# RAND
Safety and Justice Program

# *Autopsy as a Service – Distributed Forensic Compute That Combines Evidence Acquisition and Analysis*

## Presentation to OSDFCon 2016

**Dan Gonzales, Zev Winkelman, John Hollywood, Dulani Woods, Ricardo Sanchez, Trung Tran**

## October 2016

# *Objective and Background*

- **RAND has been funded by the National Institute of Justice to accelerate the processing of digital forensics data**

- **Objective: Develop a Digital Forensics Compute Cluster (AutopsyCluster)**
  - Based on open source, state of the art software
  - Reduce processing time and storage costs

- **We have chosen Autopsy as a core component of AutopsyCluster**
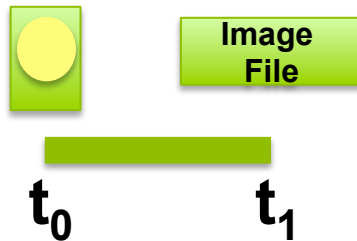  - "Autopsy as a Service"

RAND

# *Vision*

- **Provide law enforcement with a cost effective and efficient digital forensics analysis capability**

- **Combine data ingest and analysis steps to speed up the digital evidence analysis process using**
  - **Distributed computing tools**
  - **Cloud computing services**

- **Approach designed to**
  - **Reduce infrastructure cost**
  - **Stand up infrastructure only when needed**
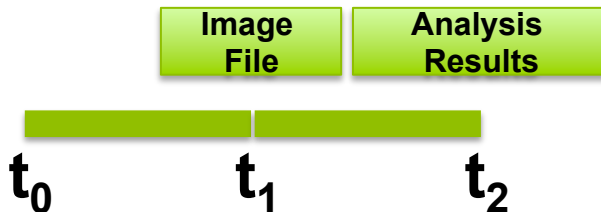  - **Access infrastructure to perform multiple analyses in parallel**

RAND

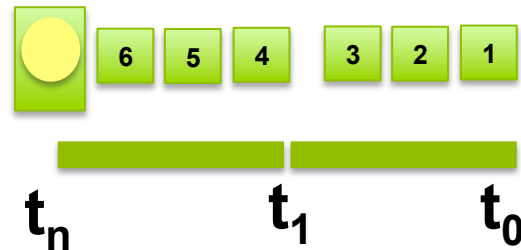# *To implement the Vision We Stream Data into the Cloud*

## Old Way

- **Step 1: make copy**

| Image File |

$t_0$       $t_1$

- **Step 2: analyze image on standalone workstation**

| Image File | Analysis Results |

$t_0$      $t_1$      $t_2$

## New Way

- **Step 1: start stream**

| 6 | 5 | 4 | 3 | 2 | 1 |

$t_n$      $t_1$      $t_0$

- **Step 2: process stream on the fly in micro batches**

Byte 0          Byte N

| File 1 | Unallocated |
| File 2 | File 3 |

Batch 1 @t1    Batch 2 @t2    Batch N @tn

**If we can keep up with the data coming off the disk, we are processing as fast as is physically possible**

RAND

# *Outline*

- **Objectives and vision**

- ➡ **Architecture**

- **Initial results**

- **Lessons Learned**

- **How to use AutopsyCluster**

- **Beta testing**

RAND

# The Forensics Analysis Functions of AutopsyCluster are Based on Autopsy[a]

- **Basis Technology has developed a version of Autopsy for collaborative forensics analysis over a network[b]**
  - **We chose this version because it is designed to work over a network with supporting servers**

- **AutopsyCluster designed to run forensics processing tasks in parallel at near "streaming speed"**
  - **Speed at which disk blocks are read from evidence disk**
  - **With dc3dd with USB 3.0 this is about 15 MBps**

- **We modified the *Autopsy* so it is a streaming application**
  - **Integrated with Apache Spark[c] (cluster computing framework) and Apache Kafka[d] (messaging)**

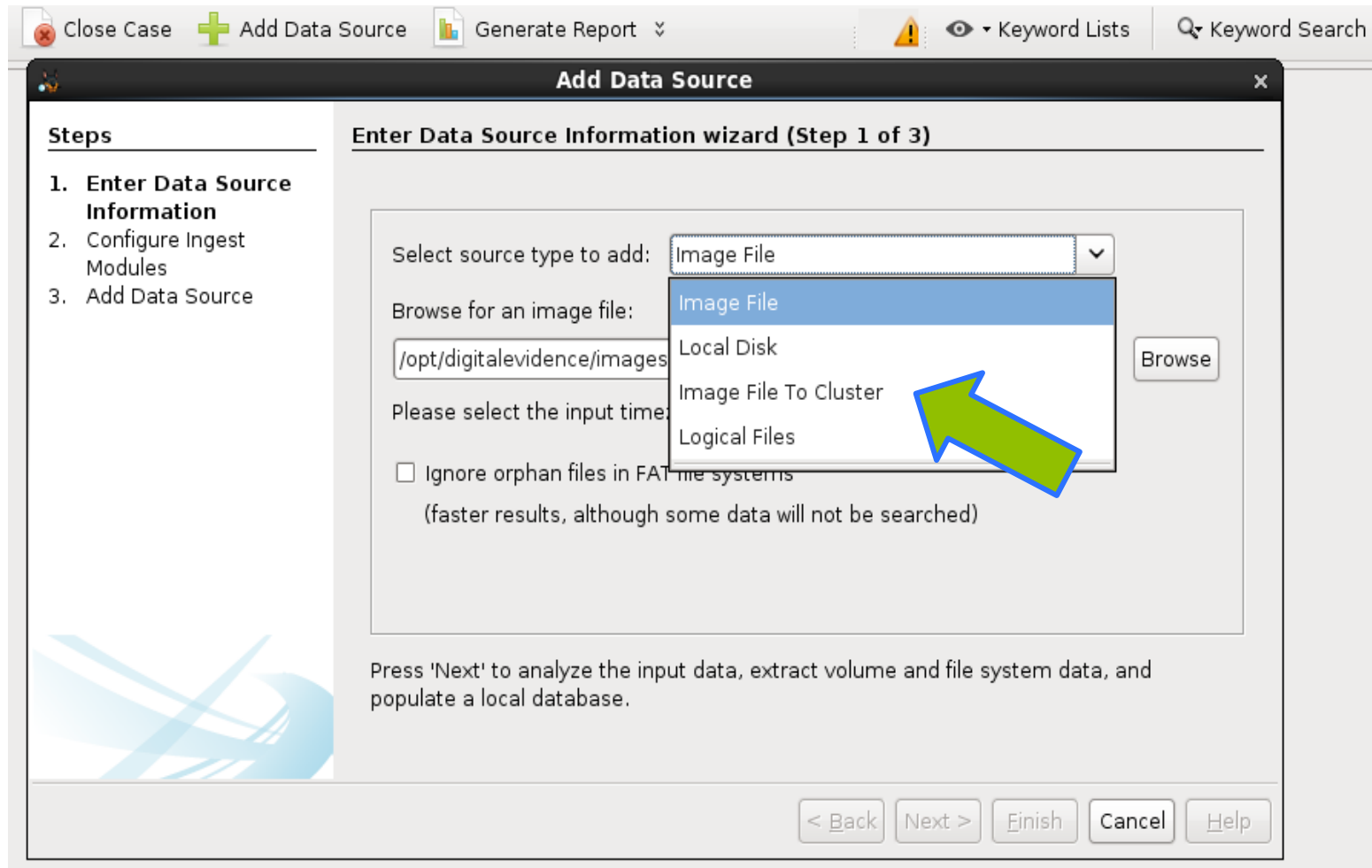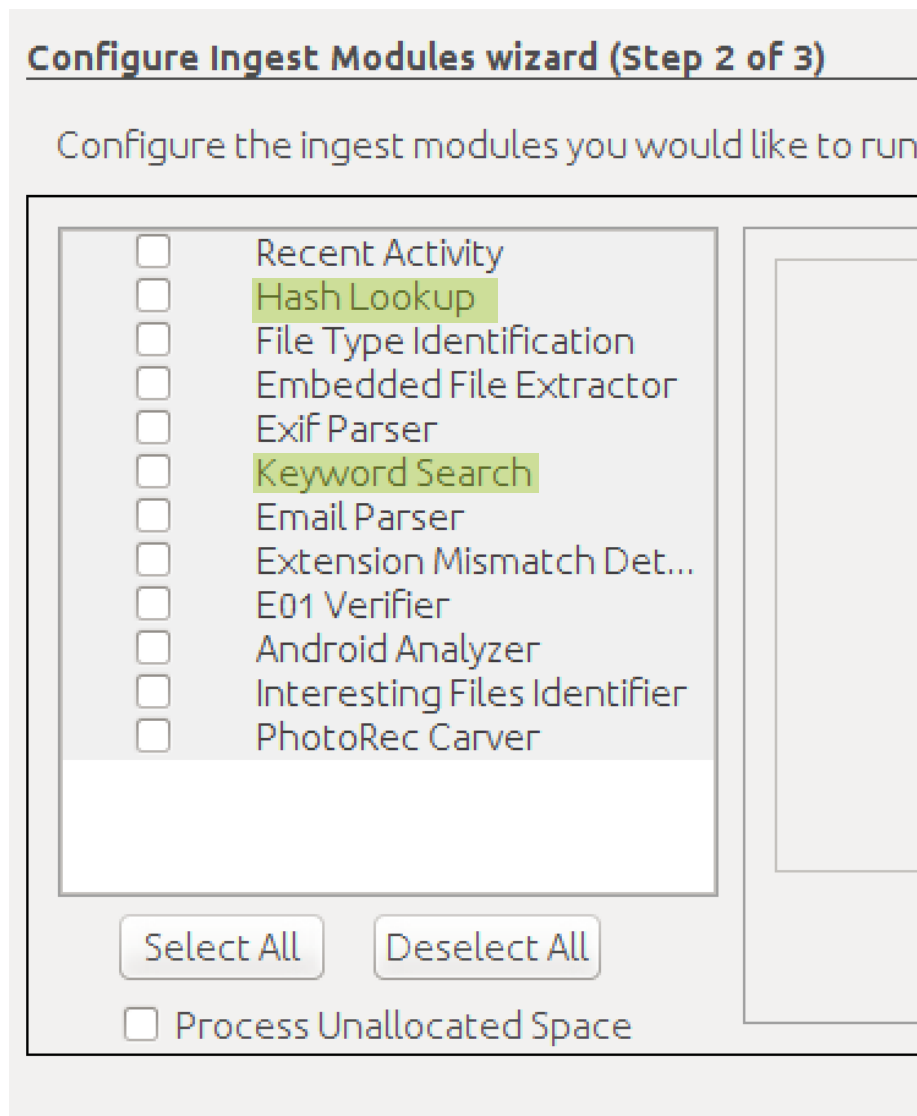- **Autopsy analysis modules read from the stream**

**Autopsy Sleuth Kit**

Kafka

a http://www.sleuthkit.org/autopsy/
b https://github.com/sleuthkit/autopsy
c http://www.sleuthkit.org/autopsy/
d http://www.postgresql.org/

RAND

# User Interface for Autopsy Streaming Branch

# Autopsy Modules For Autopsy Streaming Branch

## Configure Ingest Modules wizard (Step 2 of 3)

Configure the ingest modules you would like to run

- Recent Activity
- Hash Lookup
- File Type Identification
- Embedded File Extractor
- Exif Parser
- Keyword Search
- Email Parser
- Extension Mismatch Det...
- E01 Verifier
- Android Analyzer
- Interesting Files Identifier
- PhotoRec Carver

Select All    Deselect All
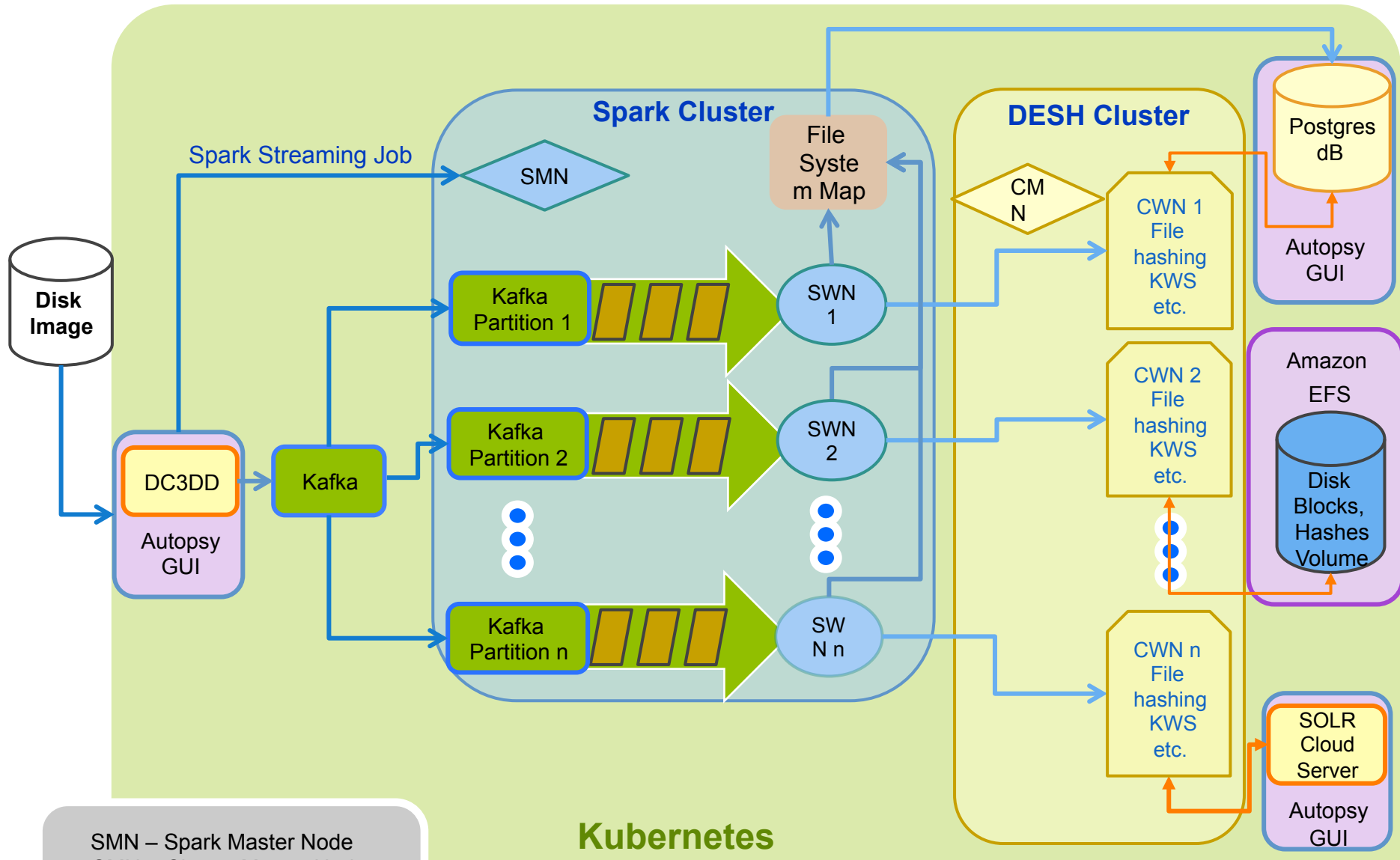
☐ Process Unallocated Space

**Currently Working in Spark:**
- "Hash Lookup"
- "Keyword Search"
- Hardcoded configurations
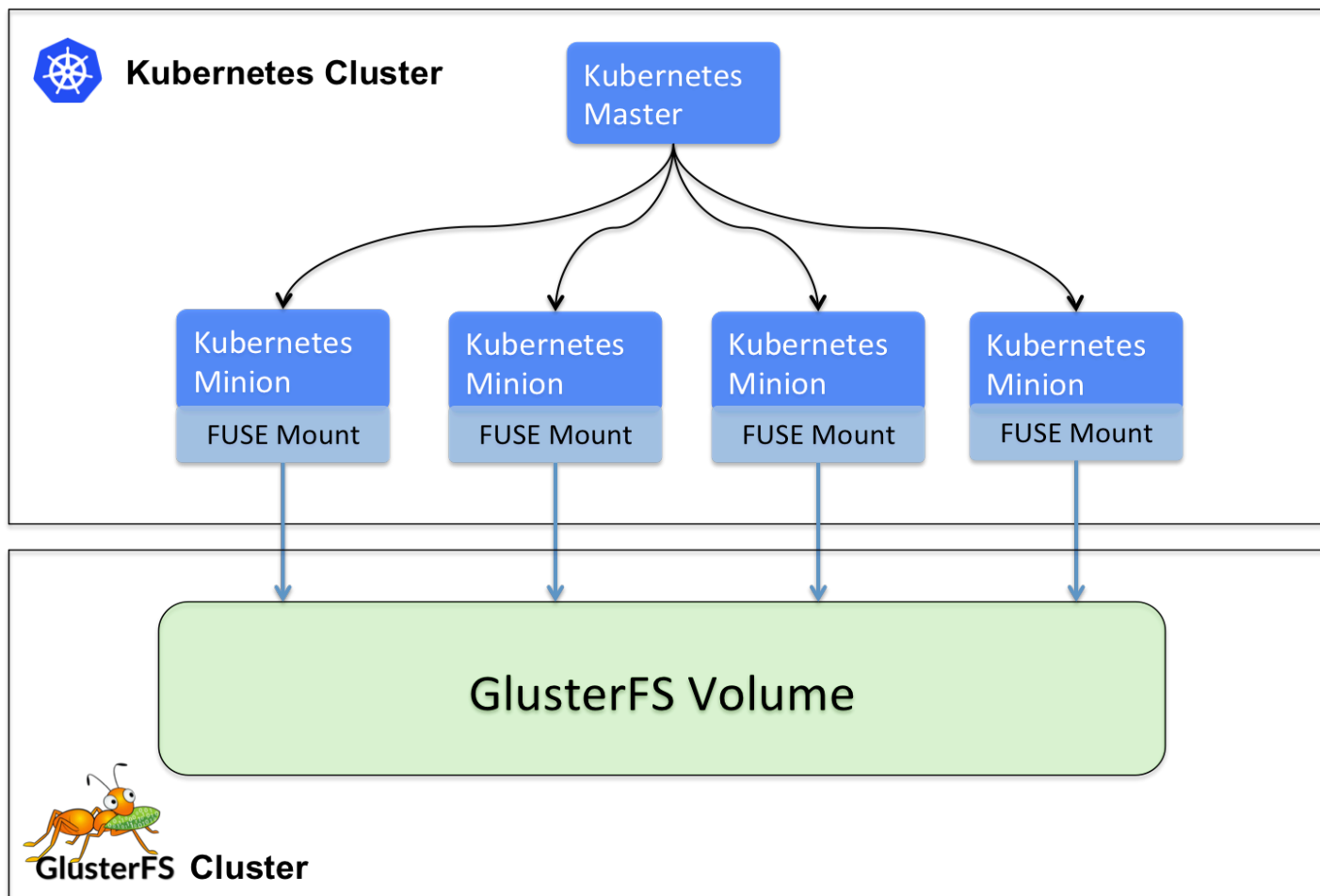
**Next Steps:**
- Remaining modules starting with "Interesting Files Identifier"
- Implement configuration of modules with Autopsy UI

RAND

# AutopsyCluster Architecture



SMN – Spark Master Node
CMN – Cluster Master Node
SWN – Spark Worker Node
CWN – Cluster Worker Node
KWS - Key Word Search

Gonzales and Winkelman-9  October 2016

# *Kubernetes + File Volumes*

# *AutopsyCluster Kubernetes Dashboard*

Workloads

## Replication controllers

| | Name | Labels | Pods | Age | Images |
|---|---|---|---|---|---|
| ✓ | activemq | name: activemq | 1 / 1 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | desh-worker-controller | component: desh-worker | 6 / 6 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | kafka | name: kafka | 1 / 1 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | nfs-server | role: nfs-server | 1 / 1 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | postgres | name: postgres | 1 / 1 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | solr | name: solr | 1 / 1 | an hour | gordianknot.rand.org:5001/desh |
| ✓ | spark-master-controller | component: spark-master | 1 / 1 | an hour | gordianknot.rand.org:5001/gcr.i |
| ✓ | spark-worker-controller | component: spark-worker | 6 / 6 | an hour | gordianknot.rand.org:5001/gcr.i |

## Pods

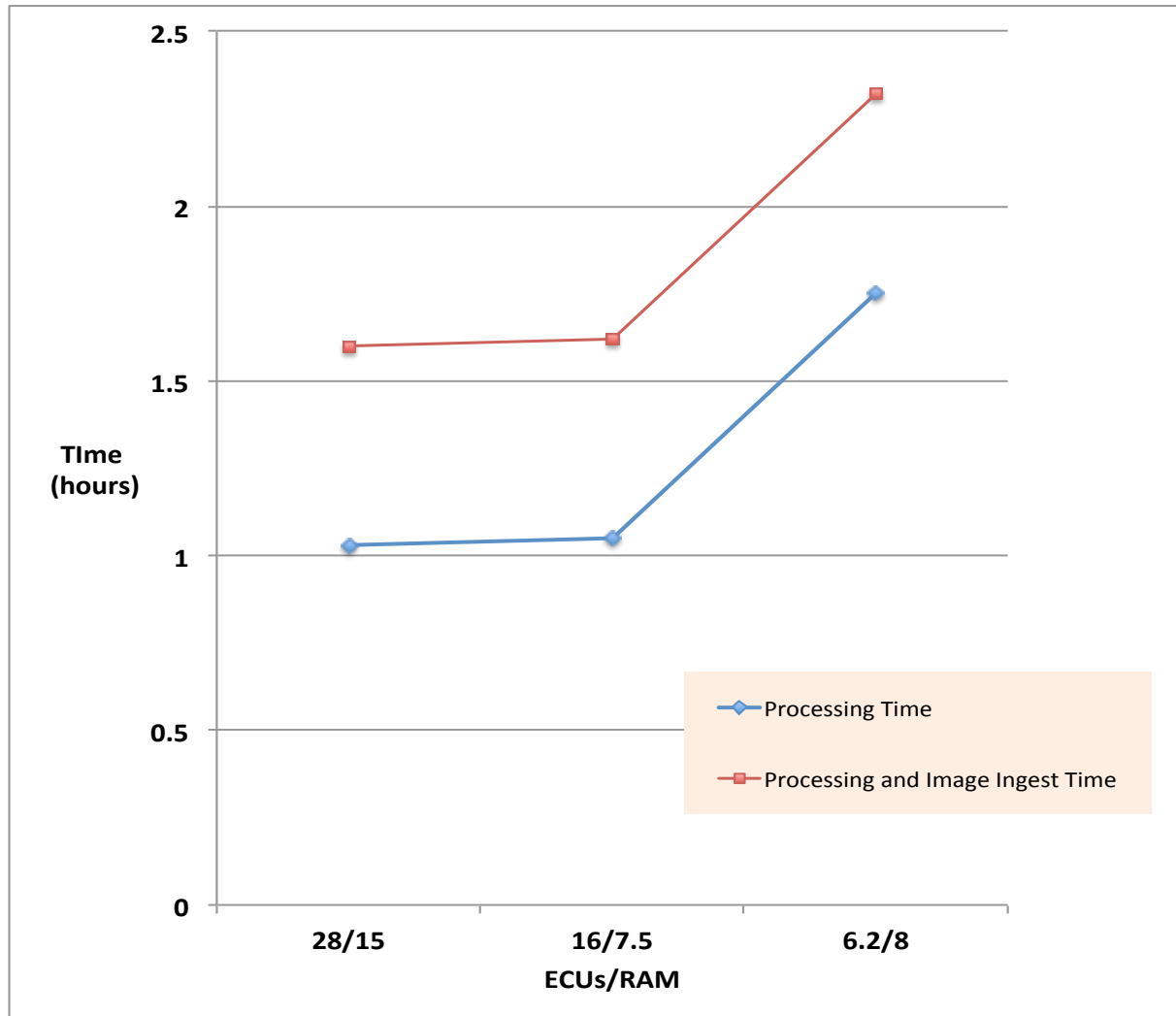| | Name | Status | Restarts | Age | Cluster IP | CPU (cores) | | Memory (bytes |
|---|---|---|---|---|---|---|---|---|
| ✓ | activemq-r9ybs | Running | 0 | an hour | 172.18.5.4 | | 0.001 | 18 |
| ✓ | desh-worker-controller-auu06 | Running | 0 | an hour | 172.18.1.7 | | 0.151 | 1. |
| ✓ | desh-worker-controller-d9b8i | Running | 0 | an hour | 172.18.2.6 | | 0.107 | 2. |
| ✓ | desh-worker-controller-iuztt | Running | 0 | an hour | 172.18.0.6 | | 0.136 | 2. |
| ✓ | desh-worker-controller-u2tlh | Running | 0 | an hour | 172.18.3.6 | | 0.102 | 2. |
| ✓ | desh-worker-controller-xyloz | Running | 0 | an hour | 172.18.5.5 | | 0.12 | 1. |
| ✓ | desh-worker-controller-yom7e | Running | 0 | an hour | 172.18.4.6 | | 0.138 | 1. |
| ✓ | kafka-ri2gi | Running | 0 | an hour | 172.18.5.6 | | 0.003 | 9. |
| ✓ | nfs-server-9ptkv | Running | 0 | an hour | 172.18.4.3 | | 0 | 15 |
| ✓ | postgres-utodq | Running | 0 | an hour | 172.18.1.5 | | 1.084 | 17 |
| ✓ | solr-eh53b | Running | 0 | an hour | 172.18.0.3 | | 0.025 | 28 |
| ✓ | spark-master-controller-hi7zq | Running | 0 | an hour | 172.18.2.4 | | 0.006 | 39 |
| ✓ | spark-worker-controller-c5nca | Running | 0 | an hour | 172.18.0.5 | | 0.008 | 6. |
| ✓ | spark-worker-controller-jh4nu | Running | 0 | an hour | 172.18.5.2 | | 0.016 | 2. |
| ✓ | spark-worker-controller-lkzij | Running | 0 | an hour | 172.18.4.5 | | 0.007 | 6. |
| ✓ | spark-worker-controller-pic1z | Running | 3 | an hour | 172.18.1.6 | | 0.008 | 6. |
| ✓ | spark-worker-controller-poiur | Running | 0 | an hour | 172.18.3.5 | | 0.007 | 6. |
| ✓ | spark-worker-controller-s16ge | Running | 0 | an hour | 172.18.2.5 | | 0.008 | 6. |

RAND

# *Outline*

- **Objectives and vision**

- **Architecture**

➡ - **Initial results**

- **Lessons Learned**

- **How to use AutopsyCluster**

- **Beta testing**

RAND

# *Forensic Images We are Using In Performance Testing*

| Image | Size | Source |
|---|---|---|
| Rhino Hunt | 250 MB | NIST (CFReDS) |
| Data Leakage | 20 GB | NIST (CFReDS) |
| NPS DOMEX Users, 2009 | 40 GB | Digital Corpora |
| NPS 1weapondeletion, 2011 | 75 GB | Digital Corpora |
| NPS 2weapons, 2011 | 253 GB | Digital Corpora |
| NPS 2 TB, 2011 | 2 TB | Digital Corpora |

- **Initial tests conducted on**
  - **Stand alone machines**
  - **A typical RAND server (Digital Evidence)**
  - **Amazon Web Services (AWS)**

RAND

# Stand Alone Autopsy Results on AWS Windows Virtual Machines (VMs)
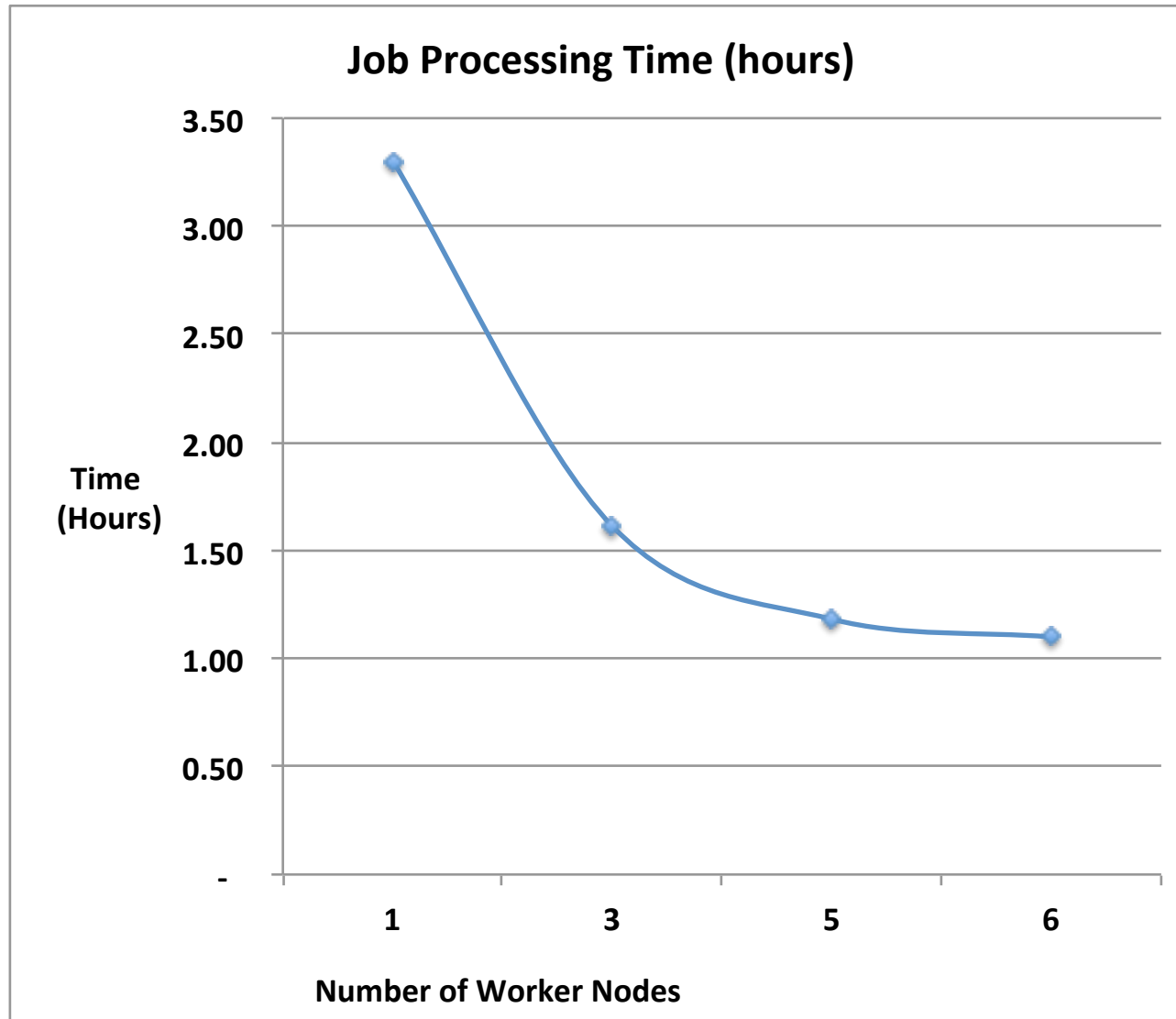


**40 GB Hard Disk Image**

**Ingestion, hashing, Key Word Search**

ECU = Elastic Compute Unit = 2007, 1 GHz CPU

- **Autopsy performances varies based on machine capabilities**
- **All results are for raw HD images already ingested in cloud**

RAND

# *AutopsyCluster Results on a Single Server for a 40 GB Hard Disk Image*

**Job Processing Time (hours)**

Time (Hours)

3.50
3.00
2.50
2.00
1.50
1.00
0.50
-

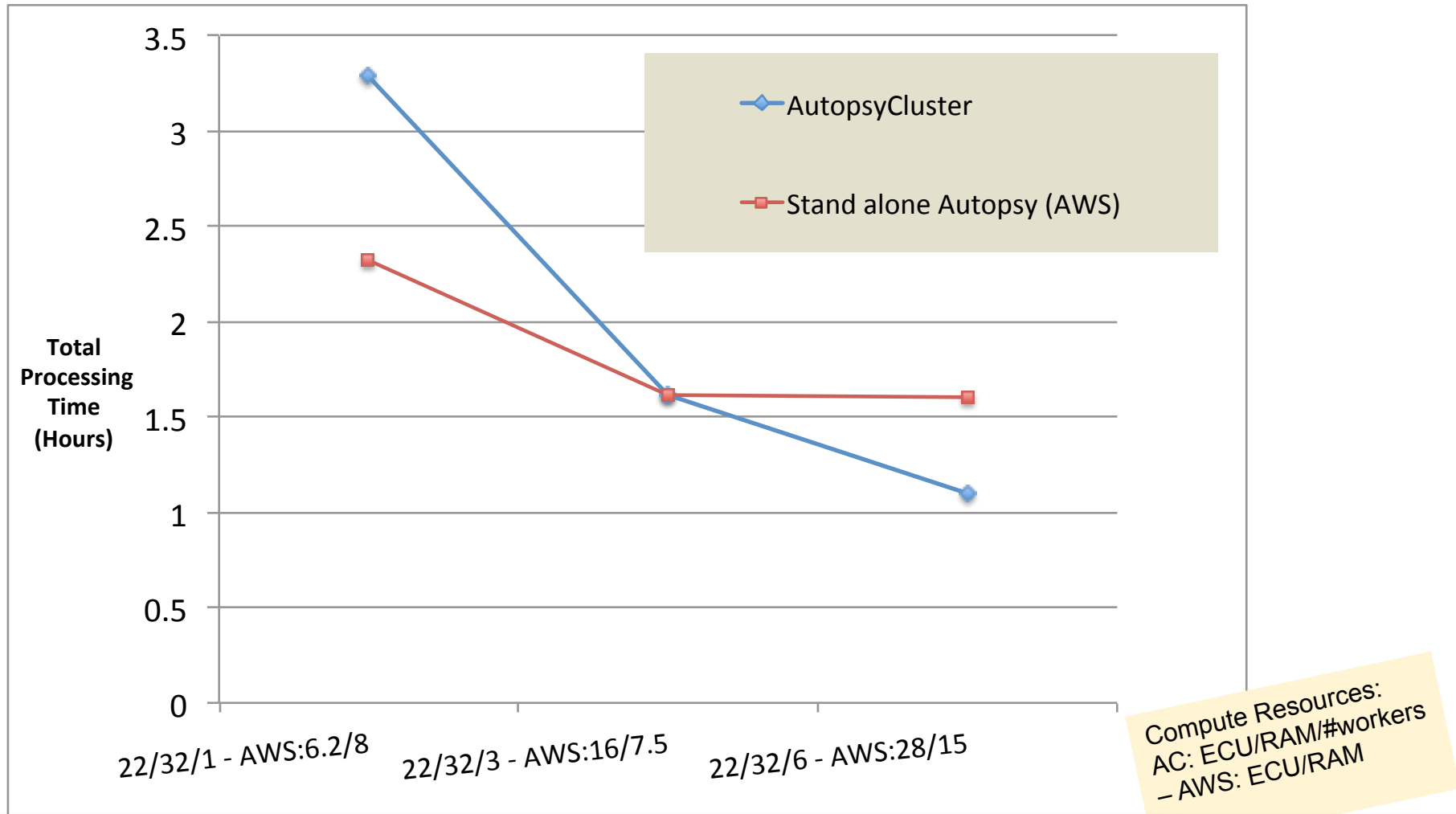1        3        5        6

**Number of Worker Nodes**

**Local server equivalent To 22 ECUs with 32 GB RAM (22/32)**

**Ingestion, hashing, Key Word Search**

**Performance roughly Comparable with stand alone Autopsy With 5 or more worker nodes**
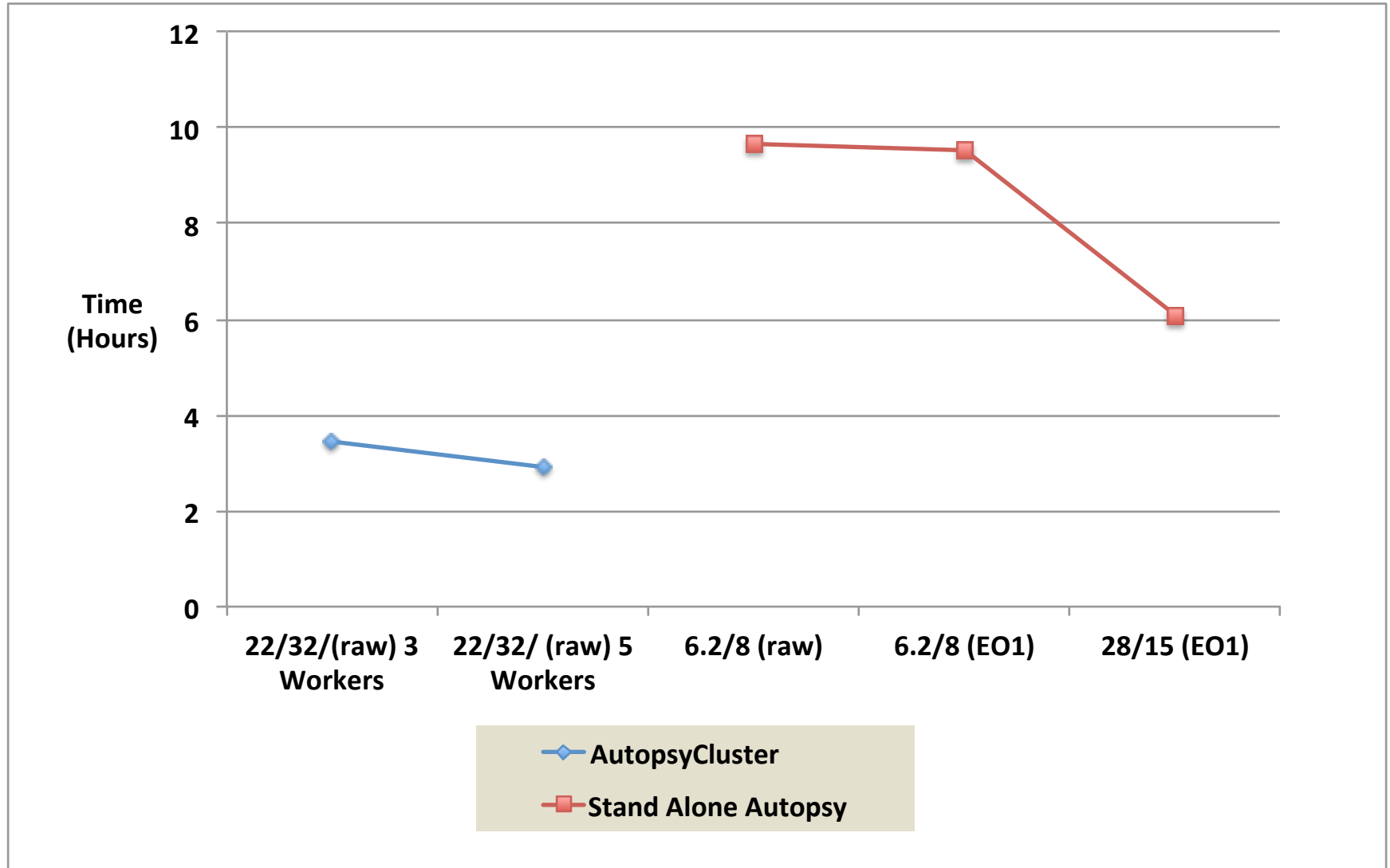
**Number of worker nodes constrained by memory limitations on specific server used**

RAND

# *Stand Alone Autopsy (SAA), AutopsyCluster (AC) Performance Comparison for a 40 GB Drive*

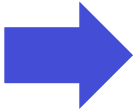

Compute Resources:
AC: ECU/RAM/#workers
– AWS: ECU/RAM

- **As Worker nodes are added to the Server AutopsyCluster Performance improves; With 6 worker nodes AutopsyCluster is faster than Autopsy**

RAND

# Stand Alone Autopsy and AutopsyCluster Results on AWS for 75 GB Disk Images



RAND

# *Outline*

- **Objectives and vision**

- **Architecture**

- **Preliminary test results**

➡ - **Lessons learned**

- **How to use AutopsyCluster**

- **Beta testing**

RAND

# *Moving to the Cloud Can Present a Number of Challenges*

- **Good communications links to the cloud are essential for good performance**

- **Testing at RAND showed that communications links to AWS were frequently congested, adding time delays**

- **It is possible to purchase a direct link to AWS for many ISP links, which may improve performance significantly**
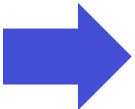
RAND

# *Outline*

- **Objectives and vision**

- **Architecture**

- **Preliminary test results**

- **Lessons Learned**

- **How to use AutopsyCluster**

- **Beta testing**

RAND

# *Four Ways to Use Fully Operational AutopsyCluster*

- **Acquire and ingest locally on a single machine**
  - **Advantage is acquisition and analysis at the same time**

- **Acquire locally and ingest on local private distributed computing (e.g., on premises datacenter)**

- **Acquire locally, ingest remotely (e.g., cloud) and transmit via streaming**

- **Ship drive(s) to cloud service provider for remote acquisition, and multiple side-by-side ingest "jobs"**
  - **We plan to investigate feasibility with AWS**

RAND

# AutopsyCluster Provides Scalable Options for Data Acquisition and Ingest

| Option | Streaming | Distributed | Cloud |
|---|---|---|---|
| Autopsy Standalone | No | No | No |
| AutopsyCluster on premise single machine | Yes | No | No |
| AutopsyCluster on premise data center | Yes | Yes | No |
| Autopsy on premise – remote data center | Yes | Yes | Yes |
| Ship drives for AutopsyCluster processing in Cloud | No | Yes | Yes |

RAND

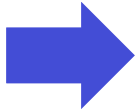# How Much Would Acquisition and Ingest of a 1TB Drive Cost on AWS?

- **Example for a 1 TB drive:**
  - Total hourly rate for 6 nodes (2 CPUs ea, 15GB RAM ea): **$1**
  - Total hourly rate for 6 Linux SSD "disks" (32 GB ea): **$0.03**
  - Total hourly rate for 2 TB of "elastic" storage (need 2x): **$0.83**
  - Run time to extract and stream 1TB at 15MB/s: ~19 hours (includes time for "setup" and "teardown" of the cluster)
- **Total "cloud" cost to acquire and ingest:**
  
  **(1 + 0.03 + 0.83**)/hour * 19 hours = **~$35**
- **Immediate access storage for uncompressed acquired image and case file data (1.2 TB): $36/month**
- **Delayed access archive storage (1.2 TB): $8/month**

RAND

# *Where Can You Get AutopsyCluster?*

- **We still have to clean up the code and document it for broader use**

- **It will be posted at**
    - **https://github.com/orgs/RANDCorporation/ AutopsyCluster**

RAND

# *Outline*

- **Objectives and vision**

- **Architecture**

- **Preliminary test results**

- **Lessons Learned**

- **How to use DIGIFORC2**
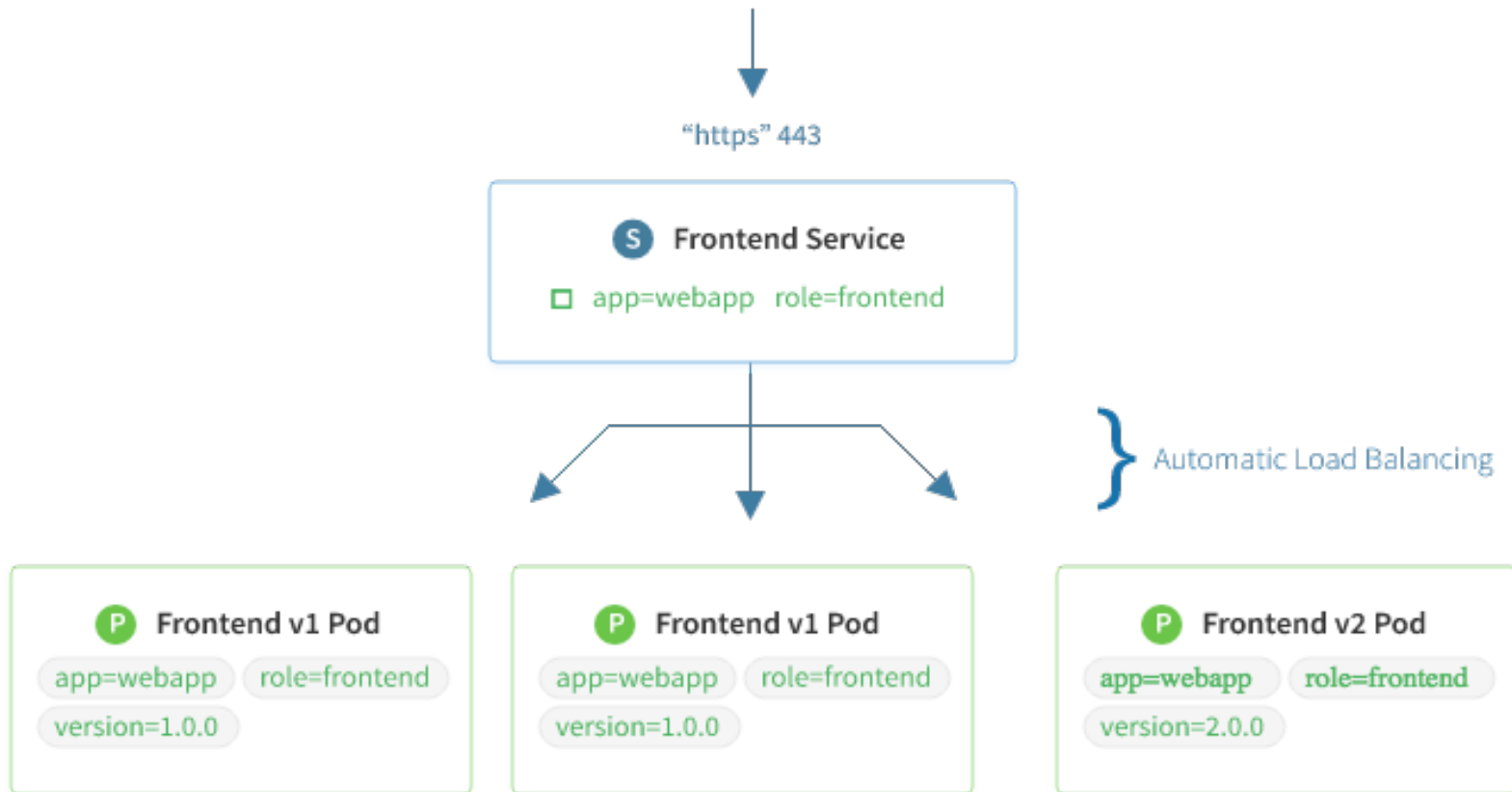
- **Beta testing**

RAND

# *We are Looking for Law Enforcement (LE) Partners as Beta Testers*

- **RAND will conduct testing, training, and evaluation with local LE**

- **Objectives of beta testing are to:**
  - **Identify performance bottlenecks found during evaluation**
  - **Provide feedback on the user interface**
  - **Simplify system configuration in response to LE feedback**

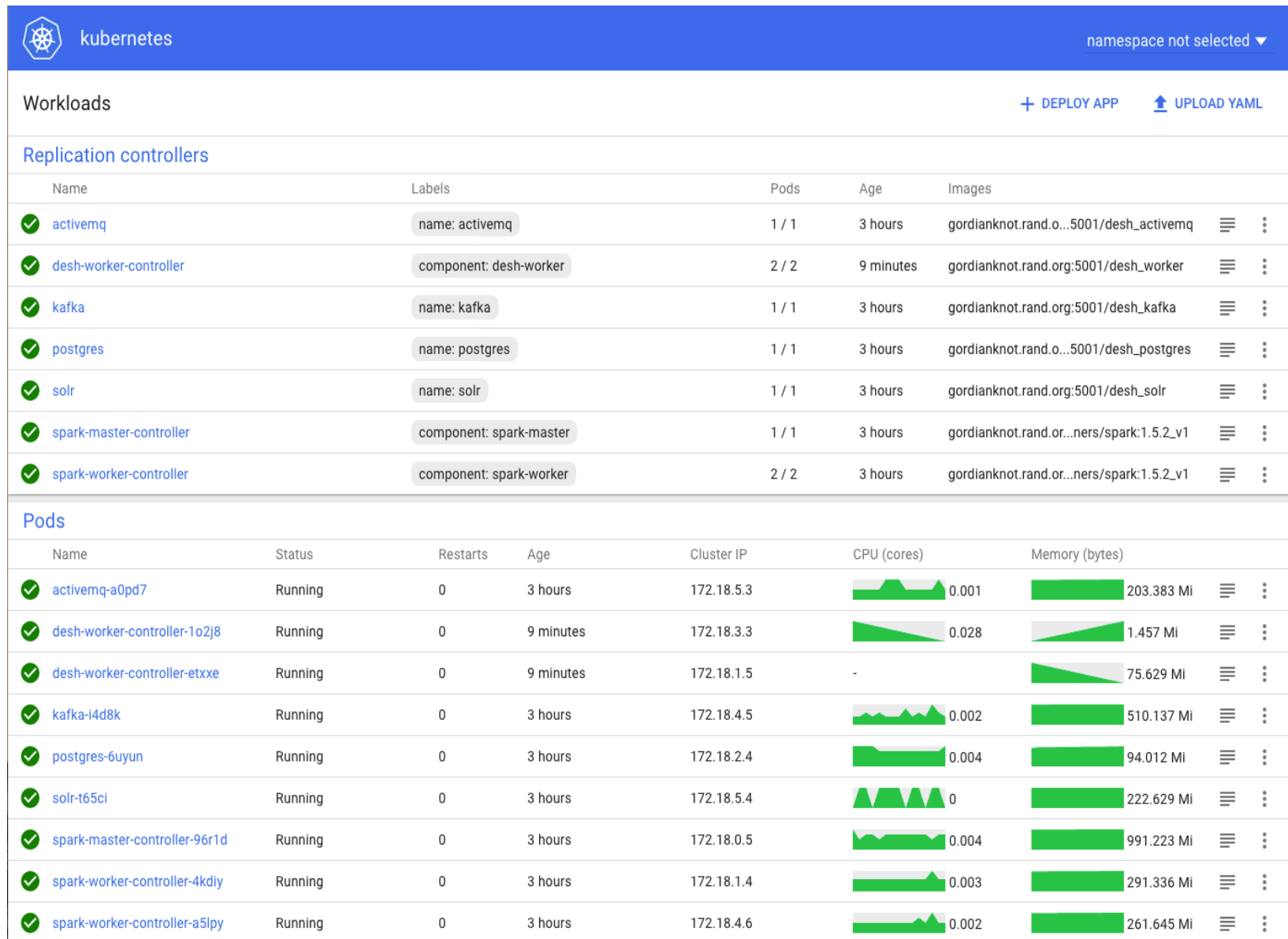- **We plan to use AWS for testing, but are open to other cloud candidates preferred by LE organizations**

RAND

# *Back Ups*

# *Kubernetes Can Provide Load Balancing*
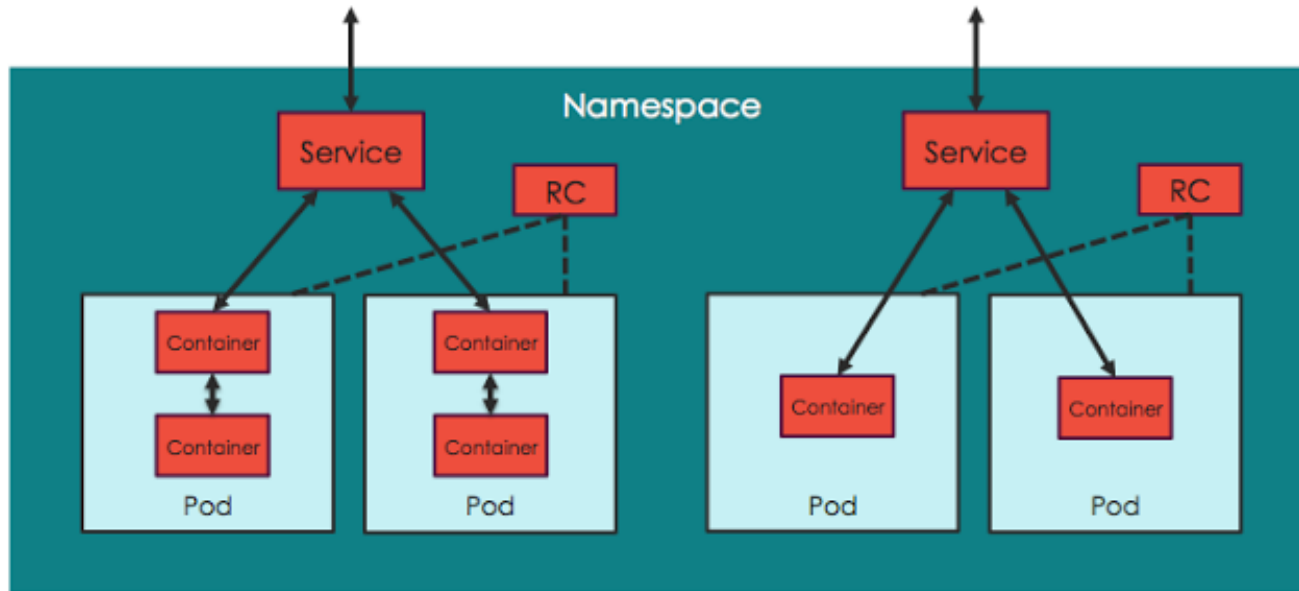
RAND

# *Overview of Project Tasks*

1. **Develop an appropriate cluster processing architecture**

2. **Integrate Autopsy with the cluster processor**

3. **Chain of custody analysis**

4. **Beta testing with law enforcement partners**

5. **Post DIGIFORC2 (Autopsy streaming branch) on Github**

RAND

# Kubernetes DIGIFORC2 Dashboard

# *Kubernetes*



- **Kubernetes is a open source platform for automating scaling and operations of containerized applications on clusters**

- **It enables applications to be scaled "on the fly"**

RAND