# Machine Learning for k8s Logs and Metrics

*AUTOMATING INCIDENT AND ROOT-CAUSE DETECTION*

Larry Lancaster
Founder and CTO
Zebrium

**Ze**

ZEBRIUM

# Machine data is my life

- NetApp - *Engineering Informatics*
- EMC / Data Domain - *Product Analytics*
- Glassbeam - *Chief Technology Officer*
- Nimble Storage - *Chief Data Scientist*
- Zebrium - *Founder and CTO*

Shrink-Wrap:

*1 incident 1 user*

*1 incident 1 monolith*

*1 incident 10 logfiles*

Log use for root-cause:

*index and search*

**Ze** ZEBRIUM

| 20 YEARS AGO | TODAY |
|---|---|
| **Shrink-Wrap:** | **SaaS:** |
| *1 incident 1 user* | *1 incident 100K users* |
| *1 incident 1 monolith* | *1 incident 100 services* |
| *1 incident 10 logfiles* | *1 incident 1K logstreams* |
| **Log use for root-cause:** | **Log use for root-cause:** |
| *index and search* | ***still index and search(!)*** |

# Complexity drives MTTR

**ZEBRIUM**

"The proportion of medium performers is up. Some are likely improved low performers, while others may be high performers who dropped as THEY STRUGGLED WITH INCREASED COMPLEXITY."

*Source: State of DevOps (2019)*

# Automation can't fix it

**ZEBRIUM**

"TIME TO RESTORE SERVICE PERFORMANCE STAYED THE SAME FOR BOTH ELITE AND LOW PERFORMERS WHEN COMPARED TO THE PREVIOUS YEAR."

*Source: State of DevOps (2019)*

**Ze**

ZEBRIUM

Autonomous RCA will save the world from the cost of complexity.

# What I want from a tool



**Automatically Detect Incidents**
Without Setting Up Manual Alert Rules

**Automatically Find Root Cause**
Without Manually Searching Across
GBs of Logs

# My requirements

- Arbitrary application
- Arbitrary runtime
- Arbitrary infrastructure
- Arbitrary environment

- Zero required tracing
- Zero required training
- Zero required alert rules

**Is it really too much to ask? :)**

# Logs are self-describing

## A free-text log tells a story:

```
[syslog] 2020-12-10 04:17:37 mars systemd[1]: Stopped PostgreSQL RDBMS.
...<191 lines>...
[jira] Caused by: org.postgresql.util.PSQLException: FATAL: terminating connection
due to administrator command
```

# People use logs for RCA

...so why aren't they better at helping us monitor?

# Log monitoring today

Ze
ZEBRIUM

Setup Agents /
Exporters /
Parsers

Configure Alert
Rules for
Known
Symptoms

Tune Alerts &
Build
Dashboards

Resolve
Incident

Get alerted or
otherwise
detect incident

Manually
Search Logs for
**Root Cause**

SLOW (MTTR)
FRAGILE (FORMATS CHG)
ANNOYING (ALERT FATIGUE)

HUMAN-DRIVEN

# What keeps logs "dumb"?

Logs are stuck in "index + search"

# Why are logs so hard?

**ZEBRIUM**

Formats change
Parses are ambiguous
Experts are needed to interpret
Apps are bespoke

# The junior SRE problem

**ZEBRIUM**

"Hey, I hadn't seen that happen before... then everything went sideways!"

--

Figure out when rare stuff and bad stuff are unusually correlated.

Ze: How it works

Complete relational structuring of logs

# Ze

**ZEBRIUM**

## No information included or required about:

- Known prefix formats

- Specific logtype keywords

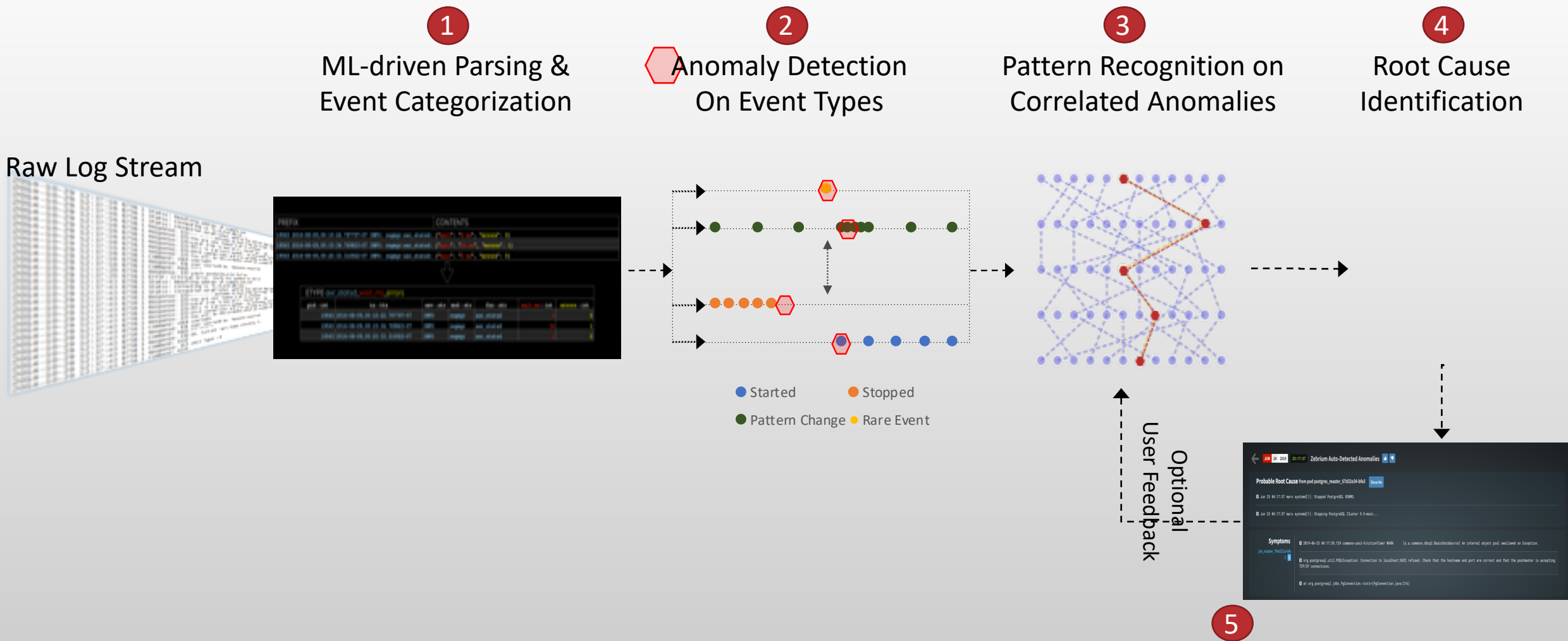- Event grammar / syntax

***We embrace free-text logs***

# Ze: How it works

Anomaly detection on relationally structured data

# Ze: How it works

## No information included or required about:

- Connectors, knowledge bases
- Specific application behaviors
- Specific semantic keywords

***Works great on bespoke app or stack***

**ZEBRIUM**

Use deep learning

Use one algorithm

Work in batch

# Use a Swiss army knife

ZEBRIUM

Structure First - Inline

Respect Pareto - multi-stage

AD/RCA gets better w / complexity!!!

GPT3/NLP requires concise RCA

**Ze**

ZEBRIUM

A picture is worth...

Ze | Overview | **Incidents** | Logs | Metrics | Alert Builder

Repeating Auto-Incidents | Show First Occ

← LIST

**Auto-Detected Incident**
Mar 18th 2020 | 12:40:10

Notes: | Did this incident detection save you time or heartache? Let us know!

Modified On: 03/18/2020 12:40:10 By: zebrium

**Possible Root Cause** | Drilldown to Incident Events →

Seen In: syslog::mars | Deployment Name: atlassian

12:40:10 | ❯ Mar 18 12:40:10 mars systemd[1]: Stopped PostgreSQL RDBMS.

12:40:10 | ❯ Mar 18 12:40:10 mars systemd[1]: Stopping PostgreSQL Cluster 9.5-main...

12:45:04 | ❯ Mar 18 12:45:04 mars systemd[1]: Starting PostgreSQL RDBMS...

**Symptoms**

All [55]

syslog::mars [1]

auth::mars [3]

bitbucket::mars [29]

jira::mars [19]

postgresexporter::mars [3]

PEAK | 03/18/2020 12:40:13.179000 pg_exporter_last_scrape_duration_seconds

PEAK | 03/18/2020 12:40:43.150000 pg_exporter_last_scrape_duration_seconds

PEAK | 03/18/2020 12:40:59.857000 confluence_one_hour_active_users_gauge

12:40:14 | ❯ time = "2020-03-18T12:40:14-07:00" level = in... "Established new...

12:40:15 | ❯ 2020-03-18 12:40:15,155 Caesium-1-3 ERROR ServiceRunner [c.a.s.cae...
8c3a'; will attempt recovery in 60 seconds

12:40:16 | ❯ org.postgresql.util.PSQLException: FATAL: terminating connection due to administrator command

12:40:16 | ❯ 2020-03-18 12:40:16,894 ERROR [hikaricp:thread-17270] org.postgresql.Driver Connection error:

12:40:19 | ❯ time = "2020-03-18T12:40:19-07:00" level = error msg = "Error opening connection to database (user = postgres%!h(MISSING)ost = /var/run/postgresql/%!s(MISSING)slmode = disable): dial unix /va
r/run/postgresql/.s.PGSQL.5432: connect: no such file or directory" source = "postgres_exporter.go:1474"

is accepting TCP/IP connections.

12:40:31 | ❯ 2020-03-18 12:40:31,880 ERROR [pool-124-thread-1] r.a.a.p.p.m.ScheduledMetricEvaluator Cannot read all projects

**Callout boxes (enlarged):**

❯ Mar 18 12:40:10 mars systemd[1]: Stopped PostgreSQL RDBMS.

❯ Mar 18 12:40:10 mars systemd[1]: Stopping PostgreSQL Cluster 9.5-main...

❯ Mar 18 12:45:04 mars systemd[1]: Starting PostgreSQL RDBMS...

PEAK | 03/18/2020 12:40:13.179000 pg_exporter_last_scrape_duration_seconds

PEAK | 03/18/2020 12:40:43.150000 pg_exporter_last_scrape_duration_seconds

PEAK | 03/18/2020 12:40:59.857000 confluence_one_hour_active_users_gauge

org.postgresql.util.PSQLException: FATAL: terminating connection due to administrator command

**Ze** Overview **Incidents** Logs Metrics Grafana    HOW-TO VIDEOS

ZEBRIUM443@ZEBRIUM.COM
TRIAL PST (-08:00)

## Nov 26th 2020

**Filter On:** First Occurrence Only    All Incident Groups    All Users    Open Incidents

⚙ **Incidents Settings**

---

**07:54:56.000000**    INCIDENT REPORT    ☑ Details...    ✕ Mute

not helpful ⭐⭐⭐⭐⭐ very helpful

**DESCRIPTION**

The root cause of the problem is that oom-killer was invoked because of a large number of allocations. The kernel's OOM killer is triggered when the system is out of memory and needs to free some memory. Since this action can kill processes, it is protected by a flag (oom_adj) which can be set or cleared by user space applications. When this flag is set, the kernel will kill processes for which there are no more than one page left in their memory cgroups (cgroups are used to control resource usage on a per-process basis). By default, Linux uses an algorithm called "RSS" (Resident Set Size) to decide whether or not to trigger the OOM killer. This algorithm calculates how much physical memory each process has reserved and compares it with its current virtual size. If there's enough memory available, then RSS will not touch any process even if they have been consuming too many resources for too long; but if there isn't enough memory available, then RSS will trigger the OOM killer and start killing processes until there's enough physical space again.

HOSTS mars    LOG TYPES kern,syslog,atlassianconfluence

---

**FIRST** 2020-11-26T07:54:56.000000  **LOGS:**kern **HOSTS:** mars  Nov 26 07:54:56 mars kernel: [2828457.044152] docker invoked oom-killer:
gfp_mask = 0x14200ca(GFP_HIGHUSER_MOVABLE), nodemask = (null), order = 0, oom_score_adj = 0

**WORST** 2020-11-26T07:55:04.763000  **LOGS:**atlassianconfluence **HOSTS:** mars  2020-11-26 07:55:04,763 WARN [Caesium-1-1]
[confluence.util.profiling.DurationThresholdWarningTimingHelperFactory] logMessage Execution time for publishing event
com.atlassian.confluence.plugins.synchrony.api.events.SynchronyStatusRestoredEvent@409ea5dd took 57123 ms (warning threshold is 5000 ms)

Zebrium - Chromium

Zebrium

portal03.zebrium.com/release-ea44_20201128094022/Logs/0005fbfc-fd00-0000-0000-01b00167fb42

Overview   Incidents   **Logs**   Metrics   Grafana   HOW-TO VIDEOS

ZEBRIUM443@ZEBRIUM.COM
TRIAL  PST (-08:00)

**Incident Report**

Search events...          Jump to time...          H  T  Y          Share

1 Watcher    Filter On:  All Logtypes   All Severities   All Labels   All eTypes   All Text   Views / Alerts

> GENERAL

> ALERT / MUTE

14 EVENTS  Seen Within A Few Seconds          Show Nearby

Mar 13 to Dec 06                                                          11-25 06:00   11-27 0

mars                                                     65,629

KERN                    Matching Events

SYSLOG              Nov 25 06:00    Nov 25 16:00    Nov 26 02:00   Nov 26 07:30   Nov 26 12:00    Nov 26 22:00

ATLAS...LUENCE      Incident Events                                                            55

09:54:56    09:55:04    Anomalies                                                          213

                    >=ERROR

3 METRICS          View All      Selected Events

node_cpu_se..._total_nice
                                    PEAK

07:54:22

node_memory_SwapFree_bytes
                                    DROP

07:54:22

process_cpu_seconds_total

Update Incident

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044152]  docker invoked oom-killer: gfp_mask = 0x14200ca(GFP_HIGHUSER_MOVABLE) , nodemask = (null) , order = 0, oom_score_adj = 0

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044153]  docker cpuset = / mems_allowed = 0

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044156]  CPU: 1 PID: 22195 Comm: docker Tainted: P O 4.15.0-122-generic #124~16.04.1-Ubuntu

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044157]  Hardware name: MSI MS-7A66/Z270I GAMING PRO CARBON AC (MS-7A66) , BIOS 1.50 04/06/2017

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044157]  Call Trace:

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044165]  ? security_capable+0x51/0x70

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044191]  RAX: 0000000000002710 RBX: 0000000000004e20 RCX: 00000000000001dc

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044192]  RDX: 000055ac17298858 RSI: 0000000000000000 RDI: 00007f62a464ae38

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044193]  RBP: 00007f62a464ae50 R08: 000000c000000180 R09: 0000000000000001

✓  ⚡ ❯ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044193]  R10: 0000000000000000 R11: 0000000000000202 R12: 0000000000000000

✓  ⚡ ❯ 2020-11-26 07:54:57.000000 syslog  mars kernel:  [2828457.044386] Out of memory: Kill process 18860 (oom_test) score 817 or sacrifice child

✓  ⚡ ❯ 2020-11-26 07:54:57.000000 syslog  mars kernel:  [2828457.044390] Killed process 18860 (oom_test) total-vm:54212848 kB, anon-rss:29072816 kB, file-rss:1260 kB, shmem-rss:0 kB

✓  ⚡ ❯ 2020-11-26 07:54:57.000000 syslog  mars kernel:  [2828458.143874] oom_reaper: reaped process 18860 (oom_test) , now anon-rss:0 kB, file-rss:0 kB, shmem-rss:0 kB

✓  ⚡ ❯ 2020-11-26 07:55:04.763000 atlassianconfluence  WARN [Caesium-1-1] [confluence.util.profiling.DurationThresholdWarningTimingHelperFactory]  logMessage Execution time for publishing event com.atlassian.confluence.plugins.synchrony.api.events.SynchronyStatusRestoredEvent@409ea5dd took 57123 ms (warning threshold is 5000 ms)
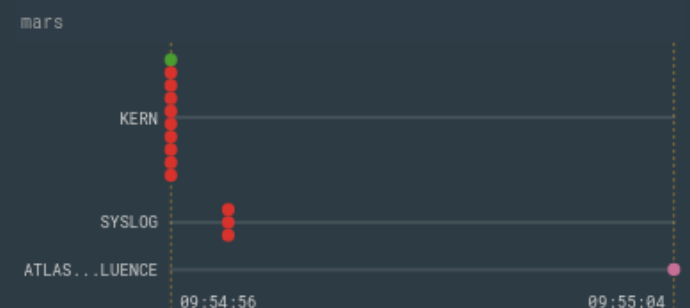
- - - End - - -

1 Active Filter

Version release-ea44_20201128094022  ·  Copyright © 2020 Zebrium, Inc. All rights reserved.

Zebrium

portal03.zebrium.com/release-ea44_2020128094022/Logs/0005fbfc-fd00-0000-0000-01b00167fb42

Overview   Incidents   **Logs**   Metrics   Grafana   HOW-TO VIDEOS

ZEBRIUM443@ZEBRIUM.COM
TRIAL   PST (-08:00)

**Incident Report**

Search events...   Jump to time...   H  T  Y   Share

> GENERAL

▼ ALERT / MUTE

1 Watcher   Filter On:   All Logtypes   All Severities   All Labels   All eTypes   All Text   Views / Alerts

For incidents of this type...

☑ Alert in Future     Mute

Zebrium Webhook   Slack

Configure Now...   Configure Now...

Incident Quality   not helpful ★★★★★ very helpful

Mar 13 to Dec 06                                            11-25 04:00  11-27 0

Matching Events

Nov 25 04:00    Nov 25 14:00    Nov 26 00:00    Nov 26 10:00    Nov 26 20:00    Nov 27 06:00

Incident Events                                                          55

Anomalies                                                                233

>=ERROR

▼ 74 EVENTS  Seen Within 3 Minutes     Hide Nearby

Selected Events

mars

ATLAS...LUENCE          2020-11-26T09:55:02.908000

ATLAS...BUCKET          2020-11-26 07:55:02,908 ERROR [Caesium-1-3] [scheduler.caesium.impl.CaesiumSchedulerService] executeClusteredJobWithRecoveryGuard
                        Unhandled exception during the attempt to execute job 'reminderTrigger'; will attempt recovery in 60 seconds

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 syslog  mars kernel: [2828457.044228] Free swap = 0 kB

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 syslog  mars kernel: [2828457.044230] [pid] uid tgid total_vm rss pgtables_bytes swapents oom_score_adj name

                                                                         50 1632 0 ( sd-pam

                                                          mask = 0x14200ca(GFP_HIGHUSER_MOVABLE) , nodemask = (null) , order = 0,

SYSLOG

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044153]  docker cpuset = / mems_allowed = 0

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044156]  CPU: 1 PID: 22195 Comm: docker Tainted: P O 4.15.0-122-generic #124~16.04.1-Ubuntu

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044157]  Hardware name: MSI MS-7A66/Z270I GAMING PRO CARBON AC (MS-7A66) , BIOS 1.50 04/06/2017

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044157]  Call Trace:

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044165]  ? security_capable+0x51/0x70

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044191]  RAX: 0000000000002710 RBX: 0000000000004e20 RCX: 00000000000001dc

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044192]  RDX: 000055ac17298858 RSI: 0000000000000000 RDI: 00007f62a464ae38

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044193]  RBP: 00007f62a464ae50 R08: 000000c000000180 R09: 0000000000000001

✓  ⚡ ⊙ 2020-11-26 07:54:56.000000 kern  mars kernel: [2828457.044193]  R10: 0000000000000000 R11: 0000000000000202 R12: 0000000000000000

Update Incident

1 Active Filter

# Join us on this journey!

```
URL: zebrium.com/how-to-try

email: larry@zebrium.com

twitter: stochastimus@twitter.com
```

Gartner COOL VENDOR 2020

DZone A DEVADA MEDIA PROPERTY
Best Log Platform for Kubernetes 2020

Forbes 2020 AI 50 MOST PROMISING AI FIRMS

Gartner Top 25 Enterprise Software Startups to Watch in 2020