**Carnegie Mellon University**
Software Engineering Institute

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

# Automated Cluster Testing and Optimization

Brad Powell
Sr. Security Engineer | CERT

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA  15213

# Document Markings

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

2

Automated Cluster Testing and Optimization

# Introduction

- How to setup an automated testing framework to get benchmarks and results that will help determine tuning parameters and improve the performance of your Spark cluster

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

3

# Development and Test Environment (DTE)

- Support the architecture, design, and test processes of the lifecycle

- Provide a baseline of technologies for prototyping and testing capabilities supporting cybersecurity use cases

- Manage a shared and collaborative environment

- Evaluate relevant technology and conduct demonstrations as appropriate to inform engineering efforts and lessen risk

- Prototype data analysis techniques using the variety of available data types and tools

- Deliver Trend Reports to capture changes in the industry/community for relevant technology spaces

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

4

Automated Cluster Testing and Optimization

# Automated Testing Tools

**Carnegie Mellon University**
Software Engineering Institute

Automated Cluster Testing and Optimization
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

5

**HiBench (Intel)** - Measure speed, throughput, and system resource utilization

- Micro benchmark workloads:
  - Sort, WordCount, TeraSort, Sleep, Enhanced DFSIO

- SQL workloads:
  - Scan, Join, Aggregate

- Machine Learning workloads:
  - Bayesian Classification, K-means clustering, Logistic Regression, Alternating Least Squares, Gradient Boosting Trees, Linear Regression, Latent Dirichlet Allocation, Principal Components Analysis, Random Forest, Support Vector Machine, Singular Value Decomposition

- Websearch benchmark workloads:
  - PageRank, Nutch indexing

- Graph benchmark workloads:
  - NWeight

- Streaming workloads:
  - Identity, Repartition, Stateful Wordcount, Fixwindow

Supported releases:
Hadoop: Apache Hadoop 2.x, CDH5, HDP
Spark: Spark 1.6.x, Spark 2.0.x, Spark 2.1.x, Spark 2.2.x
Flink: 1.0.3
Storm: 1.0.1
Gearpump: 0.8.1
Kafka: 0.8.2.2

Sample Output:

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | Type | Date | Time | Input_data_size | Duration(s) | Throughput(bytes/s) | Throughput/node | |
| 2 | ScalaSparkWordcount | 7/19/2018 | 15:45:00 | 36790 | 29.818 | 1233 | 30 | |
| 3 | ScalaSparkWordcount | 7/19/2018 | 16:02:00 | 36790 | 28.841 | 1275 | 31 | |
| 4 | ScalaSparkWordcount | 7/19/2018 | 16:16:58 | 36790 | 27.225 | 1351 | 34 | |
| 5 | ScalaSparkTerasort | 7/19/2018 | 16:27:26 | 3200000 | 36.693 | 87210 | 2180 | |
| 6 | HadoopWordcount | 7/19/2018 | 17:08:34 | 36790 | 28.971 | 1269 | 31 | |
| 7 | | | | | | | | |

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

6

## SparkBench (IBM) – Benchmarking and simulating Spark jobs

- Spark-Submit-Config:
  - SparkBench converts config files into spark-submit scripts
  - Allows multiple spark-submits in series or parallel

- Workloads:
  - Standalone Spark jobs with input/output
  - Data Generators: Graph, Kmeans, Linear Regression
  - Kmeans, Logistic Regression, Sleep, SparkPi, SQL

- Workload Suites:
  - Collections of one or more workloads
  - Control benchmark output and parallelism

- Custom Workloads:
  - Use Scala and SBT to build onto SparkBench
  - Test custom Spark libraries by including JAR

```
spark-bench = {
  spark-submit-config = [{
    workload-suites = [
      {
        descr = "One run of SparkPi and that's it!"
        benchmark-output = "console"
        workloads = [
          {
            name = "sparkpi"
            slices = 10
          }
        ]
      }
    ]
  }]
}
```

Sample Output:



**Carnegie Mellon University**
Software Engineering Institute

Automated Cluster Testing and Optimization
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

7

Automated Testing Tools

## SparkBench Config

```
1 ▼ spark-bench = {
2      spark-submit-parallel = false
3 ▼    spark-submit-config = [{
4        spark-home = "/opt/spark"
5 ▼      spark-args = {
6          master = "yarn"
7          driver-class-path = "/nfs/home/bmpowell/spark-bench_2.3.0_0.4.0-RELEASE/lib/*"
8          driver-memory = 128g
9          driver-cores = 8
10       }
11 ▼     conf = [
12 ▼       {
13           "spark.dynamicAllocation.executorIdleTimeout" = "120s"
14           "spark.executor.cores" = "5"
15           "spark.executor.memory" = "80g"
16         },
17 ▼       {
18           "spark.dynamicAllocation.executorIdleTimeout" = "120s"
19           "spark.executor.cores" = "4"
20           "spark.executor.memory" = "53g"
21         }
22       ]
```

https://codait.github.io/spark-bench/

## Workload Definition

```
23      suites-parallel = false
24      workload-suites = [
25        {
26          descr = "Packed Count with filter on appLabel"
27          benchmark-output = "hdfs:///tmp/mothra-test/mothra-appLabel-results.csv"
28          parallel = false
29          repeat = 2
30          save-mode = "append"
31          workloads = [
32            {
33              name = "custom"
34              class = "com.example.MothraFilter"
35              output = "hdfs:///tmp/mothra-test/mothra-results.csv"
36              cache = "no"
37              inputdata = "/data/mothra-ipfix/ixia-packed-applabel-test/"
38            }
39          ]
40        }
41        {
42          descr = "Unpacked Count with filter on appLabel"
43          benchmark-output = "hdfs:///tmp/mothra-test/mothra-appLabel-results.csv"
44          parallel = false
45          repeat = 1
46          save-mode = "append"
47          workloads = [
48            {
49              name = "custom"
50              class = "com.example.MothraLoad"
51              output = "hdfs:///tmp/mothra-test/mothra-results.csv"
52              cache = "no"
53              inputdata = "/data/mothra-ipfix/ixia/yaf-ixia*2018112*"
54            }
55          ]
56        }
57      ]
58    }]
59  }
```

Automated Cluster Testing and Optimization

# Mothra Refresher

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

9

## Mothra Architecture

- Facilitate bulk storage and analysis of cybersecurity data with high levels of flexibility, performance, and interoperability

- Reduce the engineering effort involved in developing, transitioning, and operationalizing new analytics

- Serve all major constituencies within the network security community, including data scientists, first-tier incident responders, system admins, and hobbyists



**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

10

## SiLK vs. Mothra Scalability

- Mothra enables more complex analyses at a scale beyond the capability of SiLK's single-node architecture

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

11

Automated Cluster Testing and Optimization

# Test Plan

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**12**

Test Plan

The goal of our testing was to identify the performance and benchmarks for the DTE cluster in the following areas:

1. Cluster Operations using pre-built, Micro and Machine Learning Workloads.

2. Mothra Dataframe Creation and Spark Query Performance.

3. Mothra Ingest Process Performance running Collector and Packer processes on a 16 core physical edge node.

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

13

## Test Environment Details

- Number of Nodes: 40+10 Virtual nodes for NameNode, YARN Resource Manager, Zookeeper, and Edge Nodes
- RAM: 256GB, Disks: 4 Disks - 600 GB, CPU: 2x8 cores
- Network: Intel Corporation 82599ES 10-Gigabit dual port
- HDP Version: HDP 2.6.4 YARN
- Spark Version: Spark 2.2.1

| Test parameters | Values |
|---|---|
| spark.submit.deployMode | client |
| spark.shuffle.service.enabled | true |
| spark.scheduler.mode | FIFO |
| spark.master | yarn |
| spark.executor.memory | 4g |
| spark.dynamicAllocation.minExecutors | 4 |
| spark.dynamicAllocation.initialExecutors | 4 |
| spark.dynamicAllocation.enabled | true |
| spark.driver.port | 36562 |
| spark.driver.memory | 8g |

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

14

# HiBench Generated Data

- Four groups: Large, Huge, Gigantic, and BigData

| | Large | Huge | Gigantic | BigData |
|---|---|---|---|---|
| **Naïve Bayes** | hibench.bayes.large.pages 100000<br>hibench.bayes.large.classes 100<br>hibench.bayes.large.ngrams 2 | hibench.bayes.huge.pages 500000<br>hibench.bayes.huge.classes 100<br>hibench.bayes.huge.ngrams 2 | hibench.bayes.gigantic.pages 1000000<br>hibench.bayes.gigantic.classes 100<br>hibench.bayes.gigantic.ngrams 2 | hibench.bayes.bigdata.pages 20000000<br>hibench.bayes.bigdata.classes 20000<br>hibench.bayes.bigdata.ngrams 2 |
| **Linear Regression** | hibench.linear.large.examples 200000<br>hibench.linear.large.features 30000 | hibench.linear.huge.examples 300000<br>hibench.linear.huge.features 50000 | hibench.linear.gigantic.examples 500000<br>hibench.linear.gigantic.features 80000 | hibench.linear.bigdata.examples 1000000<br>hibench.linear.bigdata.features 100000 |
| **Random Forest** | hibench.rf.large.examples 1000<br>hibench.rf.large.features 1000 | hibench.rf.huge.examples 10000<br>hibench.rf.huge.features 200000 | hibench.rf.gigantic.examples 10000<br>hibench.rf.gigantic.features 300000 | hibench.rf.bigdata.examples 20000<br>hibench.rf.bigdata.features 220000 |
| **K-means** | hibench.kmeans.large.num_of_clusters 5<br>hibench.kmeans.large.dimensions 20<br>hibench.kmeans.large.num_of_samples 20000000<br>hibench.kmeans.large.samples_per_inputfile 4000000<br>hibench.kmeans.large.max_iteration 5<br>hibench.kmeans.large.k 10<br>hibench.kmeans.large.convergedist 0.5 | hibench.kmeans.huge.num_of_clusters 5<br>hibench.kmeans.huge.dimensions 20<br>hibench.kmeans.huge.num_of_samples 100000000<br>hibench.kmeans.huge.samples_per_inputfile 20000000<br>hibench.kmeans.huge.max_iteration 5<br>hibench.kmeans.huge.k 10<br>hibench.kmeans.huge.convergedist 0.5 | hibench.kmeans.gigantic.num_of_clusters 5<br>hibench.kmeans.gigantic.dimensions 20<br>hibench.kmeans.gigantic.num_of_samples 200000000<br>hibench.kmeans.gigantic.samples_per_inputfile 40000000<br>hibench.kmeans.gigantic.max_iteration 5<br>hibench.kmeans.gigantic.k 10<br>hibench.kmeans.gigantic.convergedist 0.5 | hibench.kmeans.bigdata.num_of_clusters 5<br>hibench.kmeans.bigdata.dimensions 20<br>hibench.kmeans.bigdata.num_of_samples 24000000000<br>hibench.kmeans.bigdata.samples_per_inputfile 40000000<br>hibench.kmeans.bigdata.max_iteration 10<br>hibench.kmeans.bigdata.k 10<br>hibench.kmeans.bigdata.convergedist 0.5 |

# Ixia Simulated IPFIX datasets

| Filename(s) | Size | Record Count | Bytes per Flow | Description |
|---|---|---|---|---|
| /data/mothra-ipfix/live/ipfix-live-s1-20180820000030-00268.yaf | 1.66 MB | 8,746 | 189.80 | 1 hour of live ipfix data from DTE YAF sensor on 8/20 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb0-20180819000021-00110.yaf | 3.17 GB | 21,997,654 | 144.11 | 1 hour of Ixia generated ipfix data for 1 YAF sensor on 8/19 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-20180819000021-00110.yaf | 22.19 GB | 154,038,796 | 144.05 | 1 hour of Ixia generated ipfix data for 8 YAF sensors on 8/19 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-20180819*.yaf | 597.61 GB | 4,159,659,052 | 143.67 | 24 hours of Ixia generated ipfix data for 8 YAF sensors on 8/19 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-2018081*.yaf | 3.08 TB | 21,772,346,751 | 141.46 | 5.5 Days of Ixia generated ipfix data for 8 YAF sensors on 8/14-8/19 |

| Filename(s) | Size | Record Count | Bytes per Flow | Description |
|---|---|---|---|---|
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb7-20180922110219-00274.yaf | 263.78 MB | 1,088,294 | 254.15 | 5 minutes of Ixia generated ipfix data for 1 YAF sensor on 9/22 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-20180922110219-00274.yaf | 2.06 GB | 8,693,166 | 254.22 | 5 minutes of Ixia generated ipfix data for 8 YAF sensors on 9/22 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb7-2018092211* | 3.11 GB | 13,199,568 | 252.92 | 1 hour of Ixia generated ipfix data for 1 YAF sensor on 9/22 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-2018092211* | 24.89 GB | 105,639,034 | 253.00 | 1 hour of Ixia generated ipfix data for 8 YAF sensors on 9/22 |
| /data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-2018092* | 605.74 GB | 2,628,246,577 | 247.47 | 24 hours of Ixia generated ipfix data for 8 YAF sensors on 9/22 |

**Carnegie Mellon University**
Software Engineering Institute

Automated Cluster Testing and Optimization
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**16**

# Automated Custom SparkBench Testing

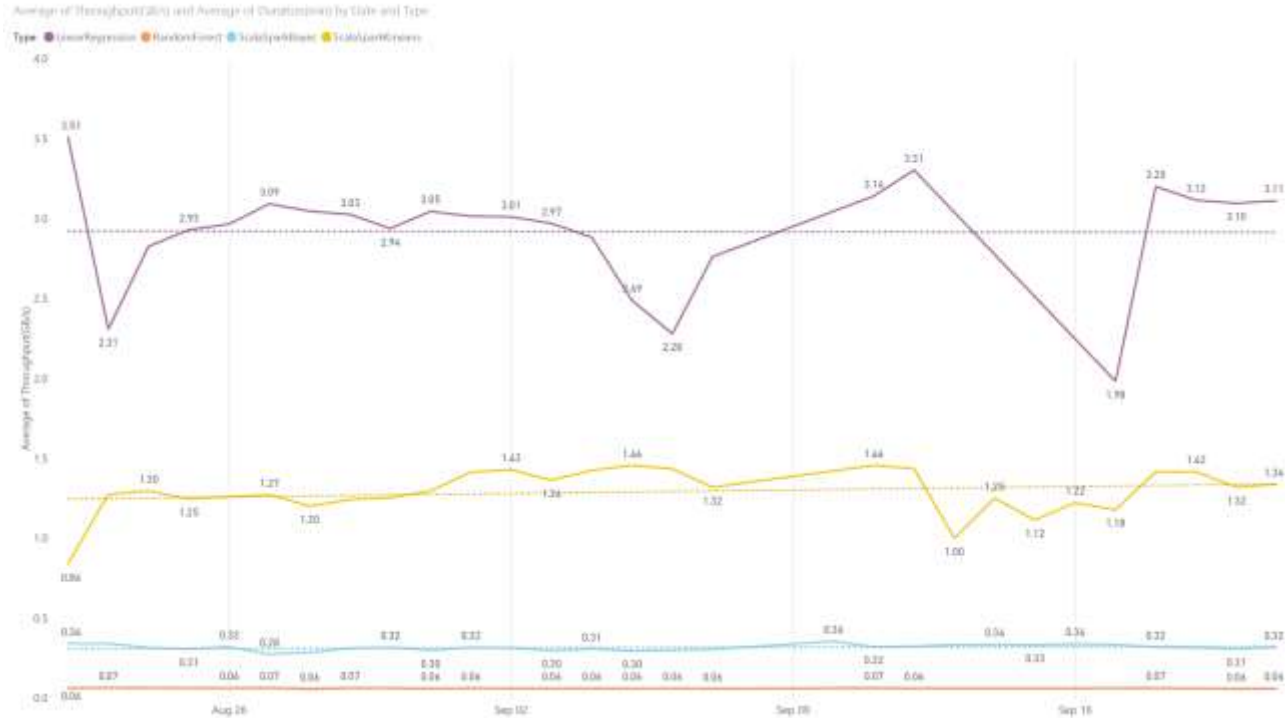| Operation | Query |
|---|---|
| Build IPFIX DataFrame (mothra), Count | val input_data_ixia = "/data/mothra-ipfix/ixia/yaf-ixia-napa_lb*-20180815*.yaf"<br>val input_df = (spark.read.<br>fields(<br>"sIP", "dIP", "sPort", "dPort", "protocol", "packets", "bytes",<br>"startTime", "endTime",<br>"dnsQName" -> "ipfix:yaf_dns/yaf_dns_qr/dnsQName",<br>"dnsQAddr" -> "ipfix:yaf_dns/yaf_dns_qr/yaf_dns_a/sourceIPv4Address" )<br>.ipfix(input_data_ixia )<br>input_df.count() |
| Simple Filter (Spark), Count | var dns_flows = input_df.filter($"dport" === 53)<br>dns_flows .count() |
| Column Selection & Display (Spark) | https_flows = input_df.filter($"dport" === 443).select(<br>"sip", "dip", "sport", "dport",<br>"protocol", "packets", "bytes", "sslCertificateHash")<br>https_flows.count() |
| Sorting (Spark) | https_flows.sort($"bytes".desc).show() |
| Aggregation (Spark) | https_flows.<br>groupBy($"dip").<br>avg("packets", "bytes").<br>sort($"avg(bytes)".desc).count() |
| SQL Query (SparkSQL) | input_df.registerTempTable("df")<br>spark.sql("""SELECT dnsQName,<br>AVG(packets) AS avg_packets,<br>SUM(packets) AS sum_packets,<br>AVG(bytes) AS avg_bytes,<br>SUM(bytes) AS sum_bytes<br>FROM df<br>WHERE dnsQName IS NOT NULL<br>GROUP BY dnsQName<br>ORDER BY sum_bytes DESC""").count() |
| Compound Query w/ Join, Filter, & Select (SparkSQL) | val bad_names = spark.read.parquet("/user/tonyc/data/sample/bad_dns_names.parquet")<br>var bad_addrs = (<br>dns_flows<br>.join(bad_names, $"dnsQName" === $"name")<br>.select("dnsQAddr")<br>.distinct).toDF("addr")<br>val pwned = input_df.join(bad_addrs, $"dIP" === $"addr").drop("addr")<br>pwned.count() |

Automated Cluster Testing and Optimization

# Results and Tuning

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public
release and unlimited distribution.]

18

## Operational ML Workloads

- Machine Learning workloads benchmark average throughput in GB/s over one month

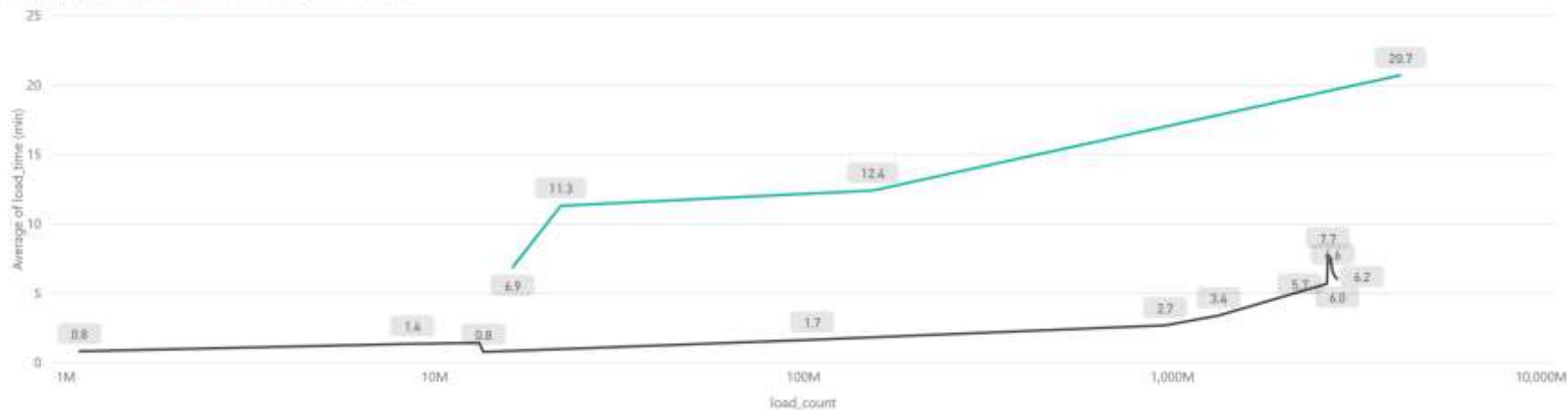**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**19**

## SparkBench Custom Mothra Workloads

- Mothra Dataframe load time in minutes by input file record count and file size below.  Graphs shows two different raw, unpartitioned file schemes.  Green is one file per hour of data and black is twelve files per hour of data.  There is a significant performance improvement when files are collected every five minutes vs. one hour.



Average of load_time (min), First load_count File Type and Average of load_time by load_count and load_count 5 Min vs 1 Hour
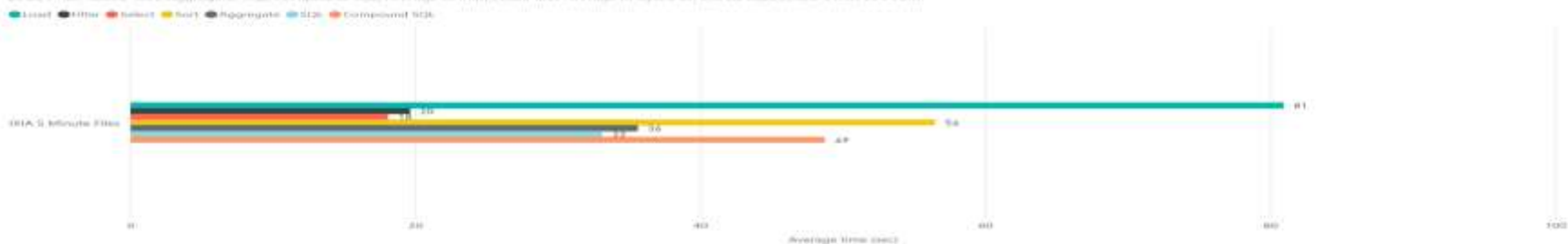
load_count 5 Min vs 1 Hour ● 00A 1 Hour Files ● 00A 5 Minute Files

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**20**
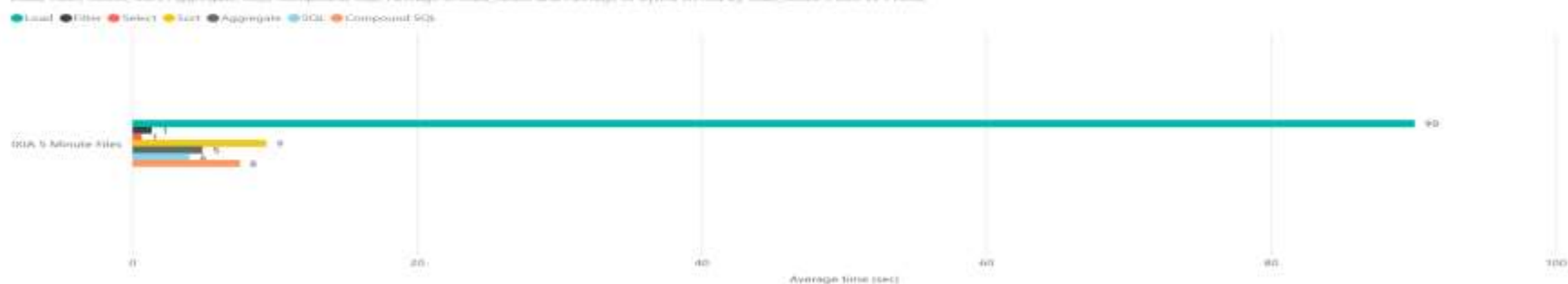
# SparkBench Custom Mothra Workloads

- Spark Submit completion times in seconds for Mothra and Spark queries.  Graph is comparing equivalent data sets with one file per hour vs twelve files per hour.  Caching in the second chart adds some overhead during load, but there is significant improvement in subsequent tasks reducing average processing time for all workloads from 293 seconds to 118 seconds a 60% improvement.



**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

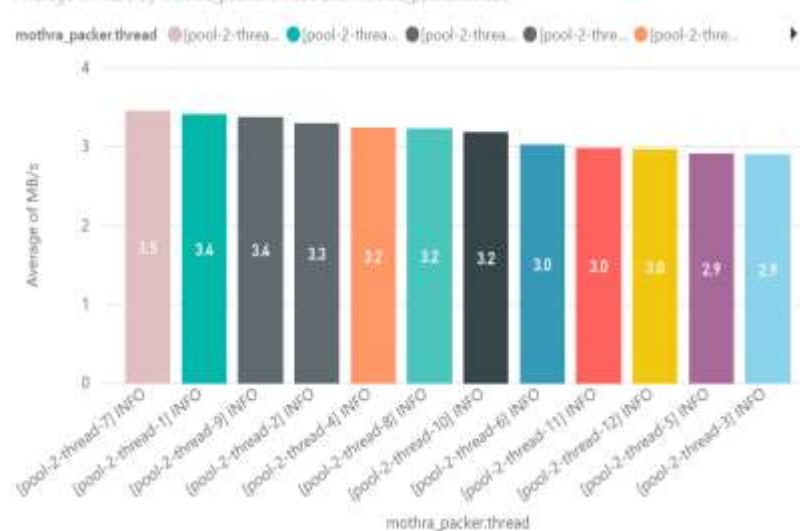[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

21

# Mothra Packer Testing

- Load times and throughput for Mothra Packer. Two sample runs of 12 max pack jobs on a 16 core physical edge node. Rwsender landed 96 files (1 hour) at once with an average of ~271 MB per file (91,346,434 records). The average throughput per process is ~3MB/s. Adding polling and flush overhead, the average of total throughput is ~27MB/s
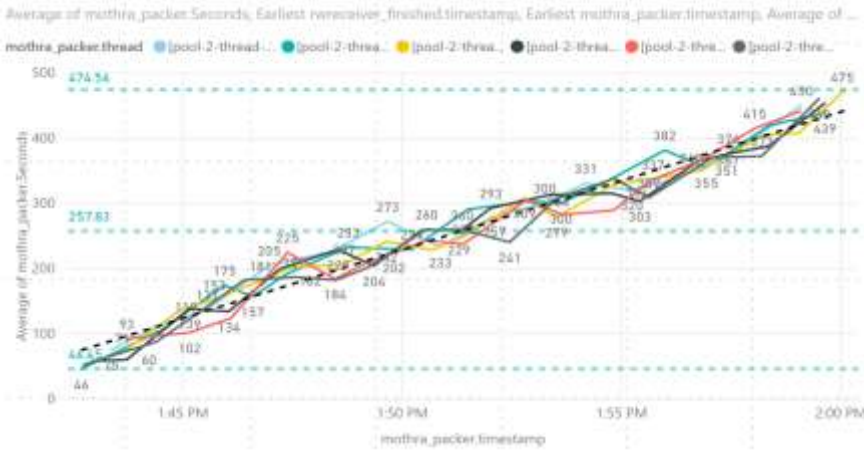
## Mothra Packer Testing

- Two sample runs 16 core edge node. Rwsender landed 96 files (1 hour) at once with an average of ~271 MB per file (91,346,434 records). The black dashed line shows the trend of completion time of each pack job.

- Test (a) shows a flat trend line which means that the jobs are keeping up with the files landing from rwsender while test (b) shows an incline trend which means that jobs are slowing over time and not able to keep up with the file ingestion. In both cases, one hour of our test data was packed in under 20 minutes, but test (a) should maintain this speed with more load, while test (b) would continue to slow as more files are landed.

(a) 12 Max Pack Jobs

(b) 6 Max Pack Jobs

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public
release and unlimited distribution.]

23
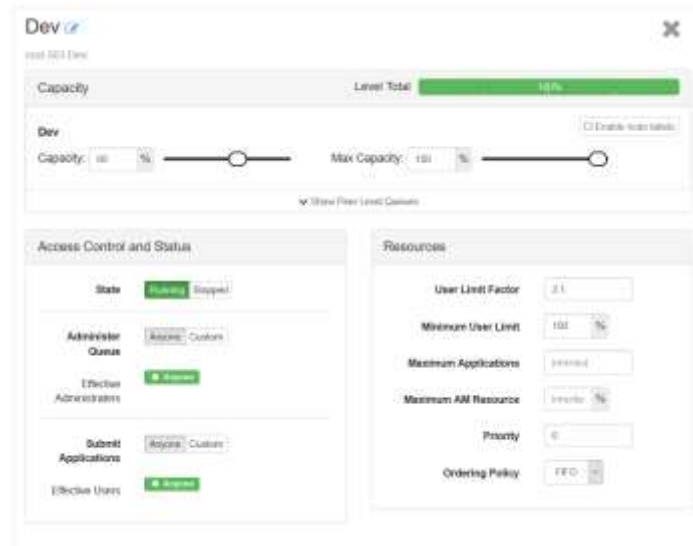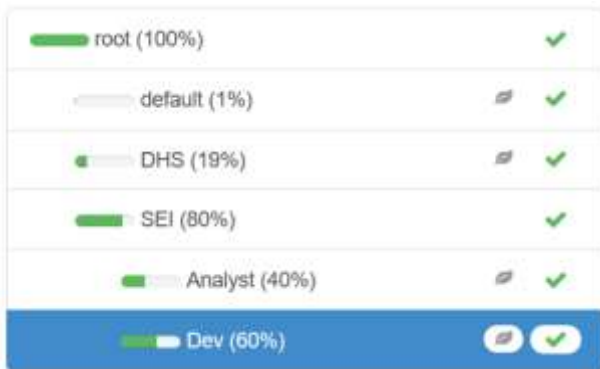
## YARN Queue Manager / Capacity Scheduler

- Certain settings needed to be changed to take full advantage of the cluster resources and utilize dynamic allocation in Spark.  Capacity and Max Capacity are not intuitive and only relate to the queue, not the whole cluster.  In order to use resources beyond the queue (80% * 60% = 48%) , User Limit Factor needs to be set above 1.

- Depending on the number of users, Minimum User Limit and Ordering Policy can be used to avoid conflicts among analysts for cluster resources.

## Spark Tuning

- Executor Cores
  - Typically no more than 5 cores can achieve full write throughput to HDFS
  - Setting cores too low (tiny executors) for large jobs on large clusters will cause garbage collection and out of memory errors
  - With executor-cores > 1, the DominantResourceCalculator must be selected for YARN

- Executor Memory
  - Calculated based on cluster size and executor-cores (Example = 6 nodes, 16 cores/node, 64gb memory/node, 5 executor-cores)

| | |
|---|---|
| yarn.nodemanager.resource.cpu-vcores * total cluster nodes = total available cores | 15 * 6 = 90 total available cores |
| total available cores / executor-cores = total available executors | 90 / 5 = 18 total available executors |
| total available executors / total cluster nodes = number of executors per node | 18 / 6 = 3 number of executors per node |
| yarn.nodemanager.resource.memory-mb / number of executors per node = memory per executor | 63 / 3 = 21 memory per executor |
| memory per executor * (1 - spark.yarn.executor.memoryOverhead) = roundDown(executor-memory) | 21 * (1 - .07) = roundDown(19.53) = 19 GB = executor-memory |

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**25**

## Spark Tuning

- Dynamic Allocation Executor Idle Timeout

  - This option controls when executors are removed once idle.

  - Losing an executor due to a timeout and starting a new one adds additional overhead to a spark job.

  - For some use cases, such as exploratory analysis in Jupyter or Zeppelin, the default timeout of 60s might be too short.

  - Finding the ideal value for this, per use case, will require an iterative process between system administrators, cluster developers, and analysts.

**Carnegie Mellon University**
Software Engineering Institute

**Automated Cluster Testing and Optimization**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution.]

**26**

Automated Cluster Testing and Optimization

# Questions?

Contact

bmpowell@cert.org