & simulate at a scaled level
- ❑ RESULTS from testing at scale
- ❑ FUTURE test improvements

# System Overview



# Complex interactions between system components

- ❑ Robustness of the system being able to handle thousands of clients/connections
- ❑ Required interaction between multiple nodes and software services
  - ❑ Health monitoring sensors to detect node status and properly handle failovers
  - ❑ IP distribution mechanism to evenly re-distribute IPs
- ❑ Multiple nodes can have direct read/write access to the same data
  - ❑ R/W conflicts must be avoided
  - ❑ Diverse protocols and the underlying filesystem must be able to access and modify data in a coordinated way

# Why thousands of users?

- ❑ Diverse environments requiring thousands of clients usage
  - ❑ NAS servers handling thousands of clients are becoming a common requirement in diverse environments ranging from public school systems to cancer research.
- ❑ Simulation of thousands of users prior to customer use
  - ❑ The need to simulate tens of thousands of active users on SMB, NFS, and Object protocols is the first step to finding issues before the customer does.

# Field problems that drove scalability testing:

- ❑ “Critical Situation” events being seen due to large number of users/clients
  - ❑ Examples of critical events:
    - ❑ SMB connection issues
      - ❑ Cleanup after connections are closed
    - ❑ Resource contention related to high number of connections
      - ❑ Database lock conflicts
    - ❑ Deadlocks during concurrent system software upgrades
      - ❑ IP management failures

# Actual field scalability issues:

- ❑ A public school system with 30,000+ students
  - ❑ First week of school, parents are logging in to see what Jack and Jill are going to have for lunch...
  - ❑ Day 2: Customer encountered a bug with CTDB not deleting closed connections.
    - ❑ Which led to >90K entries, which maxed out a DB and led to an outage.
  - ❑ RESULT> Drove investment to test with thousands of connections.
- ❑ Large research consortium:
  - ❑ One researcher was deleting a PIT (Point-In-Time copy) while another user was searching that copy → deadlocked.
  - ❑ Frequent PITs of bulk data starved other system resources.
  - ❑ RESULT> Creation of a Education Workload on a dedicated test stand to replicate the filesystem directory structure and function frequency characteristics.

# How Did We Test Thousands of Users?



# How: Testing Solutions

- ❑ Tools
- ❑ Examples of tests
- ❑ Examples of issues found at scale

# Testing Solution: Tools

- ❑ Scalability with existing hardware
  - ❑ 300-500 physical clients, even multiplied by VM's, is not sufficient.
  - ❑ More physical clients was cost prohibitive.
- ❑ Introduction of LoadDynamix\* tools
  - ❑ Simulate connections from multiple IPs and MAC addresses.
  - ❑ Able to coordinate data access among diverse clients:
    1. Create by ClientA on NodeY with SMB.
    2. Read by ClientB on NodeX with NFS.
    3. Delete by ClientC on NodeZ.

\* Disclaimer – This is not an endorsement for any one vendor

# Testing Solutions: Targeted Workloads

## □ Industry Representative Workloads

### □ Bio-Genetics

- SMB Genomic sequencers write data to be analyzed by NFS apps reads.

### □ Technology

- EDA tools generate “kabillions” (lots) of tiny temp files, access contention.

### □ Education

- Tens of thousands of subdirectories of varying depths.

## □ Workloads differentiated by:

### □ Read/Write %

### □ File Sizes

### □ Multi-protocol & Cross-protocol interaction

### □ Quantities of Directories & Connections

### □ Hardware configurations

- Size of cluster-number of nodes
- Amount of space (inodes in the filesystem) used.

# Education Workload: Directory Structure

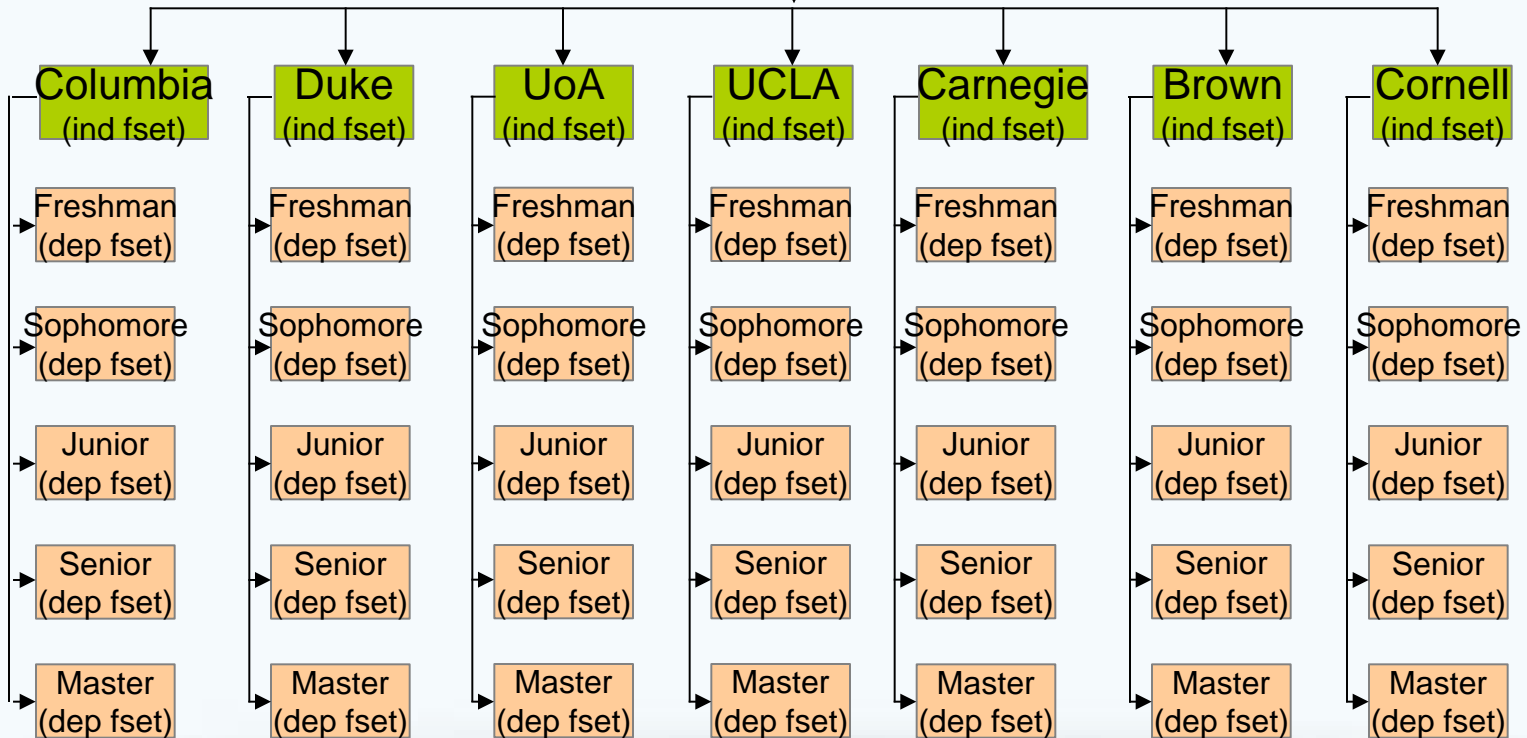
Single Filesystem

- PIT copies every 4 hours
- Nightly backup
- Hourly remote replication sync.
- Other IBM features X times per day.

Root Export

LoadDynamix

\*1K folders under each “grade level” fileset  
\*10K files for each fileset, distributed in all 1K folders



# Testing Solution: Focused Test Coverage

- ❑ CTDB / Samba – Connection Management
  - ❑ The CTDB testplan used to cover un-exercised areas
    - ❑ High # of connections / open files
    - ❑ Multi-node file access
    - ❑ Multiple sub-nets
    - ❑ Failover & Recovery
    - ❑ Vacuuming (expensive process)
      - ❑ CTDB records stress
      - ❑ Frozen DBs
    - ❑ Network Operations – IP banning & modifications

	Same Sockets	Traverse	Vacuuming	NFS Recoveries	Ban Ips	CTDB status / heavy load	Max number of nfs threads	Unbalan ced Ips	Frozen Dbs	Various subnets	M ...
CIFS NFS 02		partial				partial					
CIFS03		partial				partial					
CIFS04		partial				partial					

# Testing Solution: Additional Function Interaction

## ❑ IBM Function Integration

### ❑ Point In Time (PIT) Copy

- ❑ Ramp up frequency of Creation/Deletion and overlap timing to create contention

### ❑ Backup/Archiving

- ❑ MacOS SMB clients “surfing” for a file would trigger mass file recalls from off-line copies.

### ❑ Failovers

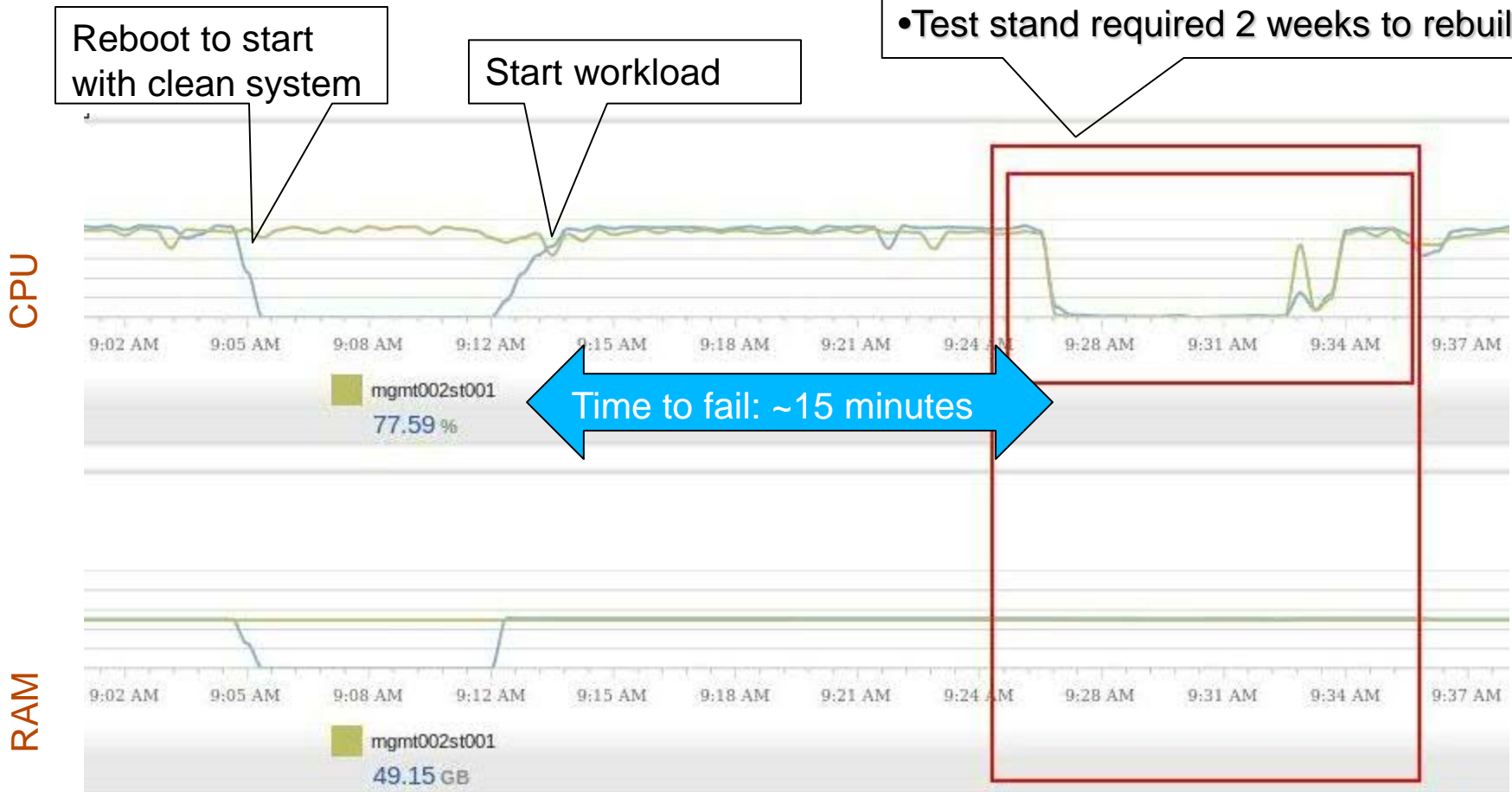
- ❑ Move service to another node in the cluster before client time-out expires.

# Testing Solution: Issues Found at Scale

- ❑ “Thundering Herd” lock contention resulting in mutex changes (multiple connections to a single file)
  - ❑ See: [Samba Mutex Exchange Presentation](#)
- ❑ CTDB Vacuuming/Recovery Master election
  - ❑ 2500 Connections PER NODE
  - ❑ Had to increase minimum recommended memory to solve.
- ❑ “High” number of IP addresses:
  - ❑ Targeting 900 IP addresses resulted in a loss of network access that persisted after cluster reboot.
  - ❑ Had to delete IP’s out of CTDB to recover.
- ❑ Memory Leaks

# Testing Solutions: Result

Catastrophic failure:  
• 2 defects in the underlying operating system.  
• IP address failover distribution defects.  
• Test stand required 2 weeks to rebuild.





# Testing Solution: Results

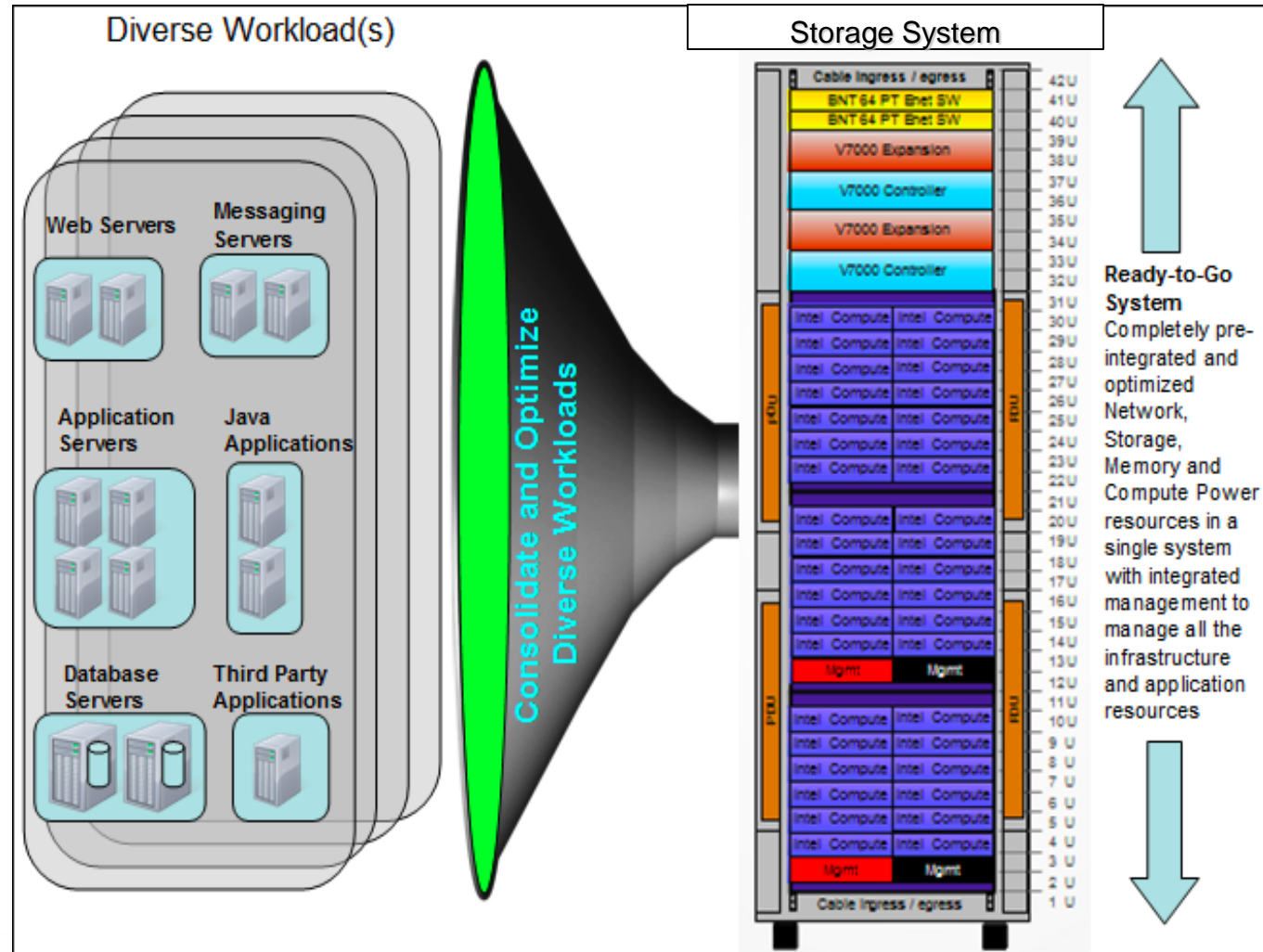
- ❑ Hardware Boundaries
  - ❑ Defined memory & CPU minimum requirements
    - ❑ Even virtualized resources are not limitless.
- ❑ Published Limitations & Best Practices
  - ❑ Set appropriate expectations
  - ❑ Guidance for sales force in sizing systems.
- ❑ Redesign of Locking Mechanism between Samba & CTDB
  - ❑ F-CTRL to Mutex
- ❑ Introduction of defensive Spectrum Scale deadlock monitoring as a new product feature.
  - ❑ Proactive monitoring of conditions that could lead to an outage.

# Introduction to Workloads

- ❑ Workloads to simulate virtual clients (ie: LoadDynamix)
  - ❑ Characteristics
  - ❑ Connections
- ❑ Workloads to simulate industry
  - ❑ Bio-Genetics
  - ❑ Technology
  - ❑ Education
- ❑ Workload Analysis Flow

# Workloads Paradigm

- Each Industry has different requirements.
- Industry characteristics can affect different elements of the storage system.
- Standard hardware storage systems must be able to manage different Industries workloads



# Differences between workloads

Industry	NFS	SMB	Cross Protocol	Locking		Nested Dirs		Dirs Number		Read Sizes		Write Sizes		Read / Write		NFS	SMB	get file attr	set file attr	ren	del	dir creates	dir del	
				NFS (only NFSv4)	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	creates	creates	file	file			
Technology	60%	40%	NO	20%	30%	1	10	0.2/conn	0.2/conn	33% <1K	33% <1K	33% <512K (50% COW) 90% file sync <1K 50% same offset	33% <1K (1M files)	33% <1K (1M files)	50/50	50/50	0.3/sec/conn	0.3/sec/conn	20X / file	0.05X / file	5%	15%	2/sec	2/sec
	NFSv4 50%	SMB1 10%		Create: unchecked 5%, guarded 55%, exclusive 40%								33% 32K (10M files)	33% 32K (25% COW)											
	NFSv3 50%	SMB2 40%											33% 32K (10M files)	33% 32K (25% COW)										

<Protocol & Version> Locking<Characteristics>  
<protocol & version> Main <write size>K  
<protocol & version> Directory Metadata  
In all scenarios

Industry	NFS	SMB	Cross Protocol	Locking		Nested Dirs		Dirs Number		Read Sizes		Write Sizes		Read / Write		NFS	SMB	get file attr	set file attr	ren	del	dir creates	dir del
				NFS (only NFSv4)	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	NFS	SMB	creates	creates	file	file		
Education	10%	90%	NO	20%	20%	5	10	1/conn	0.5/conn	50% <1K	50% <1K	50% <1K (50% COW)	50% <1K (50% COW)	50/50	50/50	0.3/sec/conn	0.3/sec/conn	10X / file	0.1X / file	10%	10%	0.5/conn	10%
	NFSv4 50%	SMB1 10%										50% 1K - 10M	50% 1K - 10M										
	NFSv3 50%	SMB2 40%										50% 1K - 10M	50% 1K - 10M										



# FUTURE: To Infinity and Beyond....

## ❑ Ultimate: tcpreplay at scale

- ❑ Record a customer's traffic and play it back from thousands of clients

- ❑ Challenges:

- ❑ Need to scrub/filter sensitive data (best not to capture it)

- ❑ Huge capture sizes

- ❑ Network replication

- ❑ Still need documentation of the customers environment

- ❑ Ex. Directory structure can not be reverse engineered from traces.

# FUTURE: Near term (This side of Infinity)

## ❑ Realistic Intermediate step: Workload Analysis

- ❑ Use Big Data Analytics to generate synthetic workloads based on trace characteristics.

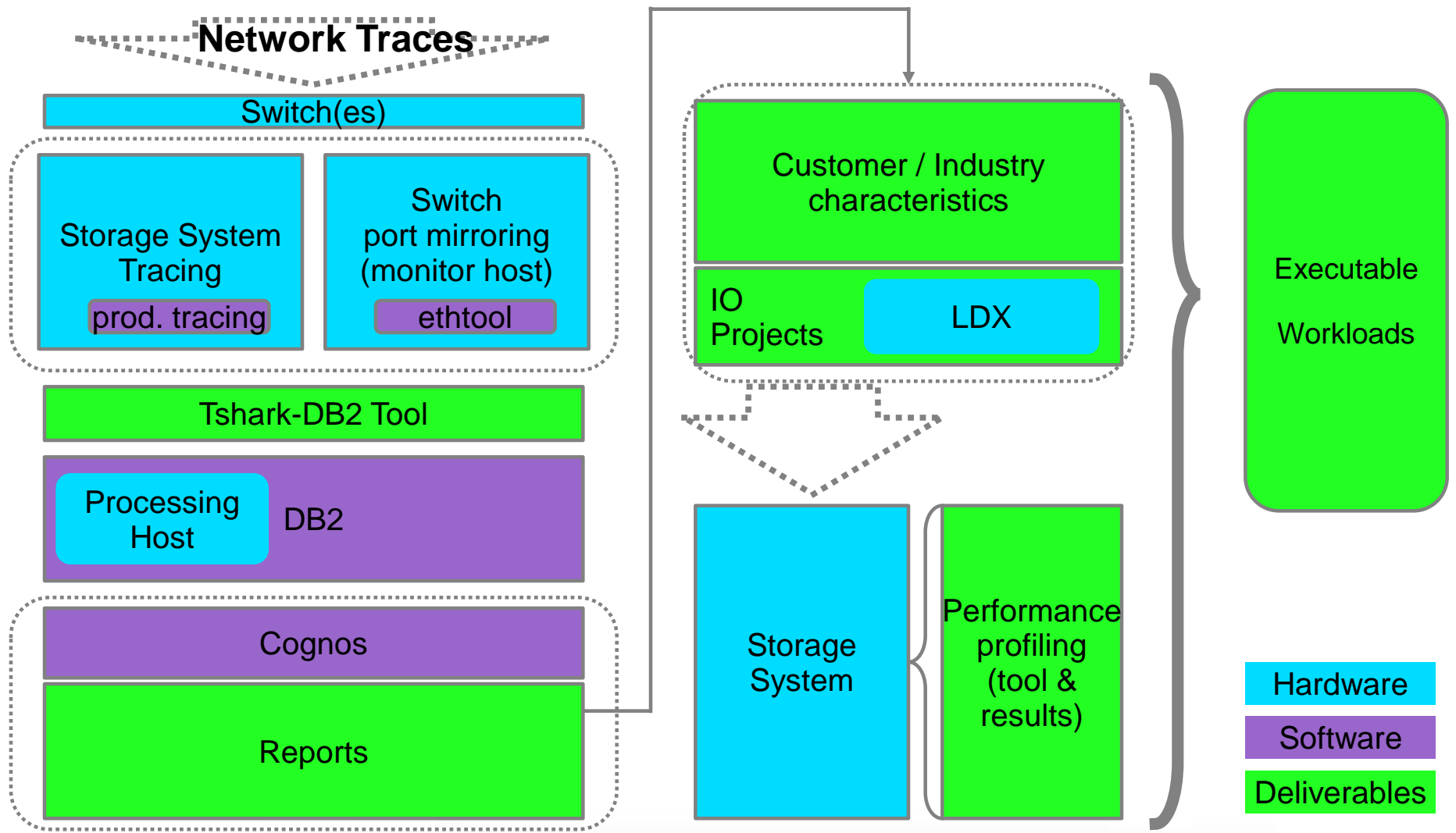
### ❑ Solution:

1. Gather traces & perform trace parsing
2. Load parsed data into analytics tools to create views
3. Develop workloads

### ❑ Challenges:

- ❑ Define indicators
- ❑ Need for analytics tools
- ❑ Needs pre-parsing of unstructured data
- ❑ Trace gathering may impact performance

# Workload Analysis Flow



# Questions?



# Acryonyms

- ❑ PIT (Point-In-Time copy)
- ❑ EDA (Electronic Design Automation)