



# Intel® Enterprise Edition for Lustre\*

Zhang HongChao ([hongchao.zhang@intel.com](mailto:hongchao.zhang@intel.com))

High Performance Data Division, Intel® Corporation

**Breakthrough Storage Performance**  
**LUG 2014**

Oct 14 2014  
Beijing, China

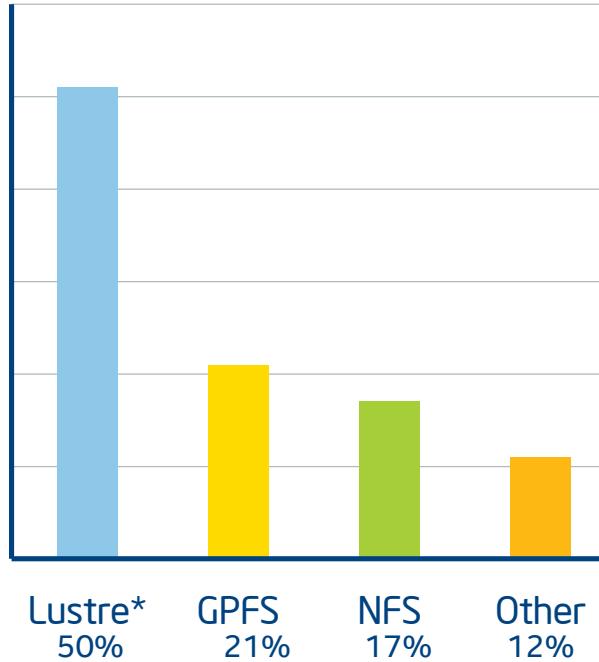


# 日程

- 英特尔Lustre\*企业版
- 英特尔Lustre集成化管理环境(IML)
- 英特尔Lustre企业版推荐实施配置

# HPC领域使用的主要文件系统

- Over 50% rely on Lustre\* today....
- Over double the share of GPFS
- NFS performance and scaling limits appeal



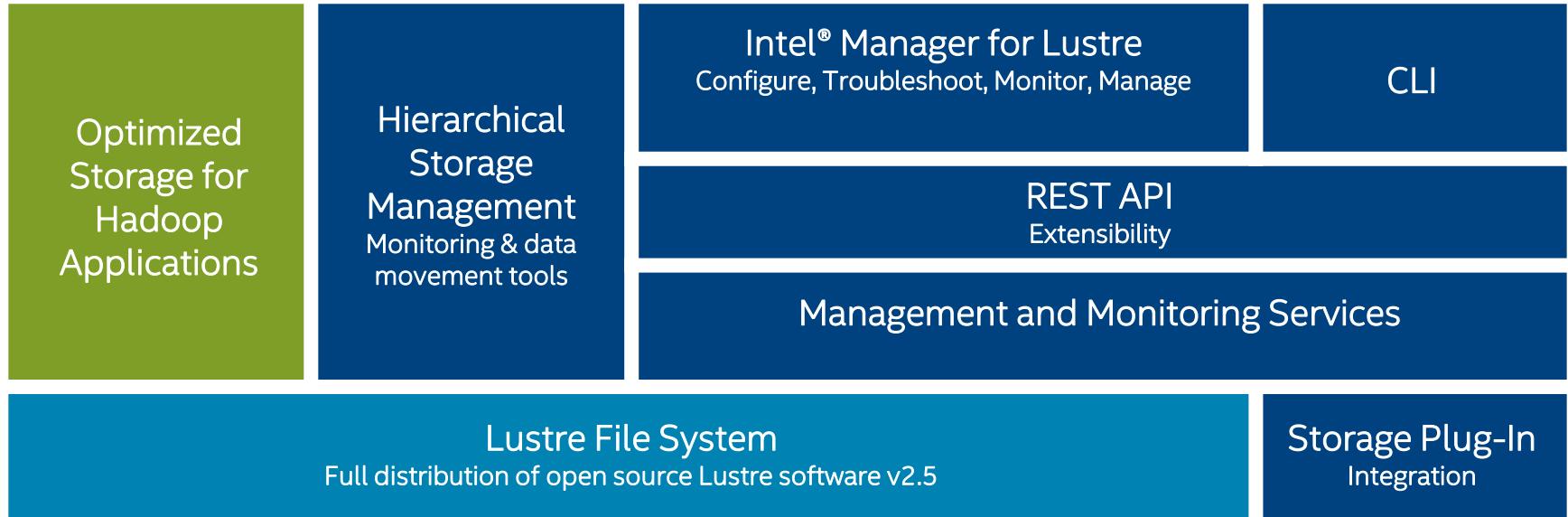
*IDC survey research, May 2014*

\* Other names and brands may be claimed as the property of others

# 英特尔Lustre\*企业版用户

- 能源, 石油和天然气勘探
  - Schlumberger (北京), Total/Fina, Shell, BP, Chevron
- 政府部门和教育科研机构
  - 上海天文台, 北京防灾大学, 北方交通大学, CEA, IU, ORNL, LLNL, ANU, SDSC, Purdue University, Harvard, Stanford
- 金融服务
  - 邮储银行, Morgan Stanley, Central Bank of Italy, National Bank of Poland
- 生命科学, 制药, 基因科学
  - Sanger Institute, Lund University, Cleveland Clinic
- 天气预报和气候建模
  - Australia Weather Bureau

# 英特尔Lustre\*企业版 V2

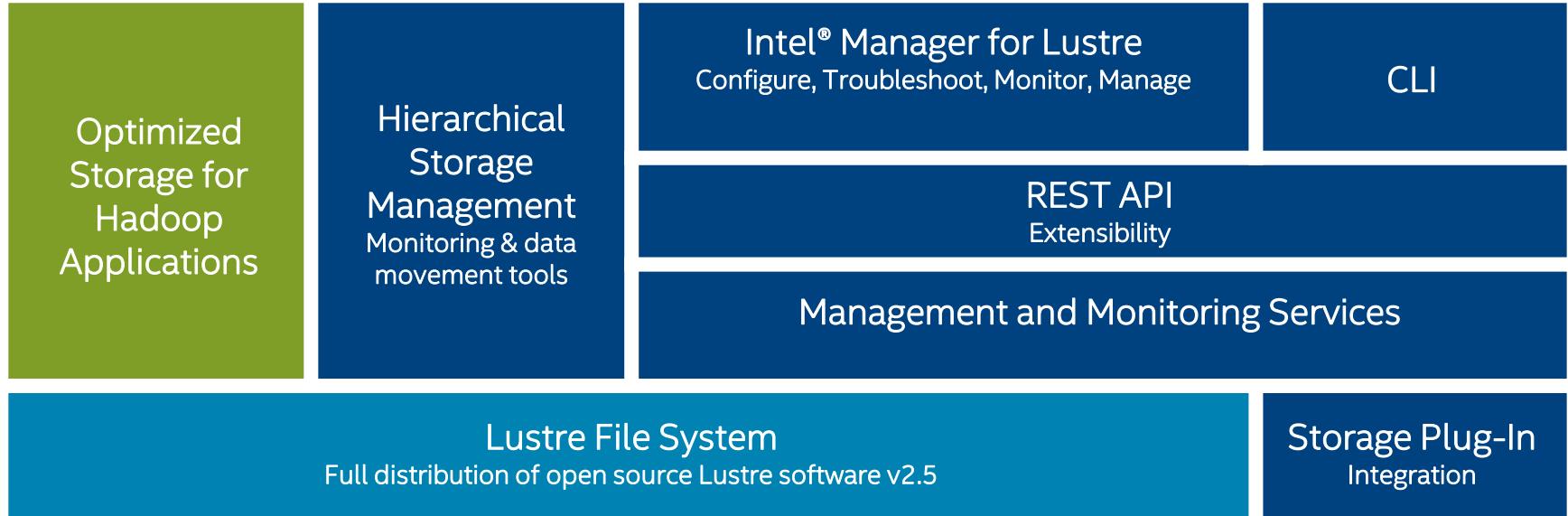


Open source base

Intel value-add for Lustre

Interoperability with Hadoop distributions for fast, shared, simple to manage storage for MapReduce applications

# 英特尔Lustre\*企业版 V2



Open source base

Intel value-add for Lustre

Interoperability with Hadoop distributions for fast, shared, simple to manage storage for MapReduce applications

# 英特尔Lustre\*企业版 V2

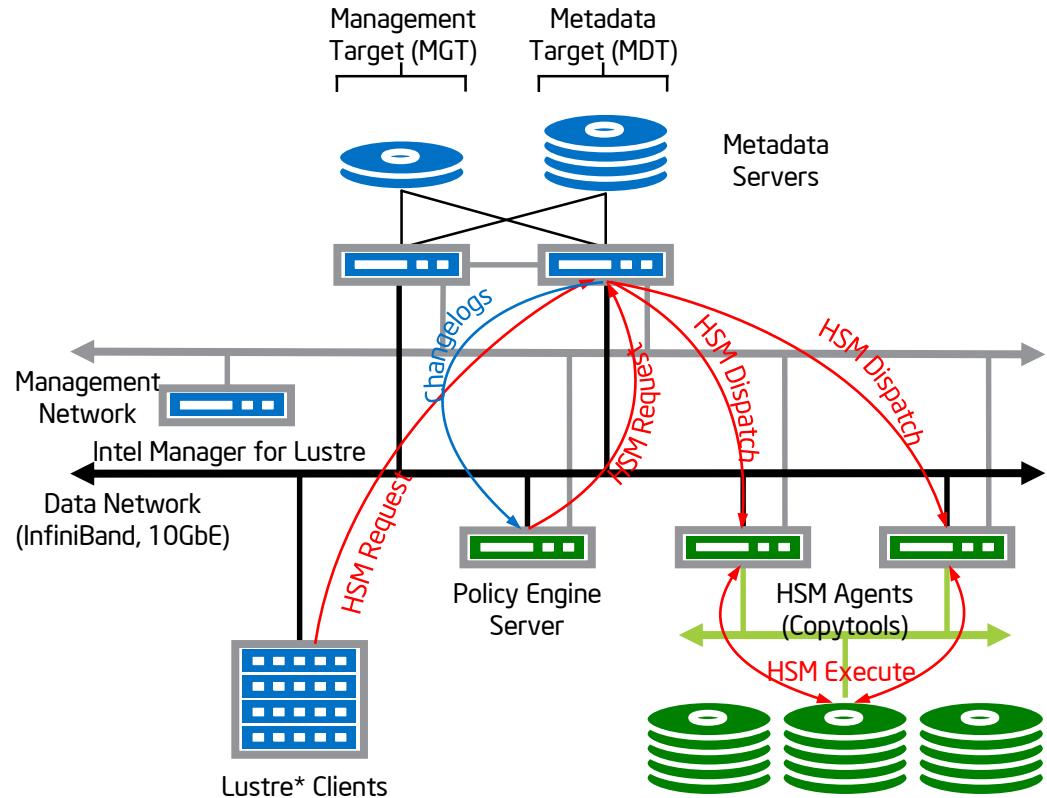
Feature or major enhancement	Benefit	Lustre 2.5	Intel® Enterprise Edition for Lustre* v2.0
Supported by Intel		✓	✓
CLI based management		✓	✓
Multi-vendor storage		✓	✓
Intelligent installation	Easily install software components from central server		✓
Intuitive configuration	Easily deploy proven, enterprise-grade storage		✓
In depth UI-based monitoring	Actively monitor previously installed Lustre solution		✓
Comprehensive UI-based management	Lowers complexity, raised productivity and overall ROI		✓
UI based configuration of hierarchical storage mgmt.	Enable HSM coverage for optimizing storage utilization		✓
Advanced and extensible charting	Detailed insights into storage in near-real time		✓
Automated configuration of high availability services	Easily configure server pairs for optimum availability		✓
Automated configuration of Lustre Networking	Ensure LNET services are configured correctly		✓
Intelligent log files	Quickly and easily understand important messages		✓
Storage plug-in infrastructure	Extends IML charting to include underlying hardware		✓
REST API	Integration with management tools		✓
Software integration layer (adapter) for Hadoop	Fast, shared, easily managed storage		✓

\* Other names and brands may be claimed as the property of others



# 完整的分层存储体系结构(HSM)

- HSM提供高效动态数据迁移机制。在同一套文件系统中整合多种级别的存储设备，通过各种策略调整热数据和冷数据在存储系统中的配置，以达到最高的性价比。
- 图形化地进行配置和管理
- 提供完备的接口和ISV应用集成
- 好处：分层存储降低总体拥有成本和简化资源共享，优化整体存储使用效率。



# Hierarchical Storage Management (HSM)

- HSM Coordinator (at MDT)

```
lctl set_param [-P] mdt.<fsname>-<MDT index>.hsm_control=enabled
```

```
lctl get_param mdt.<fsname>-<MDT index>.hsm_control
```

- HSM Agent (Copytools)

```
mount /dev/sdb /archive/demo
```

```
lhsmtool_posix --daemon --hsm_root /archive/demo --archive 1 /lustre/demo
```

- Policy Engine

```
lctl set_param mdd.<fsname>-MDT*.changelog_mask="all-XATTR-MARK-ATIME"
```

```
lctl --device <fsname>-<MDT index> changelog_register
```

```
sudo rbh-lhsm --scan -once -f /etc/robinhood.d/lhsm/demo-lustre-hsm.conf
```

```
sudo service robinhood-lhsm start
```

# HSM Agent

- Using “Archive Identifier” to distinguish different high capacity storages
- Each agent binds with one or more “Archive Identifier”
- “0” means all archives (not recommended to use)
- HSM Coordinator at MDT will choose the best agents according to the work load of each agent (by the number of pending requests)

# HSM Policy Engine Configuration

General –the mount point of Lustre\*

Log – the information about various logs

ListManager – the configuration of MySql

ChangeLog – the MDT device name and the Changelog registration ID

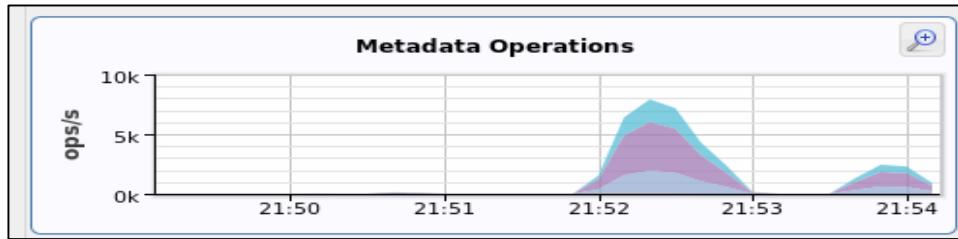
Migrate – copying files to the high capacity archive storage.

Purge – releasing files that have been archived in order to free capacity in Lustre\*.

Remove – cleaning up archive copies of deleted files.

# 针对Hadoop应用优化的高性能并行存储

- Hadoop应用的读写模式与传统的文件系统操作完全不同，特点是同时打开/关闭海量文件
- 当运行大型Hadoop作业时元数据服务器会承受很大的读写压力，经常导致过载现象发生。
- 英特尔Lustre\*企业版中提供新的扩展属性缓存机制--提高元数据服务器的性能和抗压能力。
- **好处：**为Hadoop应用提供高效，高性能，高扩展性，易管理的存储解决方案

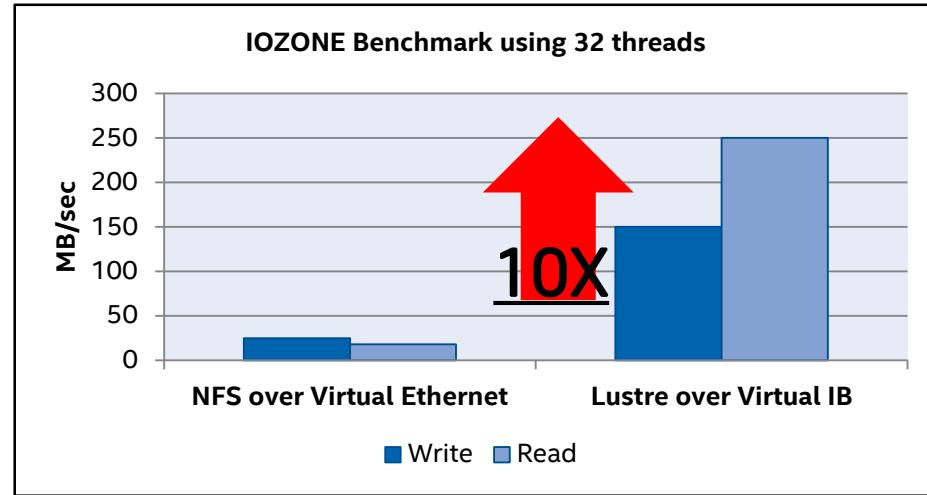


Metadata operations during a Terasort experiment<sup>1</sup> using a 3TB dataset and 240 Map-Reduce tasks

<sup>1</sup> Results have been estimated based on internal Intel analysis and are provided for informational purposes only. Any difference in system hardware or software design or configuration may affect actual performance.

# 原生Intel® Xeon Phi™ Lustre \* 客户端

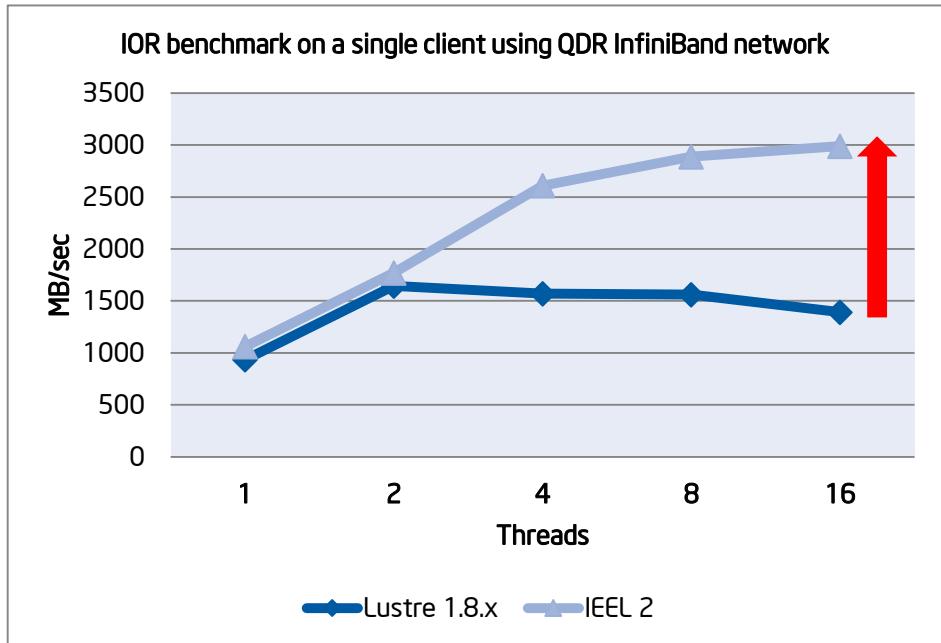
- 针对于多节点大型平行作业，例如每个节点均需要100+核
- 使在Xeon Phi™中运行的应用直接访问Lustre存储
- **好处：**为Xeon Phi™应用提供更高性能的IO访问能力



<sup>1</sup> Results have been estimated based on internal Intel analysis and are provided for informational purposes only. Any difference in system hardware or software design or configuration may affect actual performance.

# 优化单客户端读写性能<sup>1</sup>

- 在企业市场中经常会有需求为单线程单客户端的应用提供高IO服务
- 英特尔Lustre\*企业版提供针对于单客户端数据流的优化。
- **好处:** 提高商业应用软件的读写性能, 例如 : MapReduce



<sup>1</sup> Results have been estimated based on internal Intel analysis and are provided for informational purposes only. Any difference in system hardware or software design or configuration may affect actual performance.



# Intel® 集成化Lustre\*管理环境

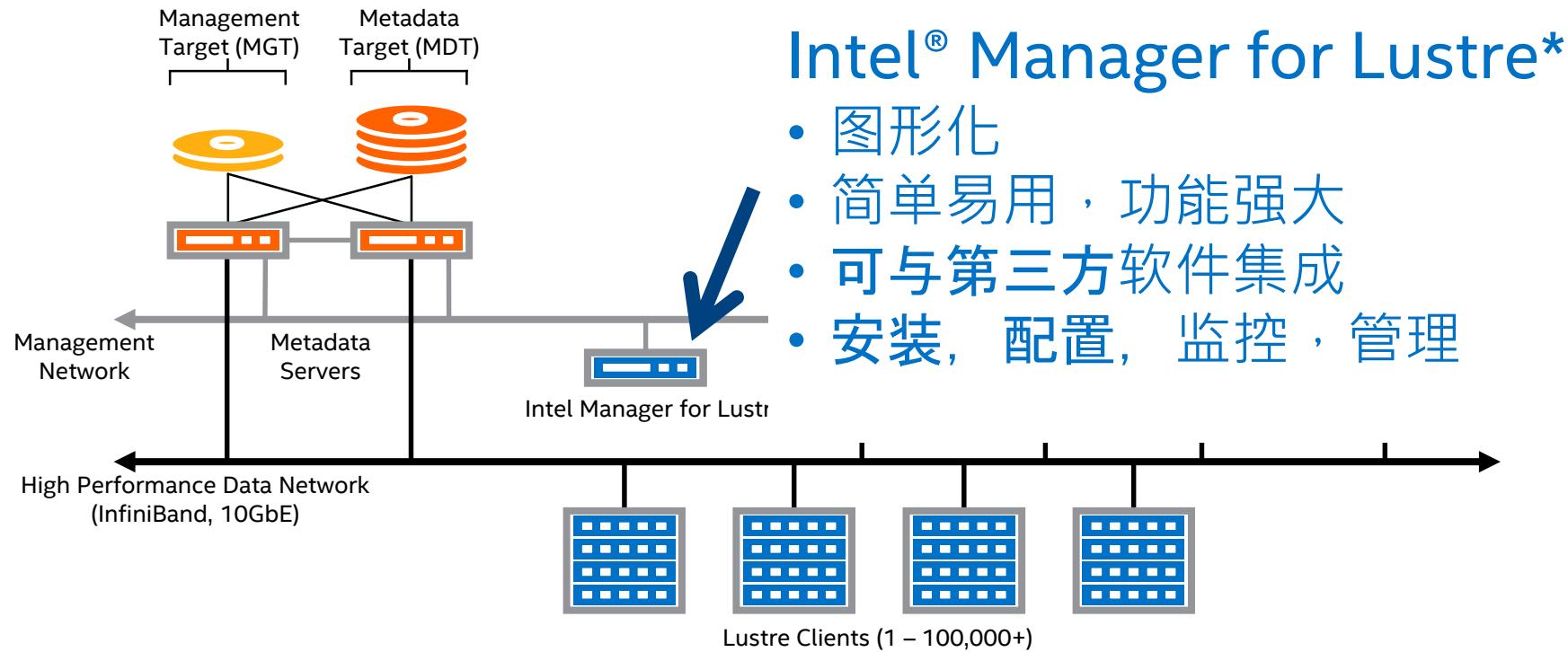
# 配置，优化，管理Lustre\*文件系统

易用、高效、功能强大的图形管理平台，同时汇集了Intel开发支持团队的十多年最佳实践精华于其中。

- 利用图形化界面指导系统管理员，高效地配置、监控及管理整个Lustre存储平台以降低管理复杂度和管理成本
- 通过各种图表揭示存储系统内部的实时性能
- 自动配置存储服务器之间的互备机制，以实现文件系统高可靠性。
- 智能化的预警设计及日志分析，使得系统管理员及时掌控文件系统的状态。

**Intel® Manager for Lustre Makes Lustre Smarter, Simpler and More Productive**

# 降低管理复杂度



# Intel® Manager for Lustre® 图形化管理界面



集成显示面板整合各种图表显示性能和资源利用率。系统管理员可以很容易地浏览文件系统，监测作业的资源利用状态以及性能。

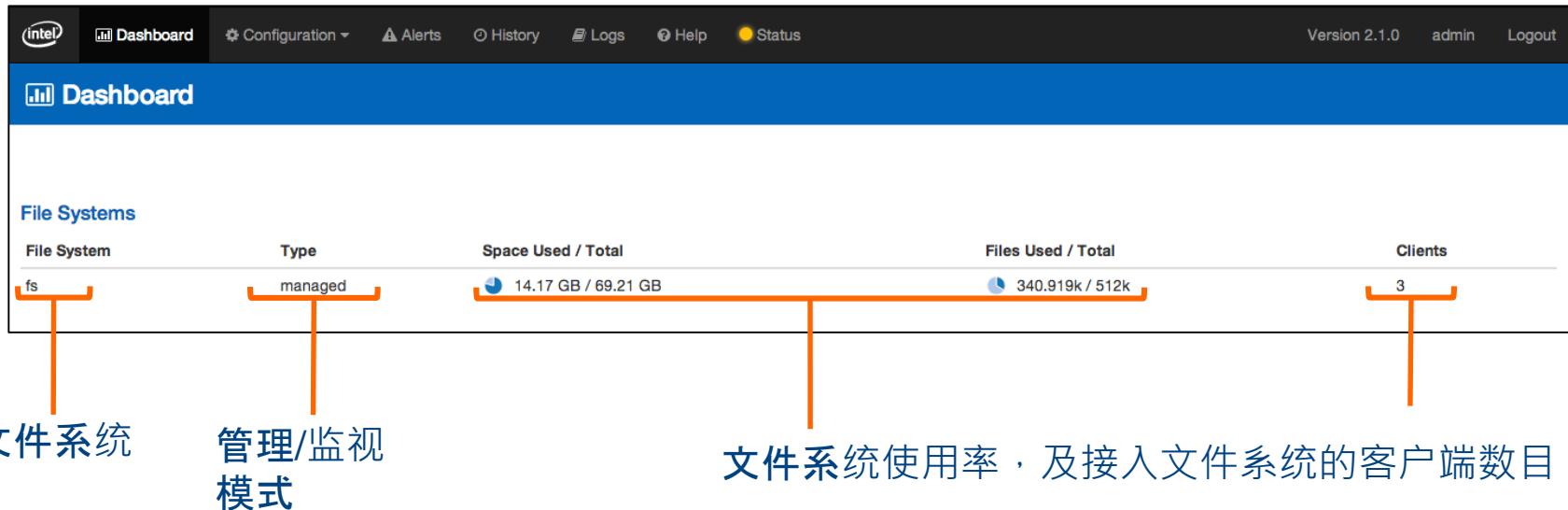


配置服务器，存储卷，  
电源管理，进一步激活  
并配置HSM服务

智能化的报警和  
日志监控机制

文件系统全局  
状态指示灯

# Intel® Manager for Lustre\* 显示面板



# 简化配置流程



**Servers:** 添加, 启动, 管理Lustre文件系统服务器

**Power Control:** 配置使用PDU或者BMC的电源管理

**File Systems:** 创建、删除、启动、停止文件系统

**HSM:** 激活分层存储服务

**Storage:** 集成存储硬件插件

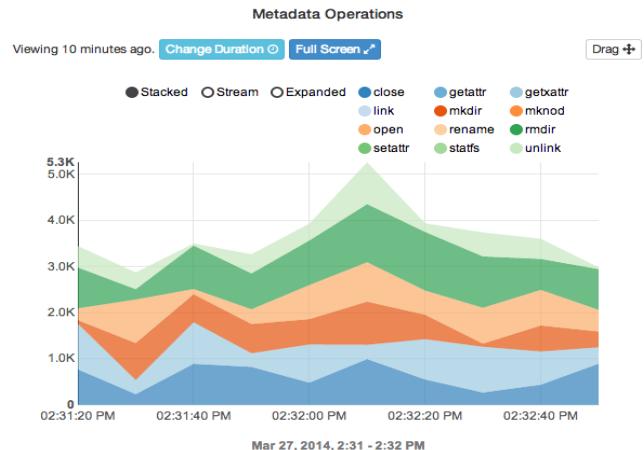
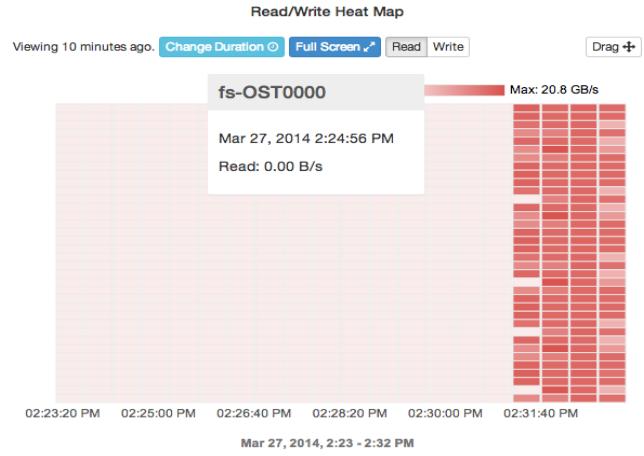
**Users:** 管理账号

**Volumes:** 配置存储卷及高可靠性互备机制

**MGTs:** 设置管理存储设备的参数

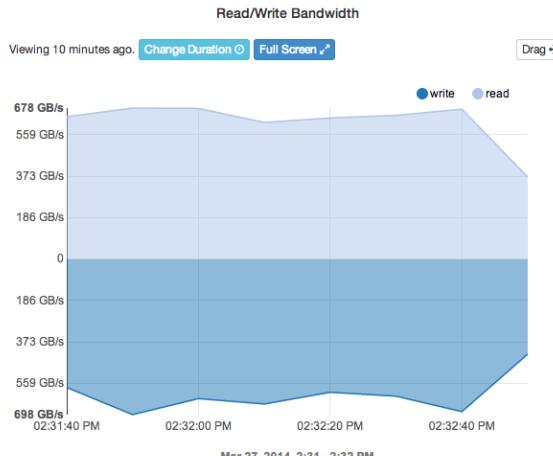
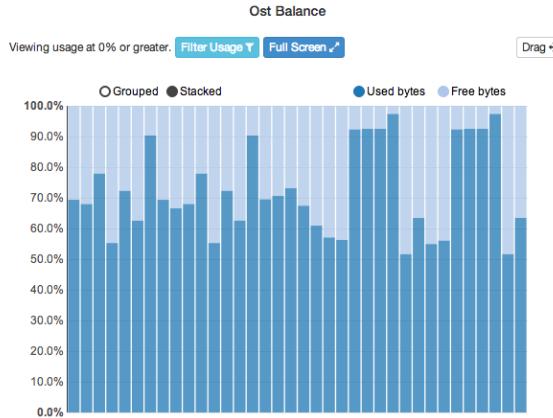
# 热点和元数据操作

- 动态显示存储资源的使用分布效果，帮助系统管理员迅速确定热点区域并进一步分析到作业以及底层存储设备读/写操作。
- 统计整个文件系统的全部元数据 IO 操作，通过各种图表以展示从全局一直细化到单类型操作的状态



# OST分布和IO统计

- 揭示底层对象存储设备使用率分布状态，帮助系统管理员直观地确认出是否需要对文件系统中的存储空间做再均衡操作。同时提供每一个OST的详细信息。
- 统计显示读写带宽，展示从整个文件系统和细化到单个OST的读写吞吐量。

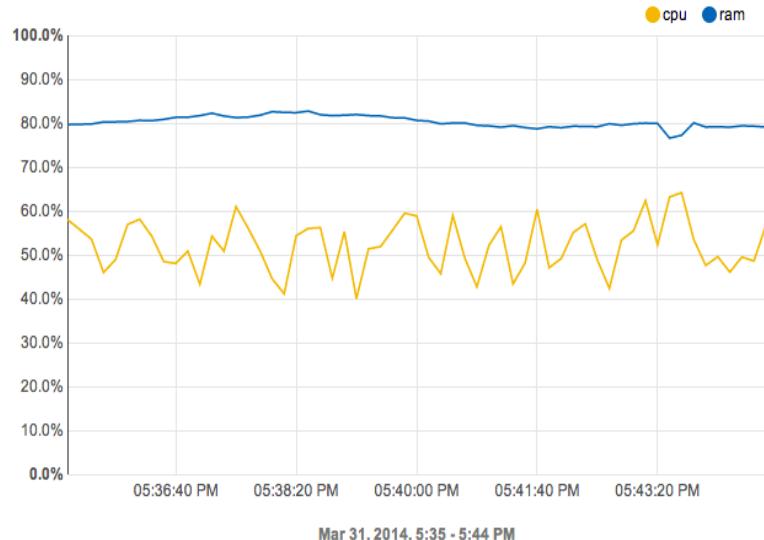


# CPU和内存利用率

Metadata Servers

Viewing 10 minutes ago. [Change Duration](#) [Full Screen](#)

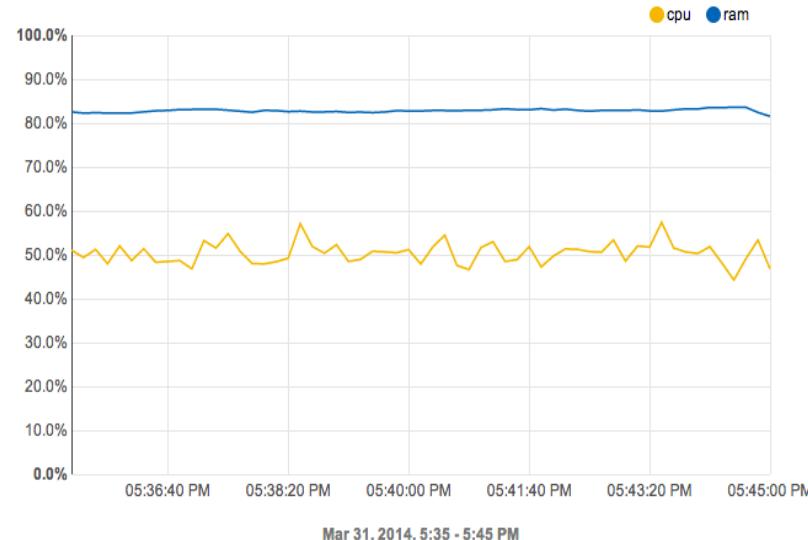
Drag +



Object Storage Servers

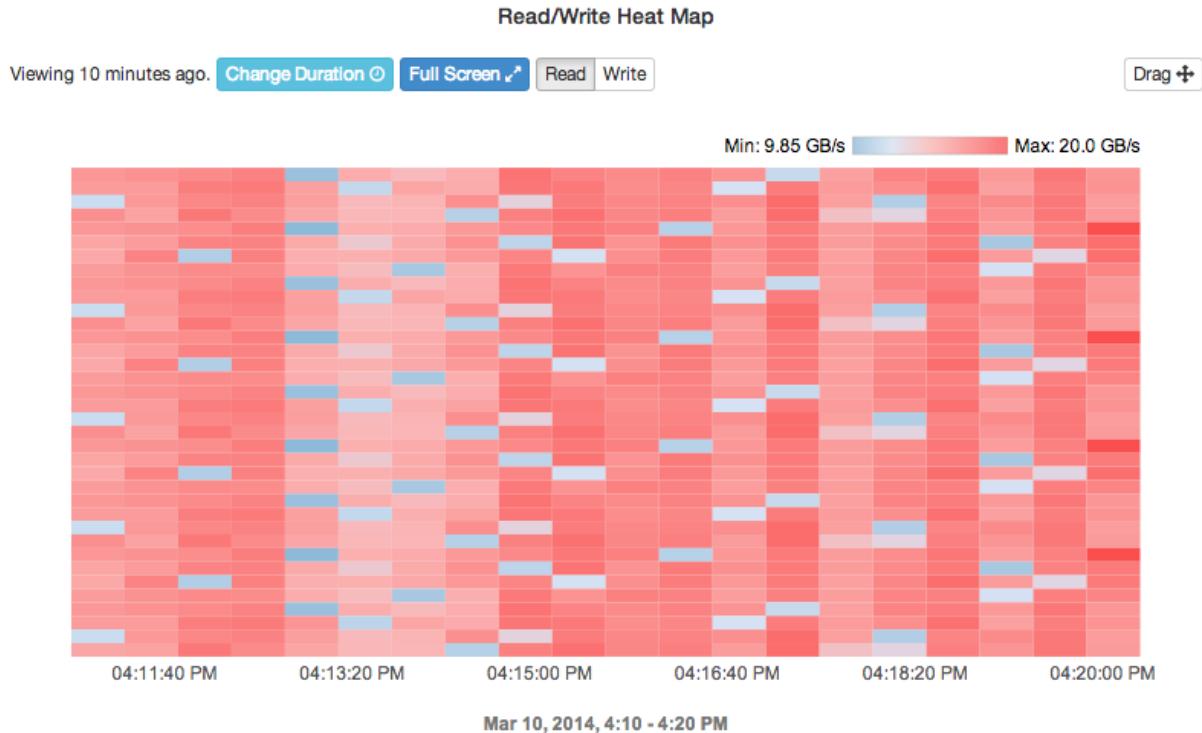
Viewing 10 minutes ago. [Change Duration](#) [Full Screen](#)

Drag +

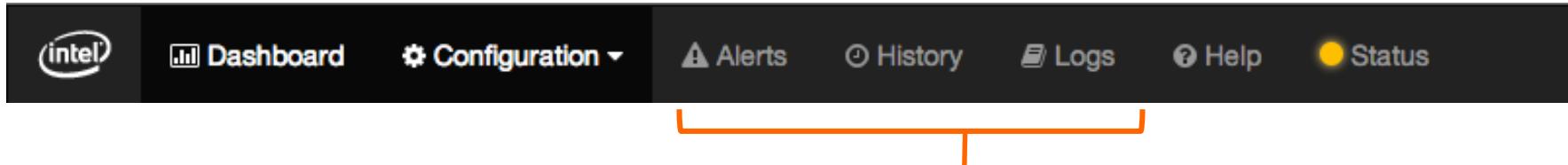


# 作业状态统计

- 针对指定作业做深度分析
- 针对指定OST的作业统计



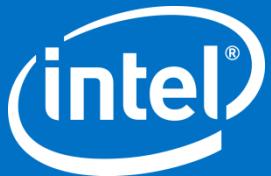
# 事件、报警及日志



简化地集中  
管理界面

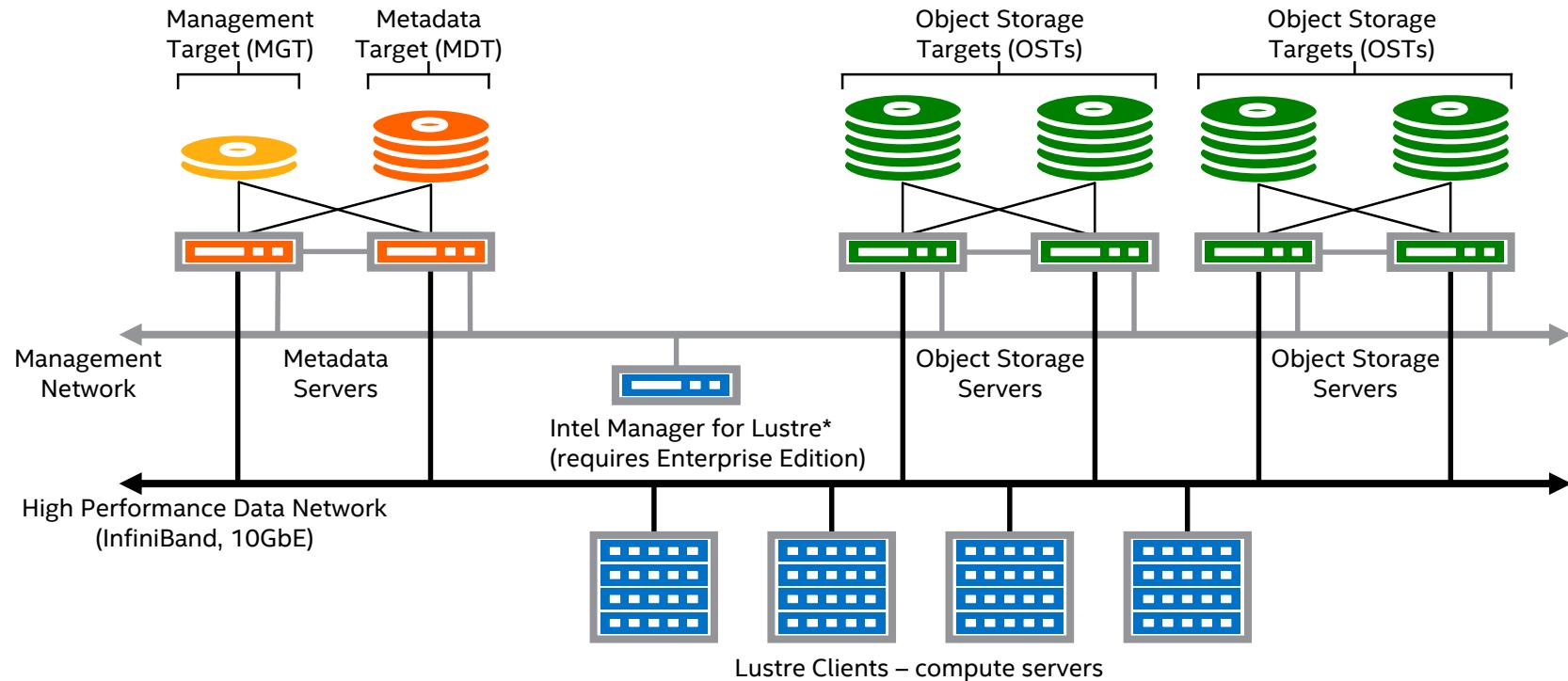
整合整个Lustre\*存储集群

- 预警
- 事件及操作历史（可设置查询时间段）
- 日志（解析过的）



# 英特尔Lustre\*企业版推荐实施配置

# IEEL推荐架构图



\* Other names and brands may be claimed as the property of others

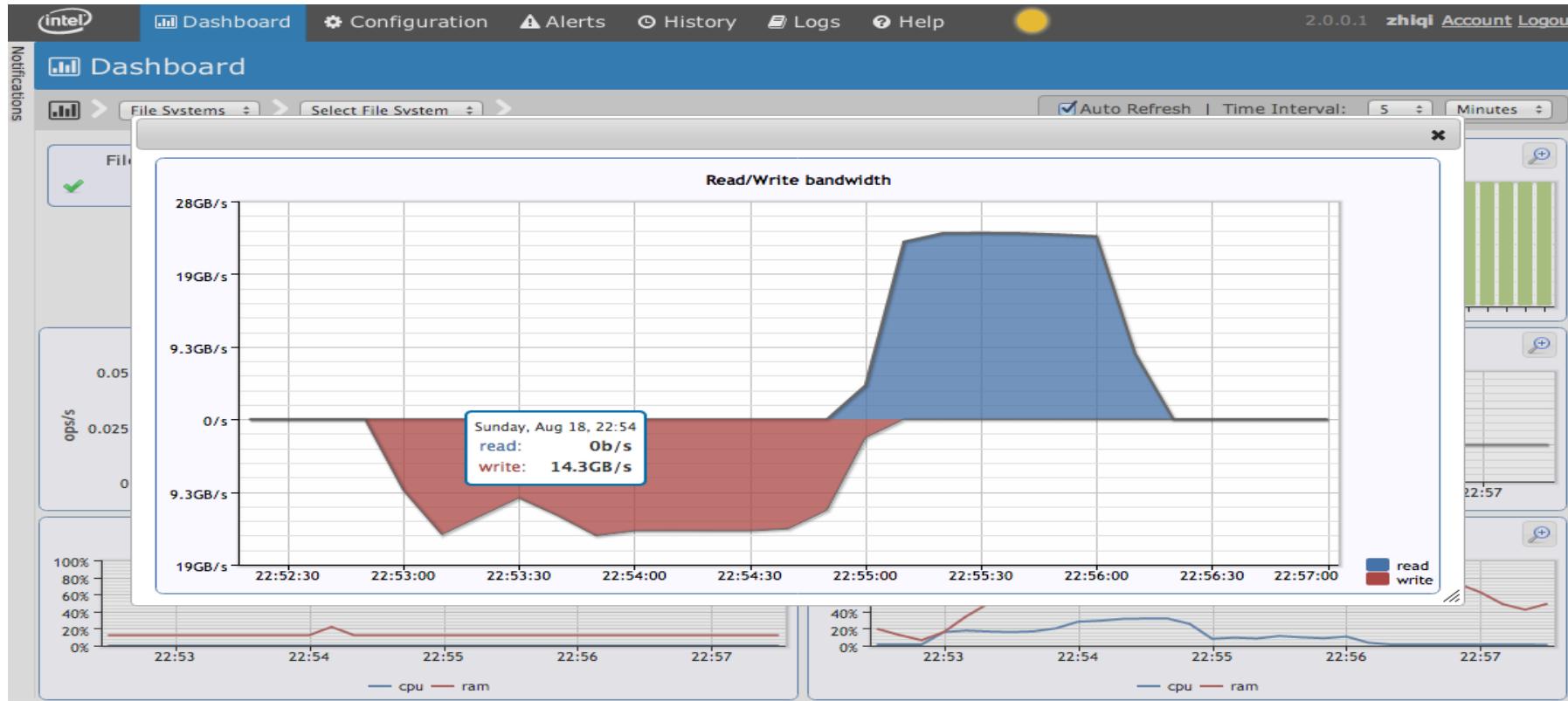
# IEEL配置案例



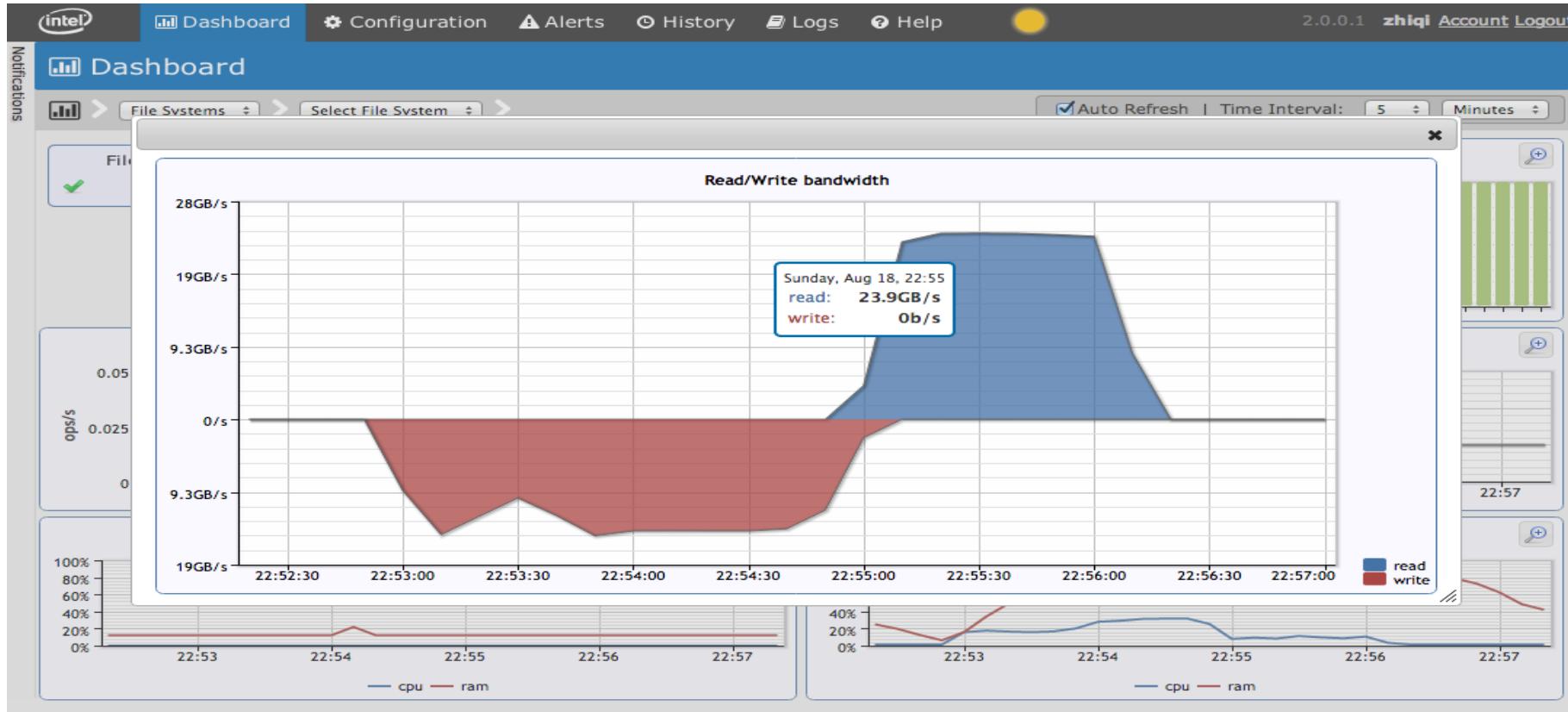
regal-ost-storage-0	28
regal-ost-storage-1	23
regal-ost-storage-2	17
regal-ost-storage-3	13
regal-ost-storage-4	7
regal-ost-storage-5	2
regal-oss00	30
regal-oss01	29
regal-mdt-storage	32
regal-mds0	34
regal-mds1	33
IB Switch	36
regal-iml	39
MGMT Net Swtich	41
MGMT Net Swtich	42
	43
	44

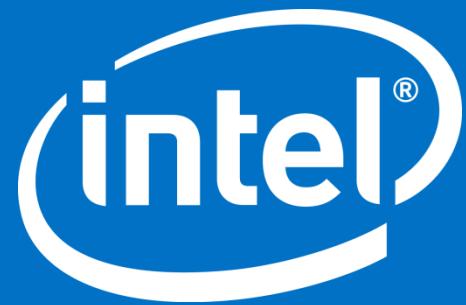
- 容量
  - 1050 TB 可用空间
  - 11亿文件
- 性能
  - 写 13GB/s
  - 读 20GB/s
- 元数据操作速率
  - Create 60K
  - Lookup 220K
  - Getattr 98K
  - Setxattr 13K
  - Destroy 33K
- 高可靠性
  - Active-Standby MDS
  - Active-Active OSS
  - Active-Active Controller
  - RAID10/RAID6

# IEEL配置案例 – 写性能展示



# IEEL配置案例- 读性能展示





Lustre.Intel.Com