# Introduction to LeoFS

*A Market Proven Parallel File System for Data Intensive Storage*

# A Supreme Parallel File System
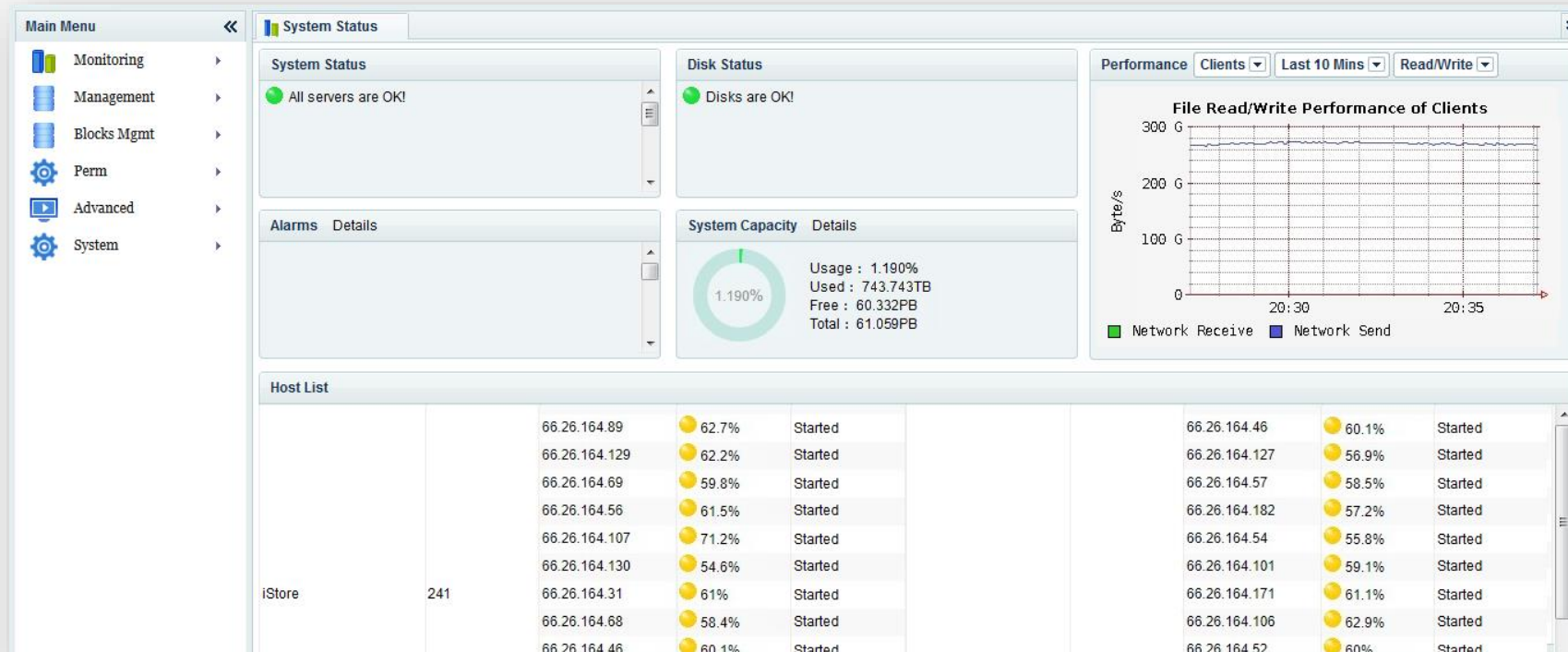
- One cluster for file, block and object-based storage

- Market proven - more than 1 EB capacity deployment

- Sustainable high performance - always saturate hardware throughput

- Reliable N+M erasure coding - up to 90% capacity utilization

- Easy scalability with no downtime or reboot

- Worry-free 24/7 customer support and management

- Great cost savings vs. traditional and open source file systems

*Current largest single cluster installed has 333 storage nodes, 95PB and over 200GB/s I/O throughput.*

*Customer Satisfaction Guaranteed*

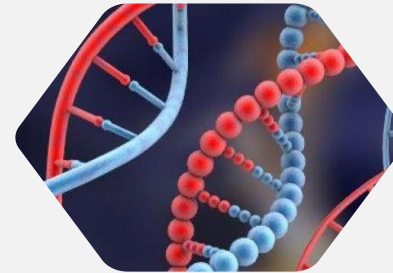*Successful deployment in Asia since 2009*

# Computational Storage

- Single server and cluster solutions, from tens of TB to hundreds of PB

- Both metadata and computing can be embedded in storage nodes

- Customer on-site 60PB cluster
  - 241 commodity servers: 4U 36-drive, dual Intel E5-2620v4, 64 GB RAM
  - 8 metadata + storage nodes: 2*480 GB SSDs, 34*8 TB SATA HDDs
  - 233 compute + storage nodes: 36*8TB SATA HDDs

# Hundreds of Customers

- Industry success
  - Oil and Gas
  - Scientific Computing
    - ✓ Genomics
    - ✓ Cryo-electron Microscopy
    - ✓ Satellite Imaginary/Observatory
    - ✓ Geographical Data and Mapping
    - ✓ Meteorology/Climate
  - Higher Education
  - Media and Entertainment
  - Telecom and Internet
  - AI and Big Data
  - Video Surveillance

# Selected References

- We take pride that most of our customers in PB usage started with only a few hundreds of TB
  - First Oil & Gas customer in 2009, I/O throughput 2x greater than StorNEXT
  - Higher Education: University of Florida, Georgia Southern University
  - Scientific Computing: Direct Electron (San Diego) on electron detection for biological molecules
  - Video Surveillance: Dante Security (New York)

# Solution Architecture

- Fully-POSIX compliant

- Native clients, all kernel modules that do not require any patches

- Best choice for combined values
  - High performance
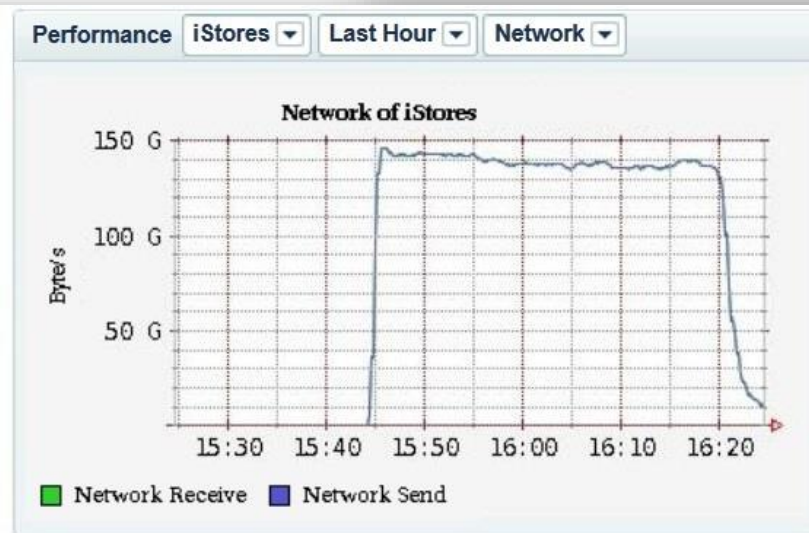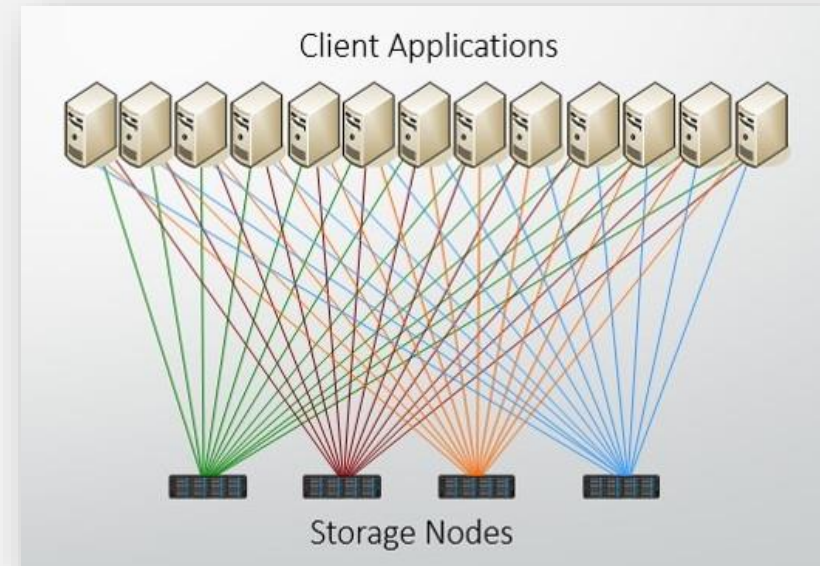  - Large and scalable capacity
  - Reliable data protection
  - Professional support
  - Affordable cost

# No Bottleneck

- Without controllers, gateways, nor distributors
  - Data files are transparently distributed over multiple nodes
  - All client applications communicate with all storage nodes
- Customer on-site 7.5PB
  - Dual 10GbE network
  - 102 nodes of 4U 24-bay storage servers
  - Total of 1,880 drives with 4TB SATA HDDs
  - Aggregated I/O close to 150GB/s
  - Average single drive write 75MB/s

# Unlimited Capacity

- Software as a service, option to choose cluster or single server products
  - Cluster starts with 2 nodes, and up to thousands
  - Single server supports 12-bay, 16-bay, 24-bay, or 36-bay drives
- System capacity and throughput always increase with additional server or drive
- Solution threshold

|  | Theoretical | Actual Deployment |
|---|---|---|
| Storage nodes | 4,096 | 333 |
| Metadata servers | 256 | 32 |
| System capacity | EB | 95PB |
| Number of files | Unlimited | 50 Billion |

# File-level N+M Erasure Coding

- Data content is distributed on a file-level across different storage nodes
  - When N+2 is applied, cluster can sustain operation up to two simultaneous failures

- Optimum data protection plans for different files
  - Better capacity utilization, up to 90% with 16+1

- With failed hardware, LeoFS rebuilds only the files that are affected, and it uses the entire cluster to rebuild
  - No downtime nor reboot
  - One TB data usually takes less than 20 minutes

- Highest level of data availability
  - System capable of self-monitoring and self-healing

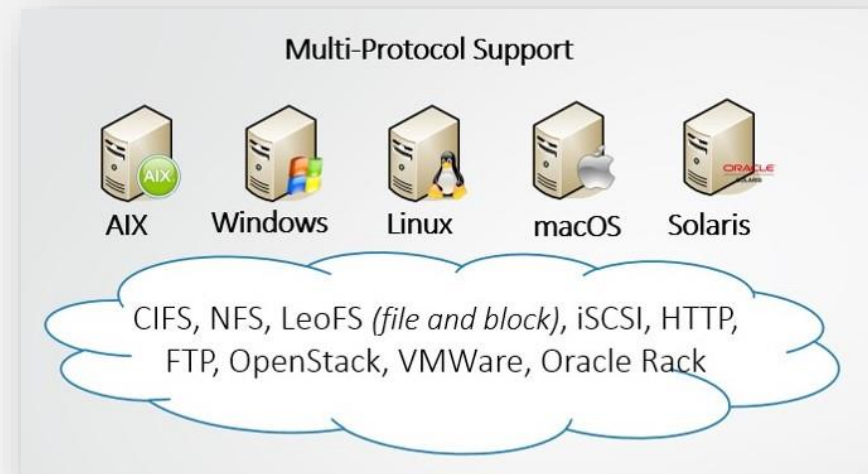# Peace of Mind

- It's CARE, MANAGEMENT and SUPPORT
  - Cluster monitoring
  - Software maintenance and update
  - High quality hardware
- With direct access to file system developers
  - Designated consultants are available 24/7
  - Most work can be done remotely
- Options to choose
  - Next Business Day Service Level Agreement
  - Re-mote or On-Site Support Warranty
  - Advanced Hardware Replacement

# Enterprise Features

- Load balance switch, hardware evenly share system workload

- Runs on platforms such as x86, OpenPOWER, ARM, and Xeon Phi

- Re-export through Samba, NFS, FTP, HTTP, LeoSAN or iSCSI

- Support for group/user ACLs and quota

- Fully active network with automatic failure detection

- Supports Infiniband, GigE, multiple subnet and bonding

- Cold data sanity check, automatic repair, no downtime

- WORM directory, avoid modification of saved data



Multi-Protocol Support

AIX   Windows   Linux   macOS   Solaris

CIFS, NFS, LeoFS *(file and block)*, iSCSI, HTTP, FTP, OpenStack, VMWare, Oracle Rack



Comp Chart

|  | Isilon Nitro | IBM Spectrum Scale | Lustre | LeoFS |
|---|---|---|---|---|
| Snapshots | Yes | Yes-Complex | No | Yes |
| Independent capacity/performance scaling | No | No | No | Yes |
| Scale to thousands of nodes | No | Yes | Yes | Yes |
| QoS | Yes | No | No | Yes |
| N+M Data Protection | No | No | No | Yes |
| Encryption | Yes | Yes | No | Coming |
| S/W only, H/W independent | No | Yes | Yes | Yes |
| IB & GbE Support | No | Yes | Yes | Yes |

# Solution Sample – University of California

- UCSF Wynton HPC center:
  - ✓ 1.2PB storage, cost about $192K USD
  - ✓ Hardware: 4 nodes of 60-bay servers, 2 nodes of metadata servers
  - ✓ Software: ZFS and BeeGFS
  - ✓ https://wynton.ucsf.edu/hpc/about/pricing-storage.html Genomics

- Competitive LeoFS cluster
  - ✓ 1.4PB storage, only $100K USD
  - ✓ Throughput: read 6.8GB/s, write 10GB/s (asynchronous)
  - ✓ Hardware: 4 nodes of 36-bay servers

    - **2 Metadata + Storage nodes**
      - CPU:Intel Xeon E5-2630 V4 x 2
      - Motherboard: Supermicro X10DRL-i
      - HBA: LSI SAS 9300-8I
      - System Disk Drive: 480GB SSD x 2 + 240GB SSD x 2
      - Storage Disk Drive: 10TB HDD x 34
      - RAM: 128GB
      - Network Port: 4 x 10 GbE

    - **2 Storage-only node**
      - CPU:Intel Xeon E5-2630 V4 x 2
      - Motherboard: Supermicro X10DRL-i
      - HBA: LSI SAS 9300-8I
      - System Disk Drive: 240GB SSD x 2
      - Storage Disk Drive: 10TB HDD x 36
      - RAM: 64GB
      - Network Port: 4 x 10 GbE

# Solution Sample – Geneva Observatory

- BeeGFS case study of Geneva Observatory
  - 4 storage nodes and 2 metadata servers, Infiniband
  - Effective 800TB, 144 drives, I/O 5-8GB/s
  - https://www.hpc-ch.org/hpc-ch-forum-yves-revaz-epfl-beegfs-the-hpc-storage-solution-adopted-at-the-geneva-observatory/
- Competitive LeoFS cluster
  - 4 nodes of 4U 36-bay storage servers, 2 metadata servers, dual 10GbE
  - Usable 1PB, 144 drives, I/O 8-11GB/s
  - No buddy mirroring, 80% capacity utilization
  - No single point of failure from drives, nodes or network
  - File-level RAID, faster data recovery
  - Better ROI



**BeeGFS**

The HPC Storage Solution Adopted at the Geneva Observatory

Yves Revaz

EPFL
ECOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

OBSERVATOIRE DE GENÈVE
FONDÉE EN 1772

The Geneva Observatory in a nutshell

Main fields of research
- Stellar Physics
- Galaxy and Cosmology
- High energy astrophysics
- Extra-solar planets

# Vs. EMC Isilon S Series

- From EMC 3D animation customer case
  - Ten nodes of Isilon S200 series
  - Raw capacity of 600 TB
  - I/O throughput 8 GB/s
- Comparable LeoFS solution
  - Commodity hardware, dual 10GbE network
  - Ten nodes of 4U 24-bay storage servers
  - Raw capacity of 960 TB
  - I/O throughput over 12 GB/s

  *60% more capacity and 50% higher throughput*

# Vs. DDN ES7K

- World's fastest Lustre appliance
  - SSD cache, SSD, performance SAS, and capacity SAS
  - 8U 144 drives, I/O up to 12 GB/s

- Comparable LeoFS solution
  - Six nodes of 4U 24-bay storage servers, 144 drives
  - Using only SATA HDDs for storage, 4 SSD drives for metadata
  - Sustainable I/O without deterioration over time

  *A better balance of Price, Performance and Capacity.*

# Vs. Lustre

- Easy installation and management
  - No problems having a large number of files in a single directory
  - No problems accessing huge amount of small files
  - No worries using ls –l

- Metadata cluster up to 256 servers, scalable performance
  - Each pair 20,000 file creates per second
  - Standalone or within storage nodes

Single Log-on GUI

# Throughput vs. Lustre and GPFS

- Lustre/GPFS*: no data on X clients and X streams
  - Six LUNs, each chassis with 30*4TB disks, total of 180 HDDs
  - Because of controller, limitation seen on 8 clients, 1 stream
    - Write: 4-5 GB/s, Read: 5-7 GB/s

- LeoFS without controller limitation
  - Eight nodes of 4U 24-bay storage servers, total of 192 HDDs
  - Dual 10 GbE network, 8 clients, 20 streams, 4TB testing data
    - Each client 500GB (greater than 64GB RAM)
    - File size 1MB, redundancy 8+3
    - Write 11GB/s, Read 7GB/s

*CERN's presentation on High Performance Storage in Science, SDC 2017*

# Vs. Ceph

- Block device SSD IOPS
- Testing hardware
  - 3 nodes of 12-bay storage servers, each with 12*240GB SSD
  - Dual 10GbE, 3 clients, block size 4KB
  - Replications: 2 and 3



WRITE IOPS(BS=4K)

READ IOPS(BS=4K)

# Products – Cluster

- **Starting minimum 2 nodes**
  - New system or add to existing LeoFS installation
  - Option to build separate instances of LeoFS
  - Hyper-convergence option
  - Integrated virtualization technology
  - Automated self-monitoring and self-healing
  - Potentially geographically redundant storage (LeoSync)
  - Supports CIFS, NFS, FTP, HTTP, ISCSI, OpenStack and Hadoop
  - LeoSAN and LeoFS private block and file interfaces
  - Fine-grained access rights
- 2U 12 bay to 4U 36 bay chassis
  - CPU: Intel Xeon E5-2630 V4 * 2
  - Mother Board: Supermicro X10DRL-i
  - HBA: LSI SAS 9300-8I
  - System Disk Drive: 480GB SSD * 2 + 240GB SSD * 2
  - Storage Disk Drive: 2TB to 10TB HDD, SATA or SAS
  - RAM: 128GB
  - Network Port: 2 * 10 GbE or 40 GbE or Infiniband

# Products – Single Server

- Benefits over typical NAS or SAN products
  - Plug and use, fully-POSIX complaint, supports NFS or ISCSI
  - Inexpensive commodity hardware, HDDs, 1 GbE or 10 GbE
  - Single directory up to 10 billion files
  - High performance, throughput > 2.5 GB/s, IOPS > 30,000
  - File-level RAID, better fault tolerance than RAID 6
  - Dynamic scalability with load balance switch
  - Faster online disk rebuild or replacement, no operation interruption
- 2U 12 bay to 4U 36 bay chassis
  - CPU: Intel Xeon E5-2620 V4 * 2
  - Mother Board: Supermicro X10DRL-i
  - HBA: LSI SAS 9300-8I
  - System Disk Drive: 480GB SSD * 2 + 240GB SSD * 2
  - Storage Disk Drive: 2TB to 10TB HDD, SATA or SAS
  - RAM: 64GB
  - Network Port: 2 * 1GbE or 10 GbE

# Thank you, let's get in touch

- No licensing fees
  - Free 1PB case study usage if replacing other parallel file system
  - Monthly service subscription with optional 24/7 support
  - Unbeatable pricing, contact us for your free solution estimate

*www.leofs.info*