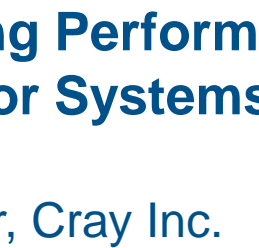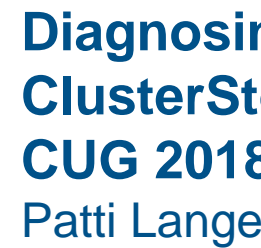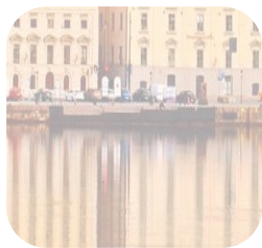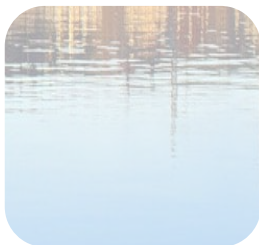# Diagnosing Performance Issues on Cray ClusterStor Systems
## CUG 2018

Patti Langer, Cray Inc.

# There must be a better way

COMPUTE | STORE | ANALYZE

# Topics

- **Overview of Cray® View for ClusterStor™**
- **Review of customer reported problem**
- **Using View for ClusterStor**
- **Summary**
- **Q&A**

# Components of View for ClusterStor

## Run-time Variability

Real-time and historical views of data to understand what is impacting a users job(s)

## Problem Resolution

A unified view of system activity provides administrators with the ability to pinpoint problem areas within their systems

## Trend Analysis

Data-driven analysis and visualization from historical data helps identify trends that can then be used to shape changes to the system

## Alerting

Threshold engine enables customized alerts based on any metric

# Customer Reported Problem

# Overview of Reported Problem

- **The problem**
  - 100% utilization of MDS which caused significant performance degradation
  - Impacting both users and system throughput

- **The complexities**
  - Site has both a Cray XC and Cray Cluster System attached to storage

- **The cost**
  - Months of time to debug and find the root cause

# Problem Identification

- **Problem isolation**
  - Several tests are run, isolating issue to stdout redirection to Lustre
  - A reproducible test case is created
  - Workaround is to redirect stdout to non-Lustre filesystem

- **Cray engineer engaged**
  - Information requested to determine MDS performance and throughput

# Problem Isolation

- ## Initial results
  - Requests were being processed….slowly
    - No lock contention
    - Request queues not backed up
  - Problem not specific to the MDS

- ## Further information requested and analyzed
  - Metadata operation statistics are collected from MDS
  - Information manually correlated with poorly performing job

# Summary of Metadata Operations for Job

| Operation | Count w/out workaround | Count with workaround | Change |
|-----------|------------------------|-----------------------|--------|
| Open | 9883 | 4135 | 239% |
| Close | 9575 | 4078 | 235% |
| Unlink | 6024 | 961 | 627% |
| Mkdir | 2000 | 4 | 50000% |
| Rmdir | 2000 | 3 | 66667% |
| Getaddr | 131598 | 31116 | 423% |
| Statfs | 2009 | 201 | 1000% |
| **Sync** | **830725** | **0** | **infinite** |

COMPUTE | STORE | ANALYZE

# Root Cause Identified

- **Large number of sync operations**
  - 4600 syncs per second
  - With a total of 830725 sync operations

- **Causing 100% utilization of the MDS**

# Challenges

- **Working in a complex environment**
  - Required running the reproducible test case several times to isolate the critical issue

- **Involvement from multiple teams**
  - Additional overhead with communication and data analysis

- **Time to root cause analysis**
  - From problem identification to root cause took 5 months

# View for ClusterStor : Bringing the Pieces Together

COMPUTE | STORE | ANALYZE

# Problem Isolation

- **The *Administrator* is notified of performance degradation**

# Overall System Performance of ClusterStor

# View for ClusterStor Home Page

COMPUTE | STORE | ANALYZE

# Job Summary Table

Copyright 2018 Cray Inc.

# Job Detail Information for 2183675



*Metadata operations increase with job start*

Lustre Job Stats for system snx11253 Job: 2183675 - Interval 1m

Job Write Rates

Job Read Rates

Job Metadata Operations

COMPUTE | STORE | ANALYZE

# Job Detail Information for 2183675



Breakdown Metadata Operations for system snx11253 Job: 2183675

Metadata Operations

| | total ▾ |
|---|---|
| 2183675 MDT0000 sync | 1.2181 Mil |
| 2183675 MDT0000 getattr | 382 |
| 2183675 MDT0000 setattr | 256 |
| 2183675 MDT0000 open | 256 |
| 2183675 MDT0000 close | 89 |
| 2183675 MDT0000 unlink | 0 |
| 2183675 MDT0000 mkdir | 0 |
| 2183675 MDT0000 rmdir | 0 |
| 2183675 MDT0000 getxattr | 0 |

# Root Cause Identified

- **Large number of sync operations**
  - ~6000 syncs per second
  - With a total of 1.2M sync operations

- **Causing 100% utilization of the MDS**

# Bringing the Pieces Together

- **Data available to the Administrator**
    - View collects and correlates information from multiple sources
    - No need for root access to ClusterStor system

- **Reduce need to run reproducible test case**
    - Information persisted and available near real-time and historical

- **Reduce need to engage an expert**

- **Reduce time from problem identification to root cause**

# Summary

- **It's all about enabling *Administrators* to better understand application storage performance**

- **View for ClusterStor enables *Administrators* to**
  - Proactively monitor and understand performance trends
  - Shorten time from problem identification to root cause
  - Improve system availability

# Legal Disclaimer

*Information in this document is provided in connection with Cray Inc. products. No license, express or implied, to any intellectual property rights is granted by this document.*

*Cray Inc. may make changes to specifications and product descriptions at any time, without notice.*

*All products, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.*

*Cray hardware and software products may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.*

*Cray uses codenames internally to identify products that are in development and not yet publicly announced for release. Customers and other third parties are not authorized by Cray Inc. to use codenames in advertising, promotion or marketing and any use of Cray Inc. internal codenames is at the sole risk of the user.*

*Performance tests and ratings are measured using specific systems and/or components and reflect the approximate performance of Cray Inc. products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance.*

*The following are trademarks of Cray Inc. and are registered in the United States and other countries: CRAY and design, SONEXION, URIKA and YARCDATA. The following are trademarks of Cray Inc.: CHAPEL, CLUSTER CONNECT, CLUSTERSTOR, CRAYDOC, CRAYPAT, CRAYPORT, DATAWARP, ECOPHLEX, LIBSCI, NODEKARE, REVEAL. The following system family marks, and associated model number marks, are trademarks of Cray Inc.: CS, CX, XC, XE, XK, XMT and XT. The registered trademark LINUX is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. Other trademarks used on this website are the property of their respective owners.*

COMPUTE | STORE | ANALYZE