# Automated AI-Analysis of the Lustre-Development Mailing List
## (and TASSI)

John Bent, Dominic Manno

May 2024

LA-UR-23-33506

# Suspected Motivation: Current Challenges with HPC Storage Systems



- Very long release cycles

- Extremely long resolution periods

- Lack of diagnostic tools in deployed systems

- Onerous requirements of arcane expertise

# Biggest challenges in distributed storage systems?

# Biggest challenges in distributed storage systems?

foreach thread:
what is topic?
was it answered?

**[analysis]**

**Los Alamos**
NATIONAL LABORATORY

# Automated AI-Analysis of Lustre-Devel



lustre_devel discussions

Around 1.8.X, Lustre became mature.

# For Those of Us Who Like Pie

# For Those of Us Who Like Bars



Percent Answered per Event Type

# Summary of Mailing List Analysis

This is not to bash Lustre!

Lustre is great!

I want to expand the analysis for more systems.

Configuration and deployment are challenging.

# Why Isn't This a Solved Problem?

- Armies of CSP developers

- Armies of Enterprise developers
  - DDN, HPE, IBM, MinIO, Panasas, Pure, VAST, Weka

# The Armies Seem to be Producing Closed Solutions



**Glenn K. Lockwood**
@glennklockwood

I thought it was important to explain that, even though it's called Lustre, we shouldn't limit our thinking on how it can be used to traditional on-prem Lustre. Its real value is in the Azure special sauce.

**Dilip Sundarraj** @dsundarraj · Aug 3
My esteemed colleague @glennklockwood describes how Azure Managed Lustre compounds the benefits of Lustre and the cloud operational semantics to deliver on some of the key requirements of parallel file systems for HPC & AI workloadstechcommunity.microsoft.com/t5/azure-high-... #azure #amlfs #lustre

The full article is a quick informative read: LINK

# How Can We Improve the Lives of Current Admins?

## TASSI: Tool for Agile Scalable Storage Infrastructure

Automation for configuration, deployment, testing of distributed storage

*"I find a bug in version V. I have a reproducer R. Someone submits patch P."*

*"How can I simply spin up a small system*
*running version V, apply patch P, and test with R?"*

*– Thomas Bertschinger, LANL*

# TASSI: Basic Workflow

**TASSI**
- Control agent

**Admin**
- Human in control

**Git Repo**
- Store configs, tests, outputs

**Lustre.org and ZFS**
- Software repos

**Virtual Cluster**
- Connected set of VMs

**libvirt**
- VM management software

**NFS**
- Stash precreated images

TASSI

Admin

Git repo    Lustre.org    ZFS    Virtual Cluster    libvirt    NFS

Los Alamos
NATIONAL LABORATORY

# TASSI Current Actual Config File

### Currently a bit obfuscated

```yaml
all:
  vars:
    ansible_user: root
    ansible_ssh_common_args: '-o UserKnownHostsFile=/dev/null'
    ansible_playbook_install: './ansible/install_all.yaml'
    ansible_playbook_config: './ansible/configure_all.yaml'
    ansible_playbook_test: './ansible/test_lustre.yaml'
    test_script: './tests/simple_mpi_ior.sh'
    vm_dir: '/mnt/usrc-storage-nfs/jbent/images'
    bootstrap_vm:
      cpus: 2
      boot_hdd_gbs: 12
      memory_mbs: 4096
      root_pwd: password
      location: 'http://mirror.centos.org/centos/8-stream/BaseOS/x86_64/os/'
      auth_keys: '/home/jbent/.ssh/authorized_keys'
    network:
      addr: '192.168.56'
    lustre:
      mgs_node: '192.168.56.10@tcp:192.168.56.20@tcp'
      version: '2.15.4-RC2'
      backfstype: 'zfs'
      patch: '/mnt/usrc-storage-nfs/jbent/patches/lustre/test_patch.patch'
      repo: 'git://git.whamcloud.com/fs/lustre-release.git'
    zfs:
      version: zfs-2.1.11
      patch: '/mnt/usrc-storage-nfs/jbent/patches/zfs/zfs_patch1.patch'
      repo: 'https://github.com/openzfs/zfs.git'
  children:
    clients:
      vars:
        configure_args: '--disable-server --enable-client'
      hosts:
        client00:
          ip: 50
          target_mount: '/mnt/lustre'
          hds: []
        client01:
          ip: 60
          target_mount: '/mnt/lustre'
          hds: []
    servers:
      vars:
        configure_args: '--with-zfs --disable-ldiskfs --enable-server'
      children:
        mds:
          vars:
            target_type: mdt
```

Los Alamos
NATIONAL LABORATORY

# TASSI Config File Essential Elements

**Host OS**
- Base image (e.g. CentOS 8)
- RAM

**Lustre Software**
- Repo location (e.g. Whamcloud)
- Version / tag
- Patches
- Backend (e.g. ZFS)

**ZFS Software**
- Repo location
- Version / tag
- Patches

**Cluster**
- Num clients, MDS, OSS, MDT, OST
- HDD sizes for each target

**Test Script**

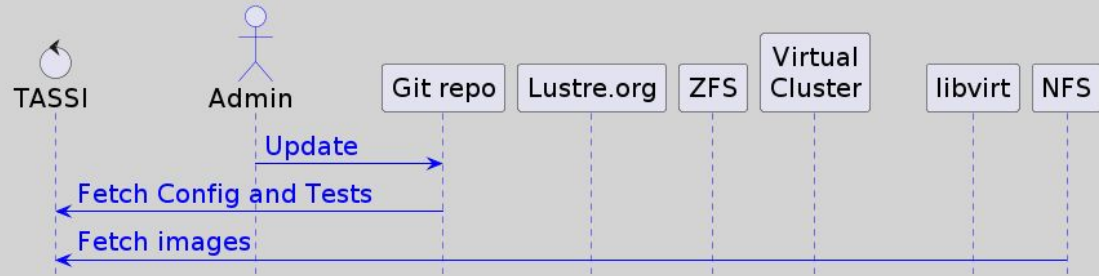# TASSI: Basic Workflow

## Step One

- Admin updates a config

# TASSI: Basic Workflow

**Step One**
- Admin updates a config

**Step Two**
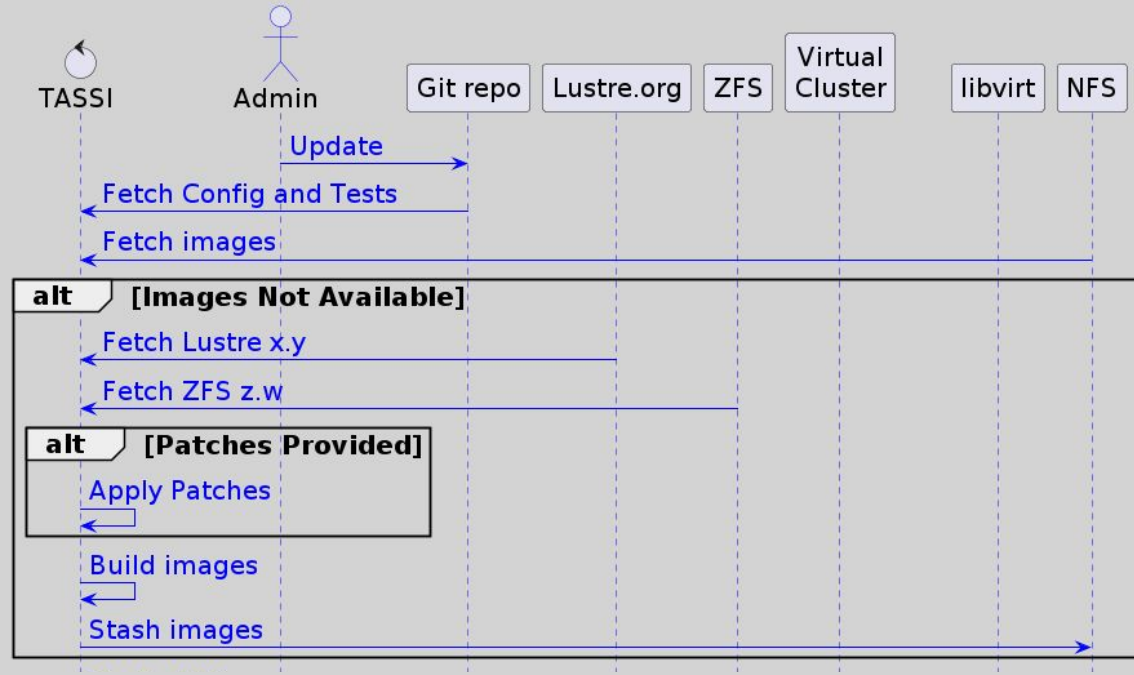- Fetch precreated images

# TASSI: Basic Workflow

**Step One**
- Admin updates a config

**Step Two**
- Fetch precreated images
- Build if not precreated
  - Apply any specified patches

*(not shown, spin up initial bootstrap VM)*
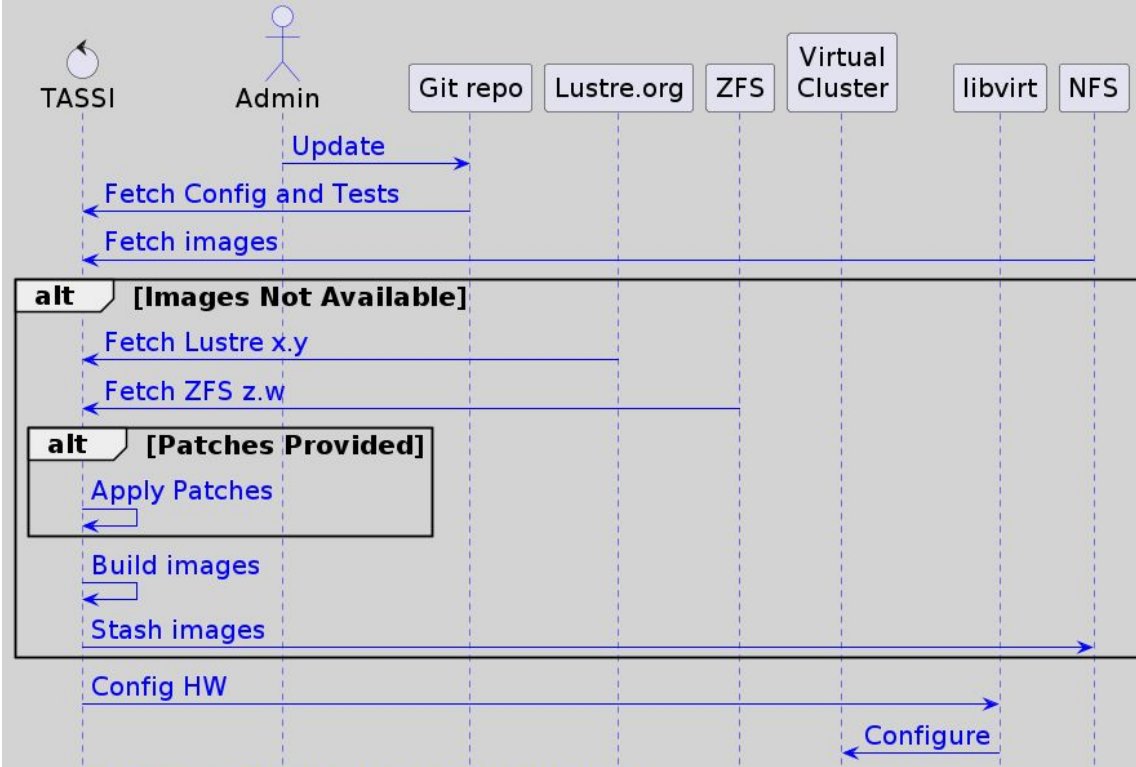
# TASSI: Basic Workflow

**Step One**
- Admin updates a config

**Step Two**
- Fetch precreated images
- Build if not precreated
  - Apply any specified patches

**Step Three**
- Setup the "physical" cluster



Participants: TASSI, Admin, Git repo, Lustre.org, ZFS, Virtual Cluster, libvirt, NFS

- Admin → Git repo: Update
- TASSI ← Admin: Fetch Config and Tests
- TASSI → : Fetch images

**alt [Images Not Available]**
- TASSI ← Lustre.org: Fetch Lustre x.y
- TASSI ← ZFS: Fetch ZFS z.w

  **alt [Patches Provided]**
  - TASSI: Apply Patches

- TASSI: Build images
- TASSI → NFS: Stash images

- TASSI → libvirt: Config HW
- Virtual Cluster ← libvirt: Configure

# TASSI: Basic Workflow

**Step One**
- Admin updates a config

**Step Two**
- Fetch precreated images
- Build if not precreated
  - Apply any specified patches

**Step Three**
- Setup the "physical" cluster
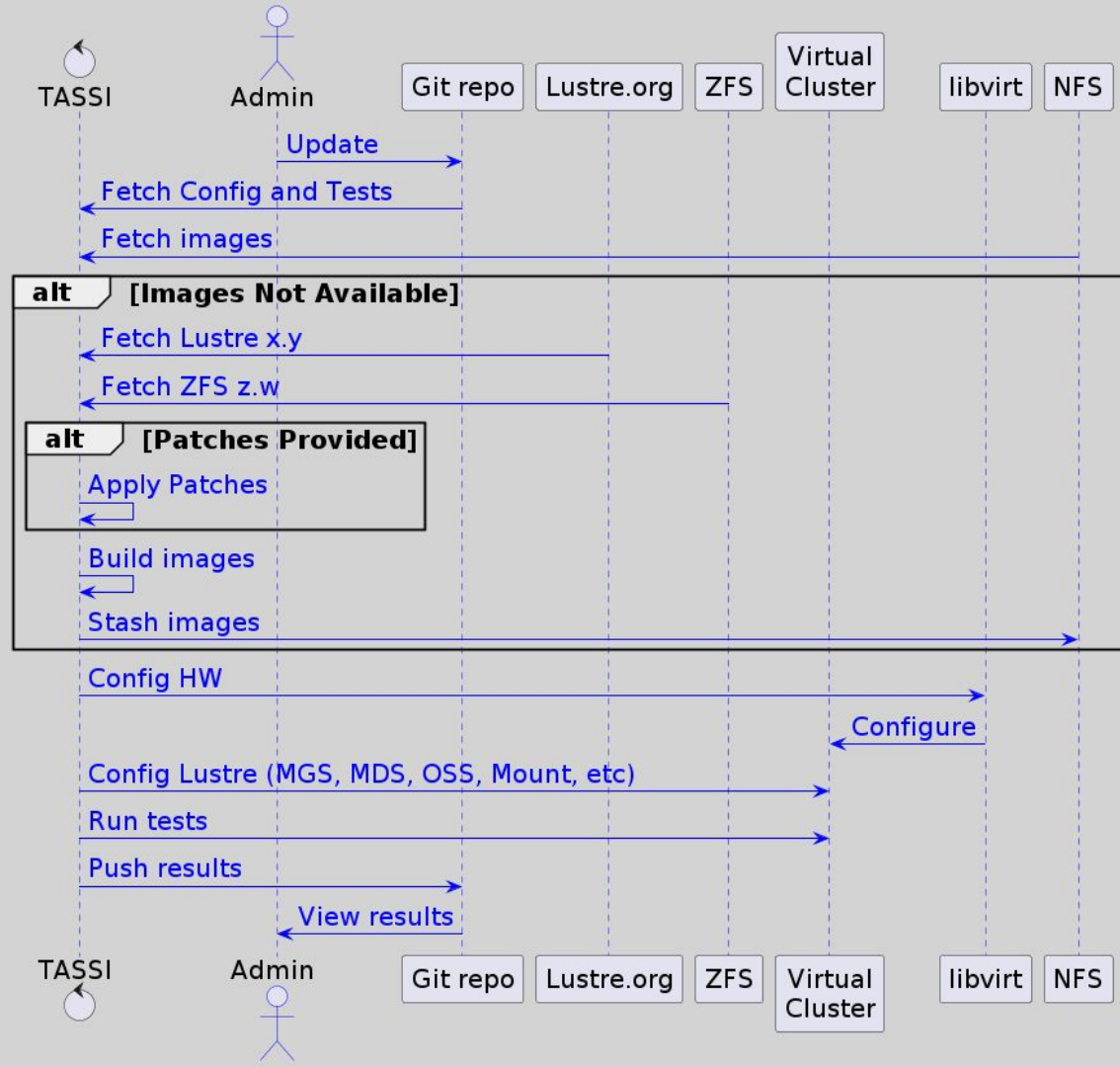
**Step Four**
- Setup Lustre

**Step Five**
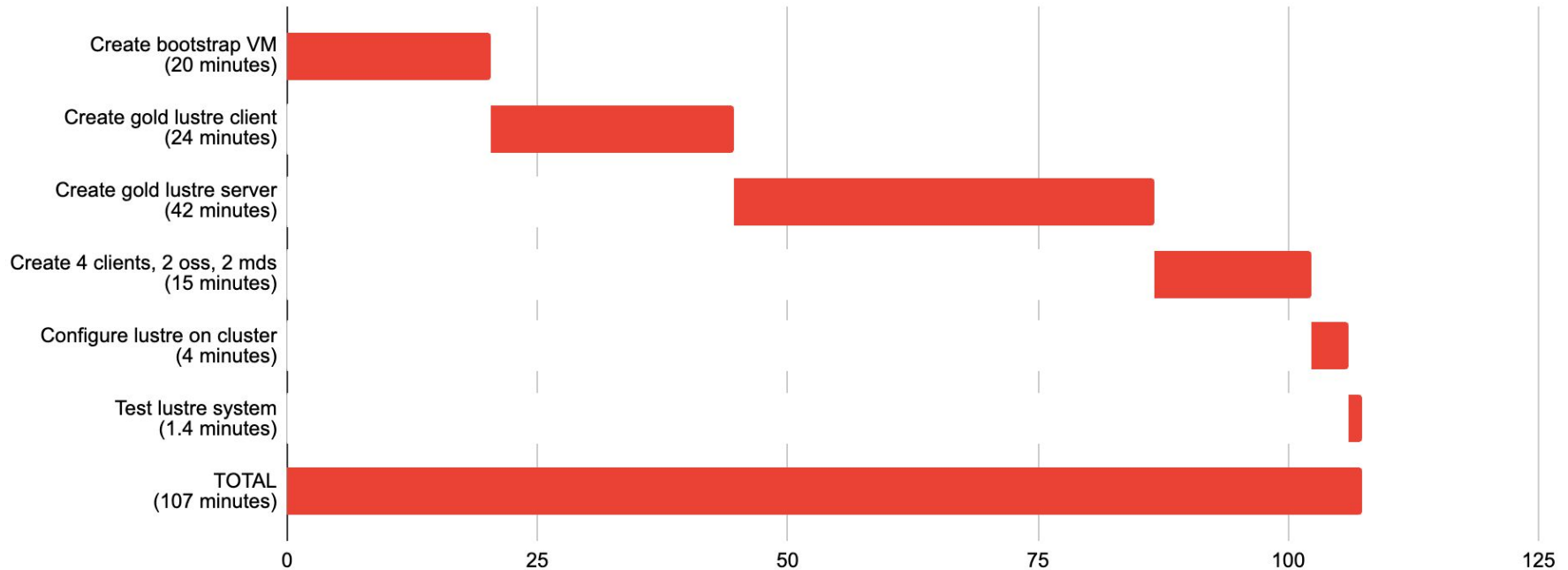- Run the specified tests

**Step Six**
- Commit the outputs

# TASSI Tools Used

- Control agent
  - Python
- Virtual Machine Management
  - libvirt/KVM/Qemu
- Lustre and ZFS installation
  - Ansible
- Config repository
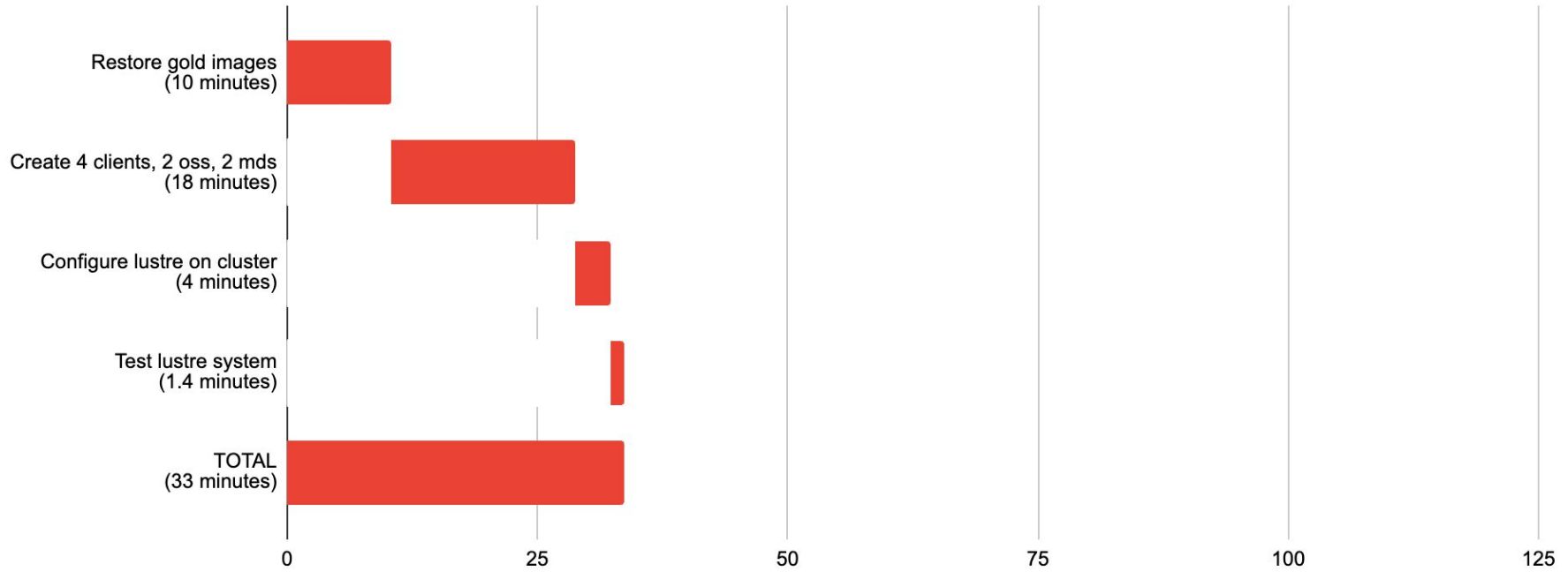  - Git
- Image stashing
  - NFS

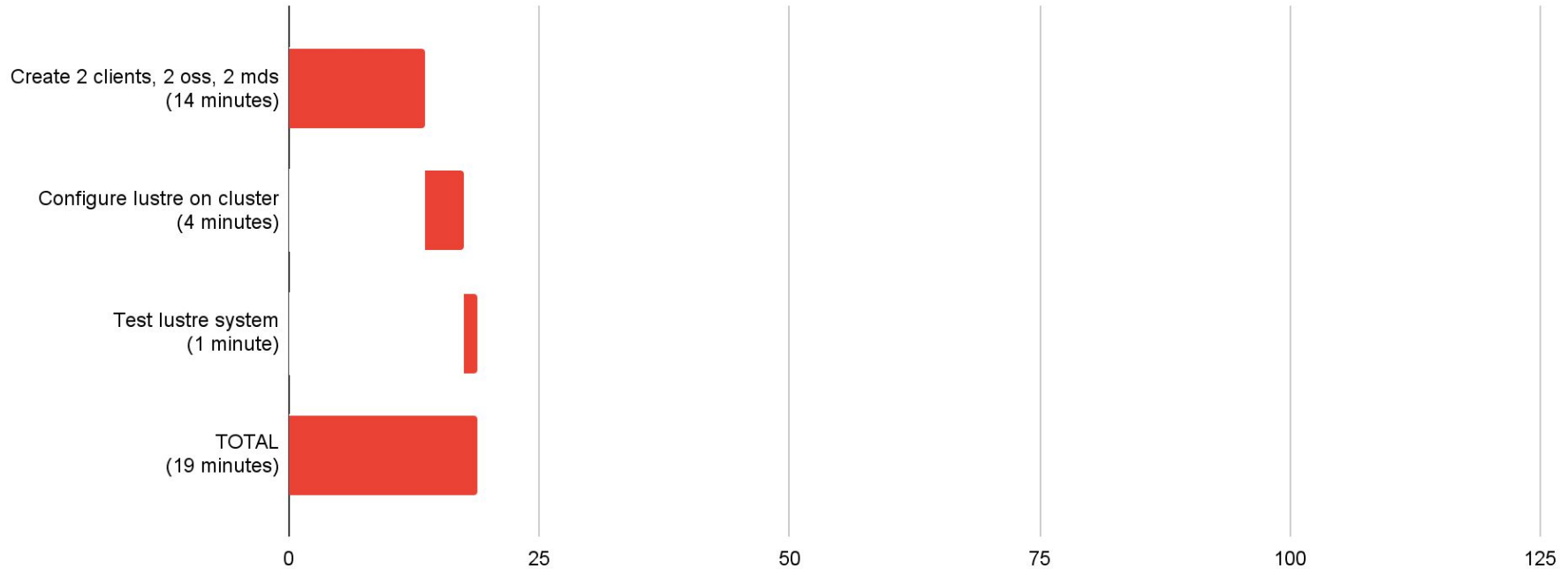# TASSI Timings - 4 Clients - No Gold: 107 minutes

## 2 MDS, 2 OSS

# TASSI Timings - 4 Clients - Cold Gold: 33 minutes
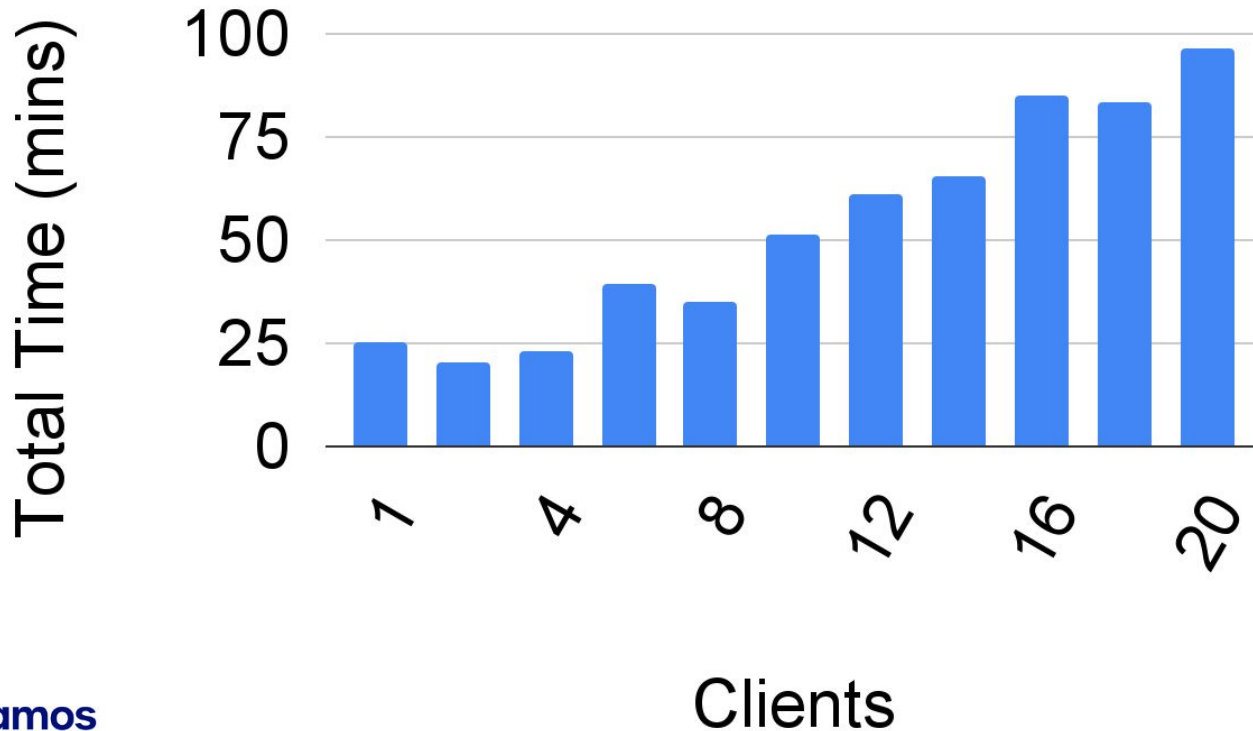
## 2 MDS, 2 OSS

# TASSI Timings - 4 Clients - Warm Gold: 19 minutes

### 2 MDS, 2 OSS

# TASSI Timings: Cold Gold Client Scaling

**2 MDS, 2 OSS**

# TASSI Future Work

- Lustre specification
  - E.g. test with ldiskfs backend also
- Performance
  - E.g. parallelize node creation, use multiple physical nodes
- CI/CD
  - E.g. convert current git manual trigger to actual CI/CD mechanism
- Cluster deployment / configuration
  - Use CSPs, LANL OCHAMI, etc to dynamically provision/configure cluster
- File system installation, configuration
  - E.g. hardware accelerators, specify disk locations
- File system support
  - Add more (HammerSpace? BeeGFS? DAOS?)
- Storage system support
  - Add KV as opposed to just file system (KV-CSD? RocksDB?)
- Decomposability and disaggregation
  - Via NVMeoF for example

# TASSI Conclusions

- Open source tools
- Flexible agile dynamic storage
- Enabling confident happy sys admins

*Creating high performance, reliable storage systems for end-users*

Email jbent@newmexicoconsortium.org and dmanno@lanl.gov for more info, feedback.

Or even to get involved!