

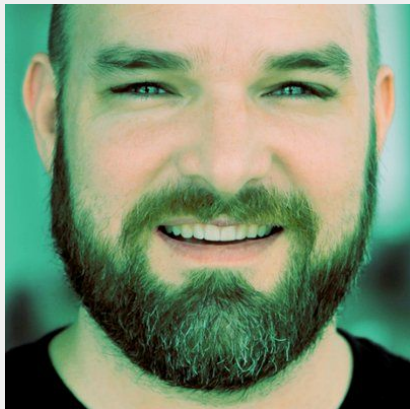


Holistic Monitoring

DevConf 2019 - Brno, Czechia

Sunday, January 27th

Who are we?



Jared Sprague
Principal Software Engineer

jsprague@redhat.com
@caramelcode



Adam Minter
Associate Project Manager

aminter@redhat.com

What we'll cover today

Problem Statement

Our Solution

Monitoring Strategy Walkthrough

How is this working for us today?

Questions!

Problem Statement

Perception & Missed Signals

Perception

Perception

The term, “Monitoring” too often coincides with a tool



- Zabbix = “Our monitoring”
 - Nagios, Solarwinds, Prometheus, etc etc
- Mis-align that tool = monitoring

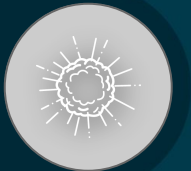
Monitoring = a wildly complex and persistent art of principals, tooling, technology, and innovation



Perception

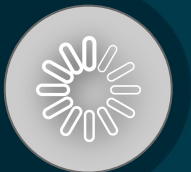
The term, “Monitoring” too often coincides with a tool

Tool Change = Really Hard



Redundant Tooling Inevitable

Technical Debt



Adoption & Usability

Missed Signals

“We didn’t catch that in our monitoring”



Noise

Aka alerts that don't matter

```
1 of 2
FRM:nagios.eng.rdu2@redhat.com
MSG:N: PROBLEM

Svc: Mysql Server

Hst: db02.db.eng.rdu2

State: CRITICAL

D/T: 01-25-2019
(Con't) 2 of 2
20:17:15

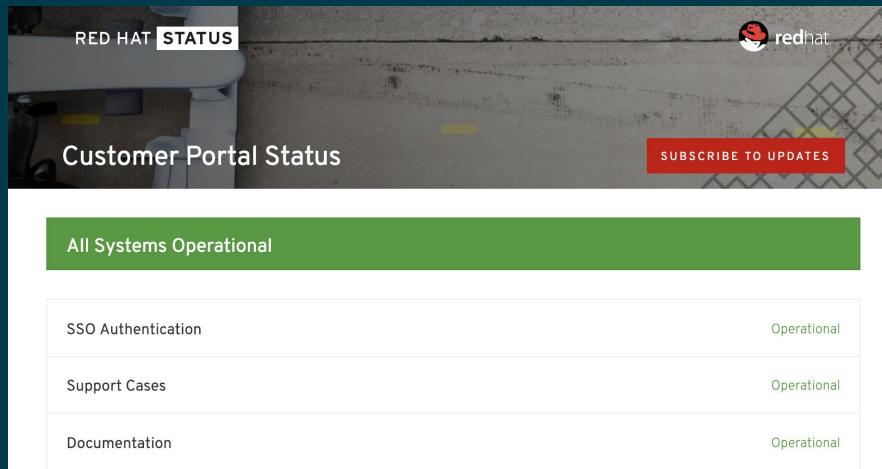
More:

SLOW_SLAVE CRITICAL: Slave IO: Yes Slave SQL: Yes Seconds Behind Master: 1000

Author:

Comment: (End)
```

Alerts



All Systems Operational	
SSO Authentication	Operational
Support Cases	Operational
Documentation	Operational

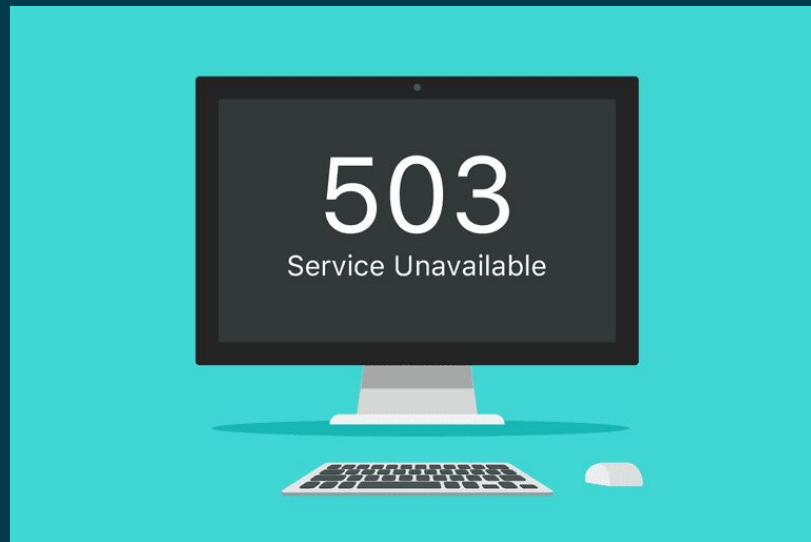
Operational State

No Signals

No alerts when it matters



Alerts



Operational State

Signal to Noise Ratio



Our Solution

Holistic Monitoring Strategy

Dashboards

- Centralized place for visual analysis of data
- Assurance that all functions can be exported and consumed here
- Provides simple digestion of the hard work done in the functions below

Logging

- A bucket for all transaction data from each function in our strategy
- Ensures we comply with data retention policies in a consistent manner
- Facilitates transparency with the data we collect from each source

Alert Orchestration • Federation of all data sources to determine if/when someone needs to act on a problem or issue

Metrics • Empowers engineers to trend the data collected from our sources and make proactive improvements

Application Performance

- Transaction traces and snapshots
- Stack specific performance metrics
- Inspection at the code level to determine performance

Availability

- Baseline determination if application is working
- Ping, API calls, or custom scripting
- Often reported through application performance indicators (KPIs)
- Examination from different geographical locations

Real User Monitoring

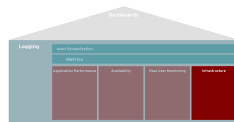
- Inspection of users interactions with a WebUI
- Sometimes Javascript that collects metrics from user's browsers
- Helps to add context to issues and scenarios afflicting a system or application

Infrastructure

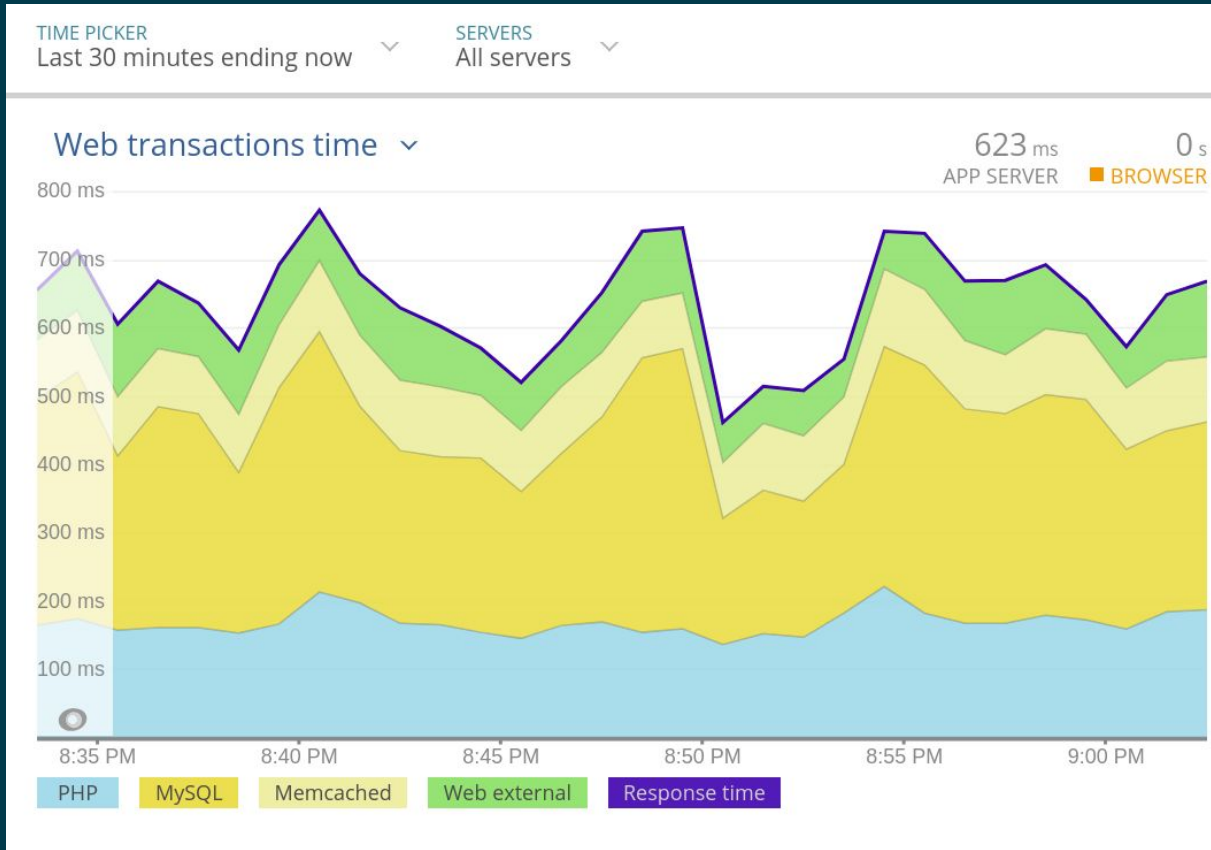
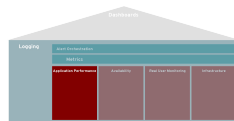
- System Resources like CPU, Memory, Disk, etc
- Network monitoring also falls into this category
- Inventory and Resource management
- The throughput of a build pipeline is one complex example

Monitoring Strategy Walkthrough

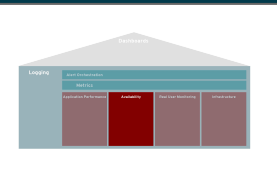
Infrastructure



APM

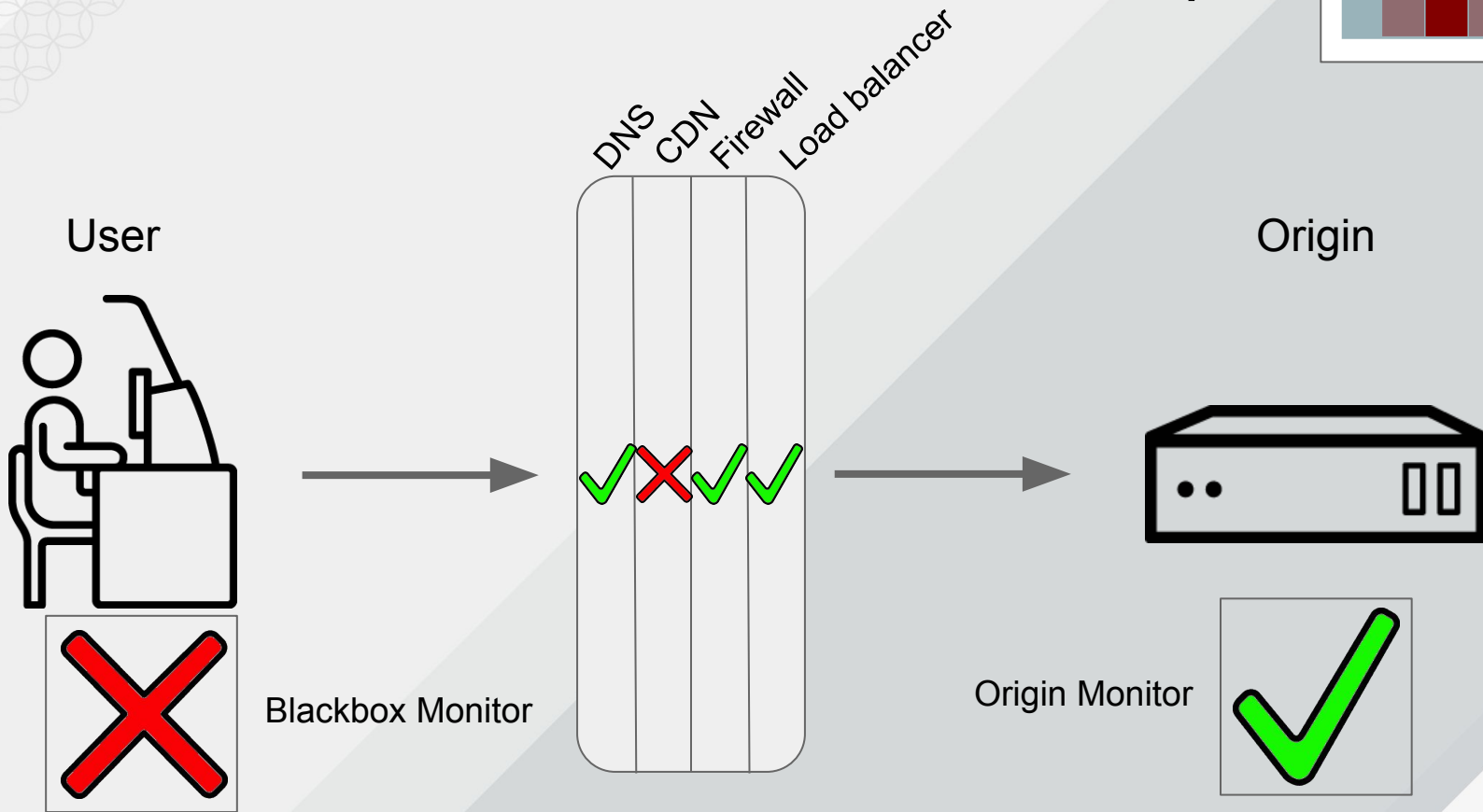


Availability

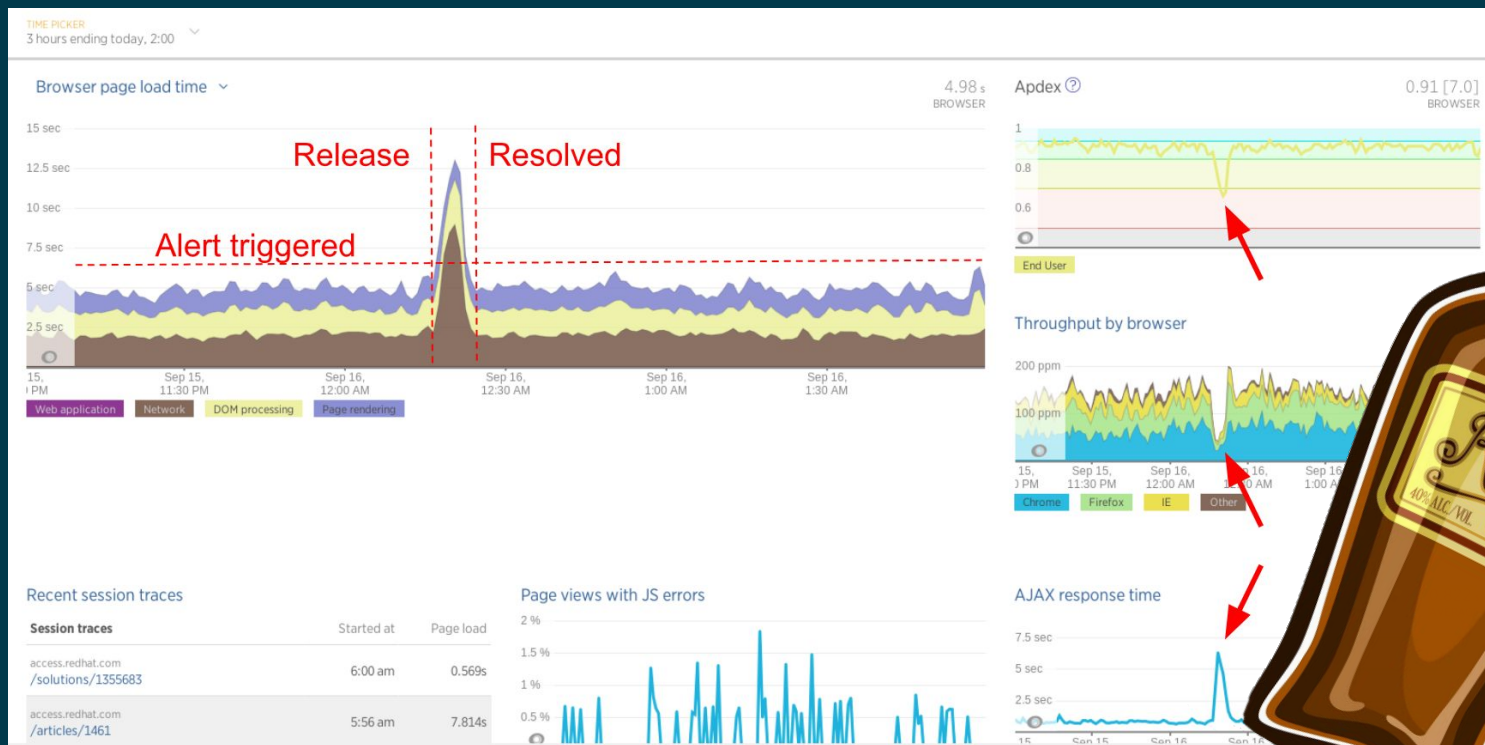
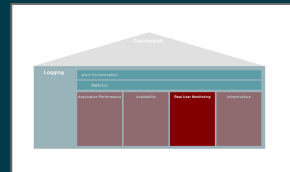


**The fallacy of origin based availability monitoring:
Customers don't use applications in the Datacenter**

User Centric Blackbox Availability



Real User Monitoring



Logging



Holistic Monitoring Strategy

Dashboards

- Centralized place for visual analysis of data
- Assurance that all functions can be exported and consumed here
- Provides simple digestion of the hard work done in the functions below

Logging

- A bucket for all transaction data from each function in our strategy
- Ensures we comply with data retention policies in a consistent manner
- Facilitates transparency with the data we collect from each source

Alert Orchestration • Federation of all data sources to determine if/when someone needs to act on a problem or issue

Metrics • Empowers engineers to trend the data collected from our sources and make proactive improvements

Application Performance

- Transaction traces and snapshots
- Stack specific performance metrics
- Inspection at the code level to determine performance

Availability

- Baseline determination if application is working
- Ping, API calls, or custom scripting
- Often reported through application performance indicators (KPIs)
- Examination from different geographical locations

Real User Monitoring

- Inspection of users interactions with a WebUI
- Sometimes Javascript that collects metrics from user's browsers
- Helps to add context to issues and scenarios afflicting a system or application

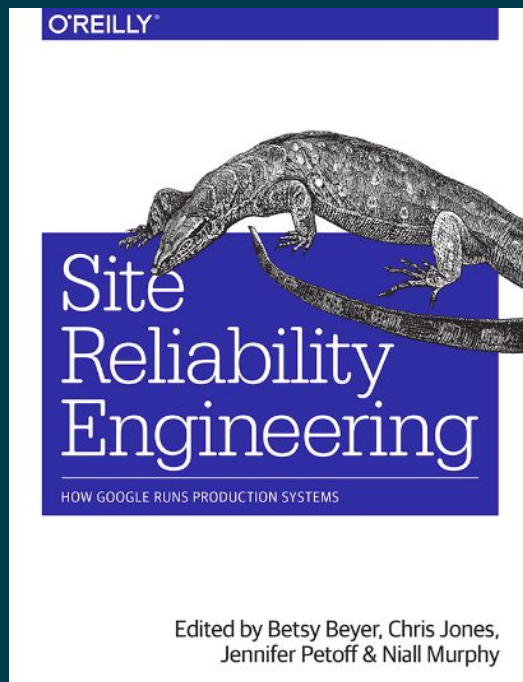
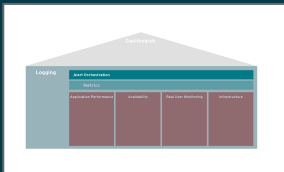
Infrastructure

- System Resources like CPU, Memory, Disk, etc
- Network monitoring also falls into this category
- Inventory and Resource management
- The throughput of a build pipeline is one complex example

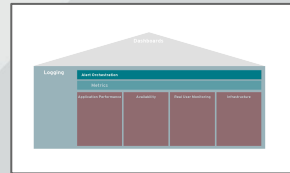
SLIs & SLOs

What are they? Why use them?

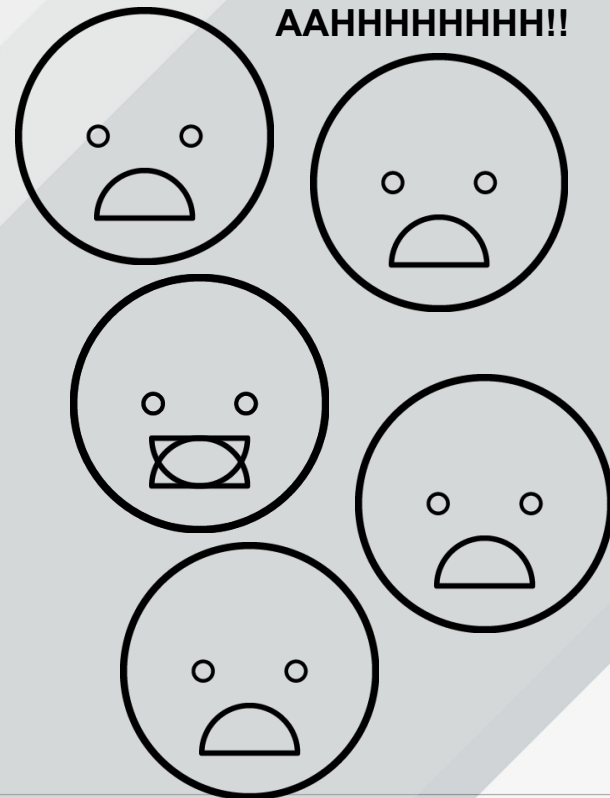
- Service-Level Indicator (SLI)
 - Ex: Error Rate.
- Service-Level Objective (SLO)
 - Ex: Error rate < 5% of 99% of requests over a 30min period
- Why use them?
 - Focuses alerts on business outcomes, and not arbitrary metrics
 - Improves signal to noise ratio
- Learn more in Google SRE book



Without Alert Orchestration



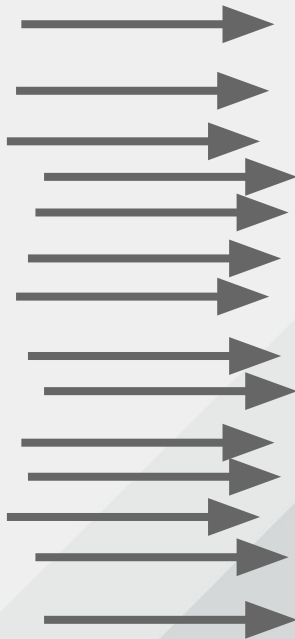
Notifications



With Alert Orchestration

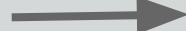


Alerts

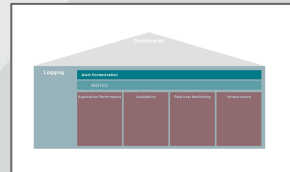


Alert
Orchestrator

Notifications



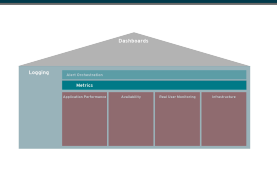
On-call



Z z



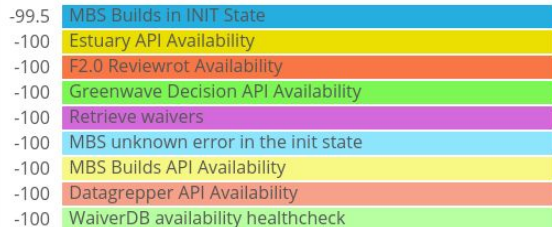
Metrics and Dashboards



Monitor Specific Uptime Past 7 Days

Since 7 days ago

Availability



7 Day Overall Uptime

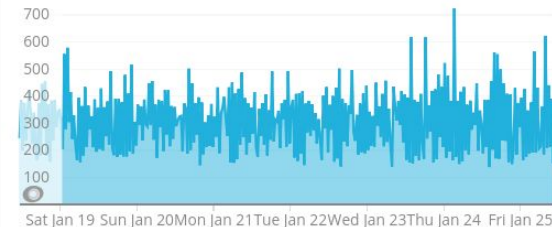
Since 7 days ago

99.94

Percentage

Greenwave Decision Latency

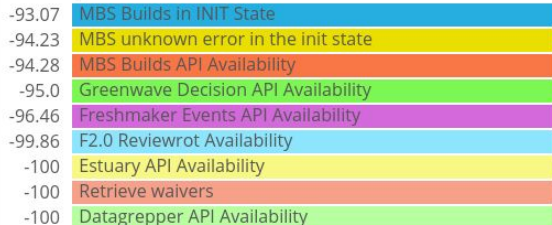
Since 7 days ago



30 days uptime by monitor

Since 30 days ago

Availability



30 Day Overall Uptime

Since 30 days ago

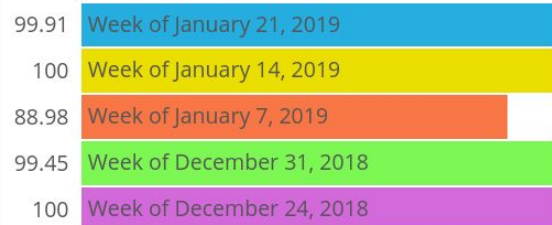
97.29

Percentage

Overall Uptime Percentage by Week

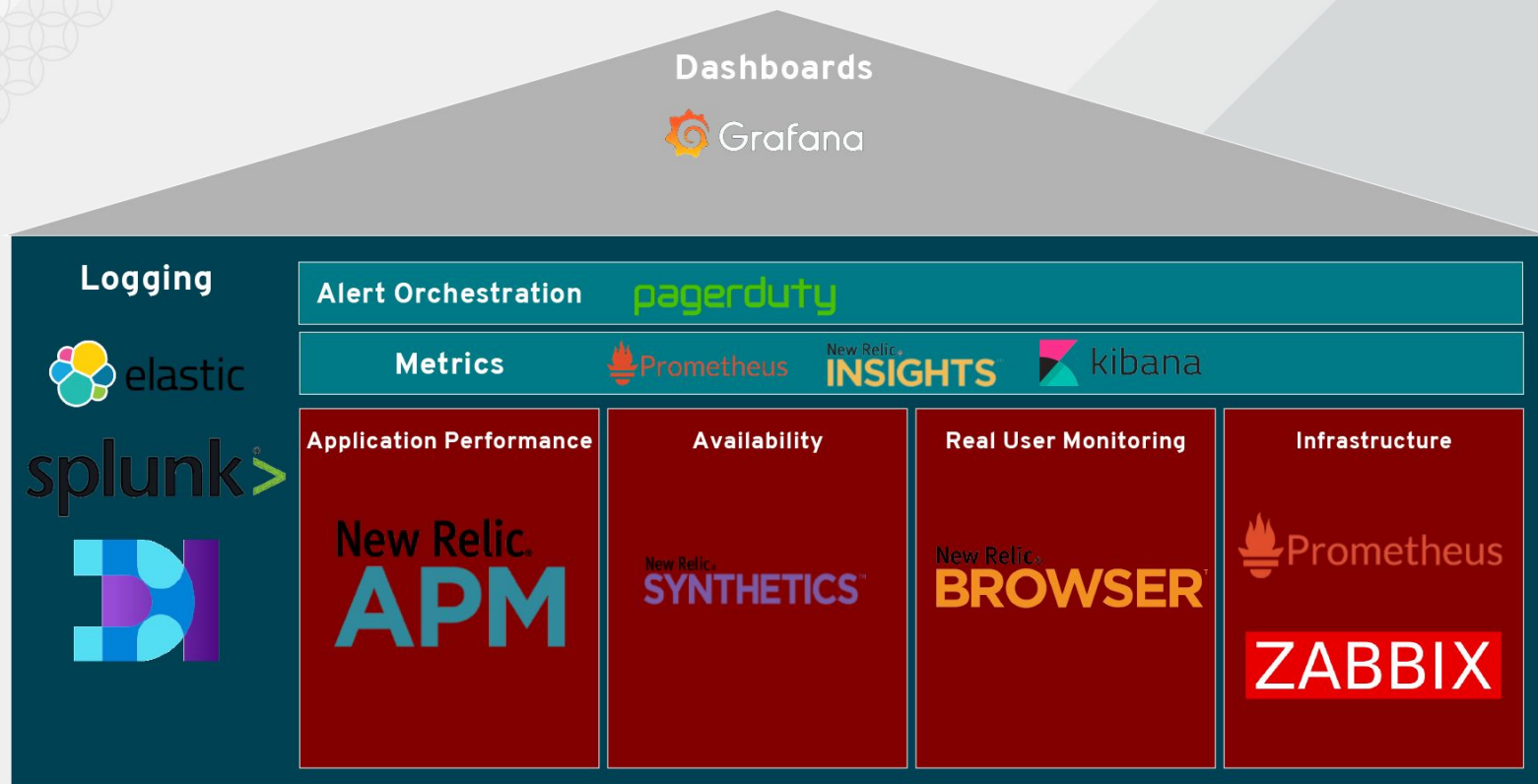
Since 1 month ago

Percentage



How is this helping us today?

Holistic Monitoring Strategy (Tooling)



Open Source Alternative (Tooling)

Dashboards



Logging



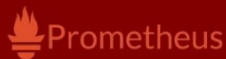
Alert Orchestration



Metrics



Application Performance



Availability



Real User Monitoring

Bucky

Infrastructure



ZABBIX



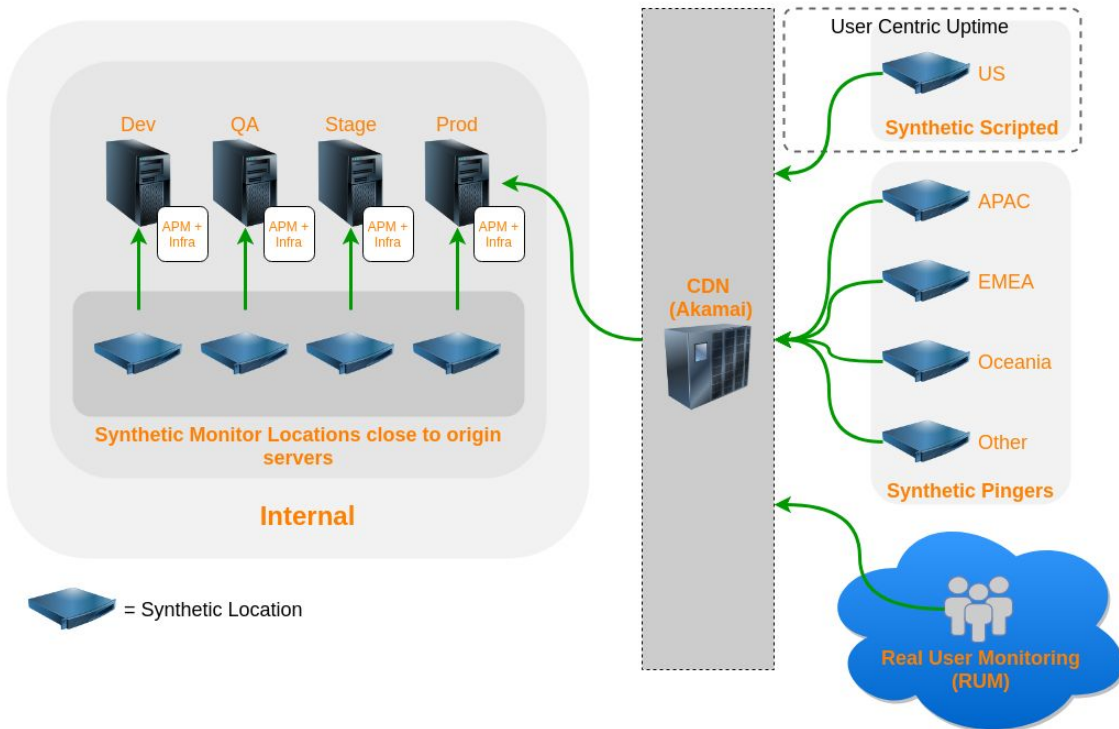
Nagios

Customer Portal Holistic Deployment

ENTERPRISE APPLICATION MONITORING APM + Synthetics + RUM + Infrastructure

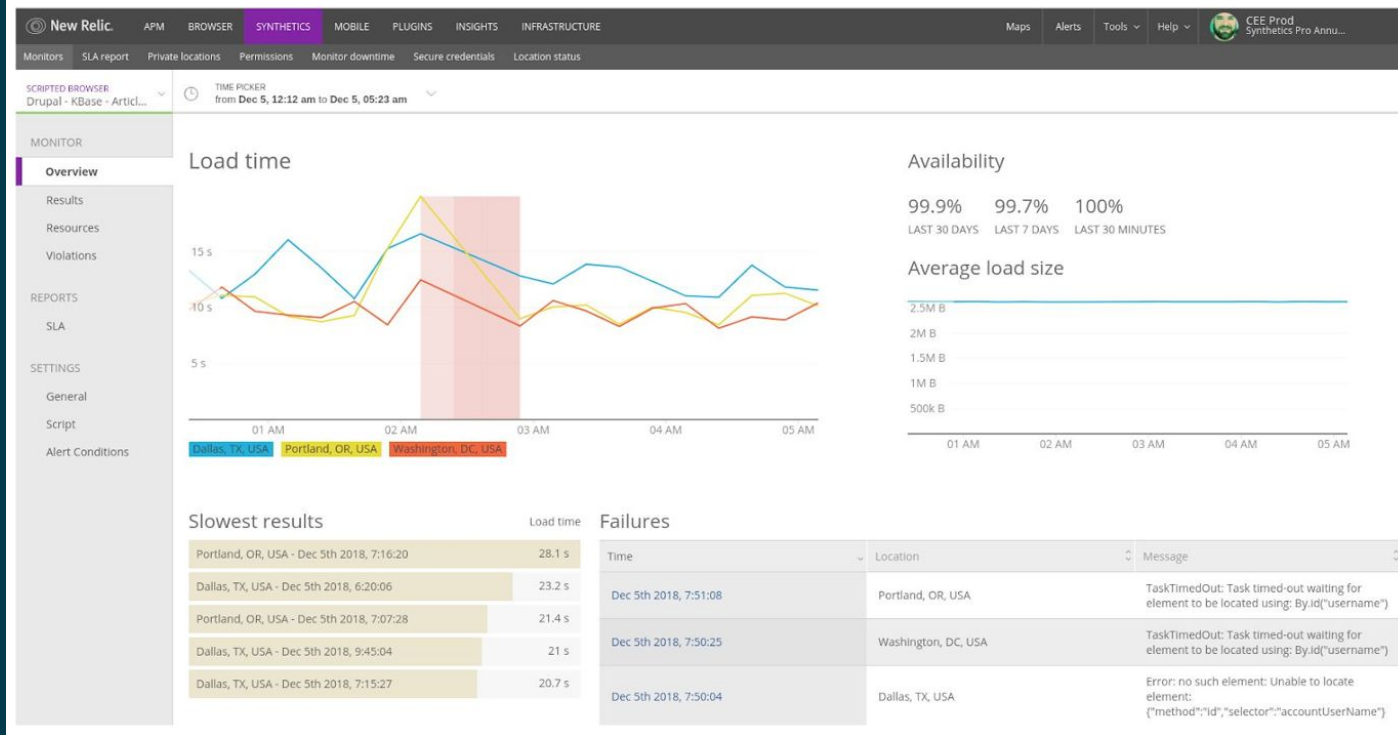
Key Features

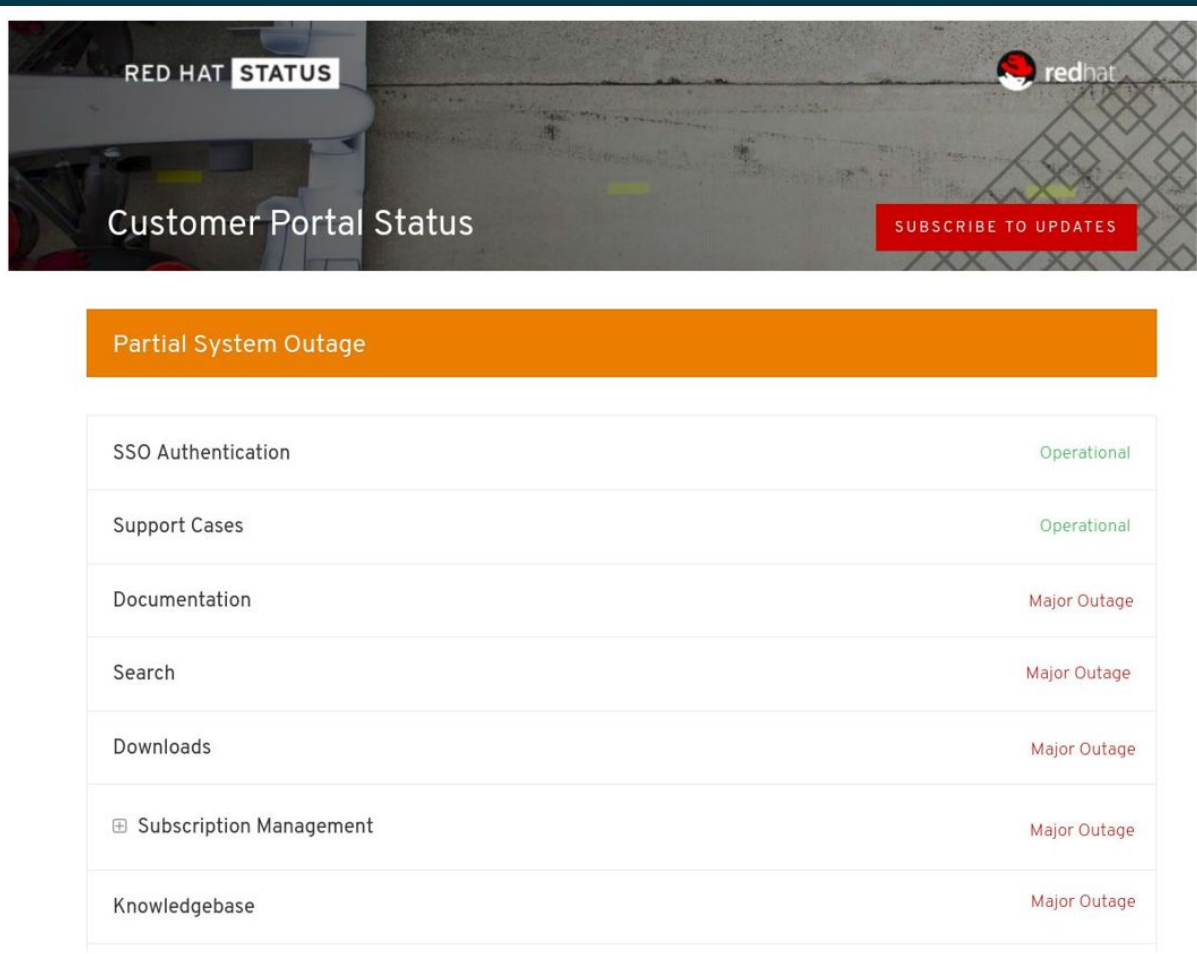
- Exceptionally accurate Uptime reports.
- Three different monitor types both internal and external provide complete picture and assist with troubleshooting problems.
- Awareness of global performance trends and isolated outages.
- Wide variety of performance metrics stored overtime from synthetic tests + Real User experience.
- Pre-prod monitoring detects problems or improvements before release.



Holistic Postmortem

At 2:19 AM EST CP On-call members were notified about failures in many New Relic synthetic availability monitors, for Search front-end, kbase, documentation, labs, container catalog and PCM.





RED HAT STATUS

redhat

Customer Portal Status

SUBSCRIBE TO UPDATES

Partial System Outage

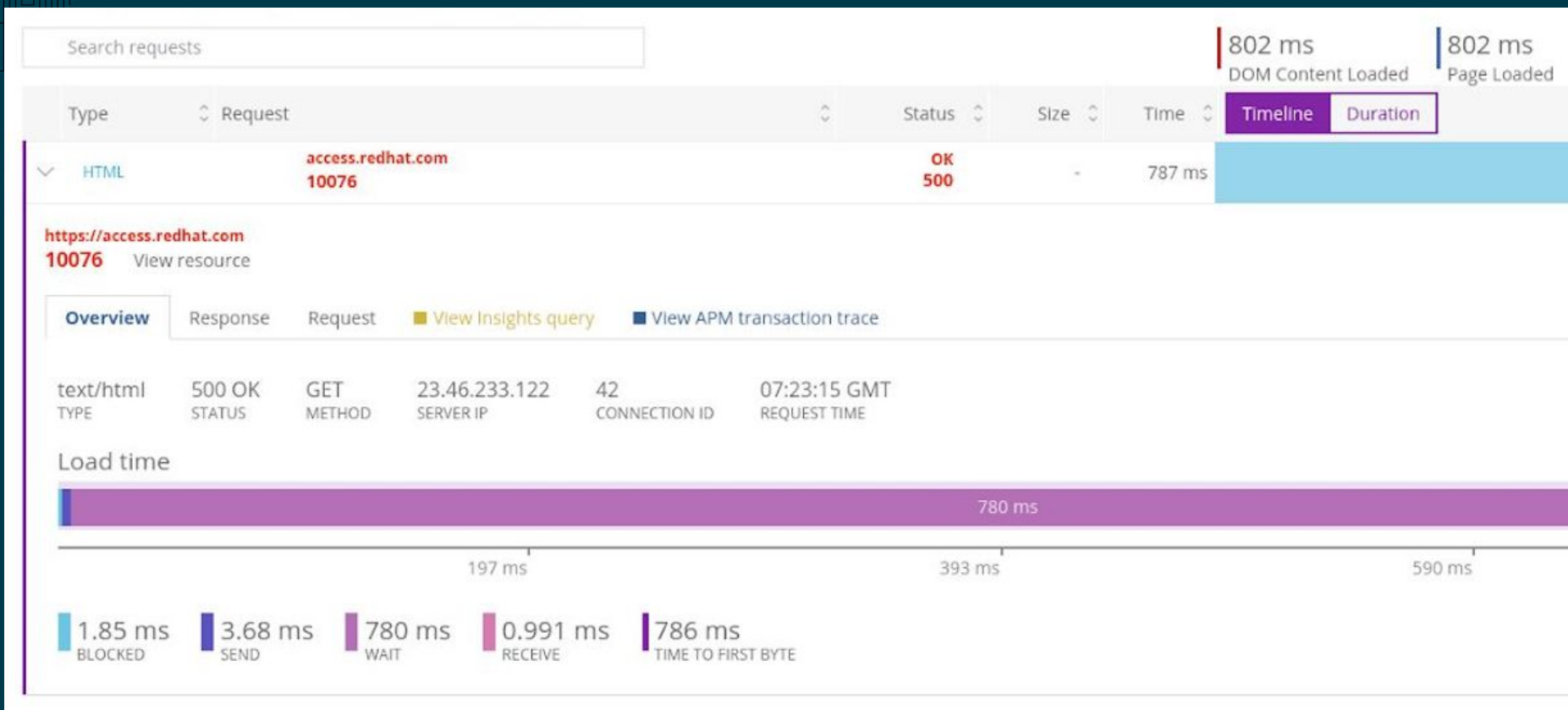
SSO Authentication	Operational
Support Cases	Operational
Documentation	Major Outage
Search	Major Outage
Downloads	Major Outage
⊕ Subscription Management	Major Outage
Knowledgebase	Major Outage

Real User Monitoring

From RUM we can see that this also affected user experience by increasing the overall page load time for end users:



Availability



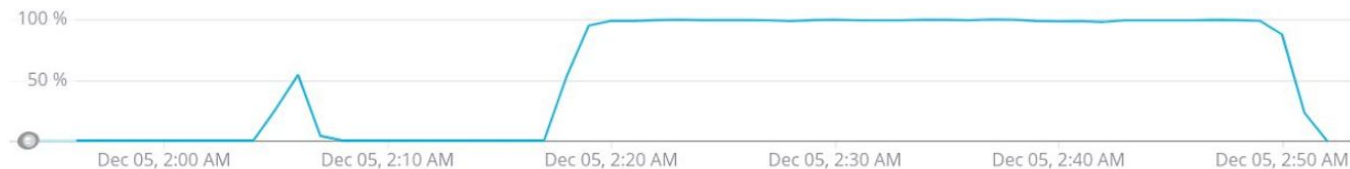
Tracing from the Synthetic check to the APM transaction trace, we see that the last call in the stack before it threw an exception was: `DatabaseConnection::__construct`

5.22%	file_scan_directory	0.545 s
2.0	file_scan_directory	0.584 s
31.0	> 6 fast method calls	0.587 s
13.0	file_scan_directory	0.626 s
3.0	variable_set	0.646 s
3.0	db_merge	0.646 s
3.0	Database::getConnection	0.646 s
3.0	Database::openConnection	0.646 s
3.0	DatabaseConnection_mysqli::__construct	0.646 s
3.0	DatabaseConnection::__construct	0.646 s

APM + Logging

Error rate ?

for all errors



Top 5 errors

by error class



● PDOException

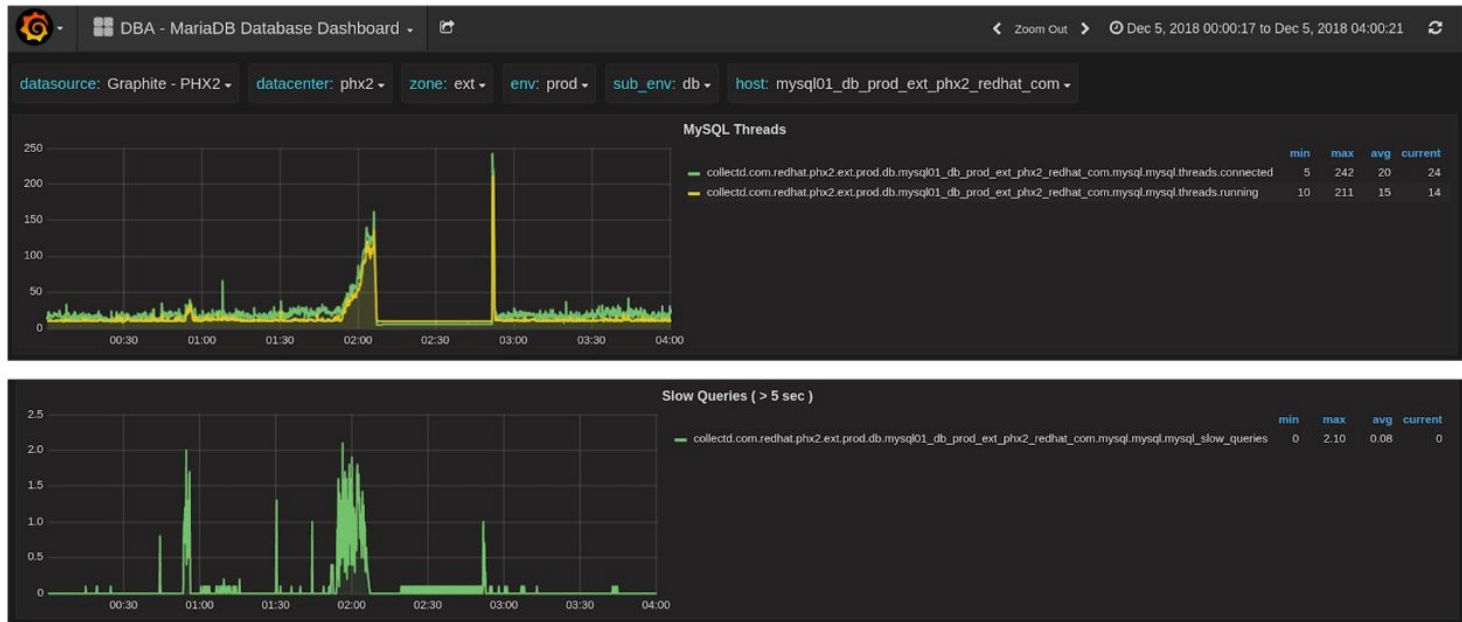
Error traces

Error frequency by error class

Error profiles

Count	Transaction name and error class	Error message	First Occurrence	Last Occurrence
9,452	/index.php PDOException	Uncaught exception 'PDOException' with message 'SQLSTATE[HY000] [2002] Connection refused' in /usr/share/drupal7/drupal/includes/database/database.inc:321	2:06 AM	2:51 AM
105	/views_page PDOException	Uncaught exception 'PDOException' with message 'SQLSTATE[HY000]: General error: 2006 MySQL server has gone away' in /usr/share/drupal7/drupal/includes/database/database.inc:2227	2:06 AM	2:19 AM

So now we know that drupal can't reach MySQL so we check the MySQL Infrastructure monitoring which is done by collectd and grafana, and we see the spike in thread count, and the host drained of memory, swap usage spike, and finally the process is killed.





Immediate Corrective Actions

Describe the immediate actions taken to stop the incident.

Restarting dead MySQL service on the mysql01.db.prod.ext.phx2.redhat.com host

Secondary Corrective Actions

- *Describe technical or process changes that are needed to prevent this issue from happening again in the future.*

Optimize slow query in drupal.

Priority	Issue	Corrective Action	Owner	Result	Completion Date	Notes
High	<u>CPDRUPAL-4059</u>	Optimize slow query in Drupal. Jason reports he's optimized it from 4 seconds to 70ms	Jason Smith	In progress		

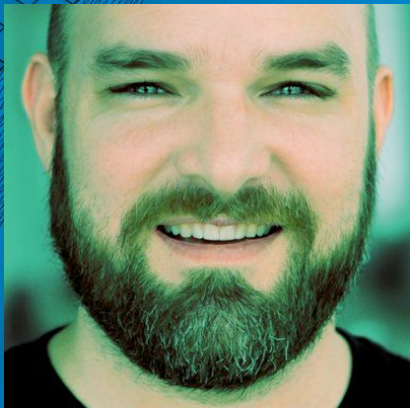
Summary

- Think of monitoring by function not by tool
- Catch signals by implementing each function
- User centric blackbox availability
- Improve signal to noise ratio with SLIs SLOs
- Make Ops happy with Alert Orchestration

Happy Customers

Questions?

THANK YOU



Jared Sprague
Principal Software Engineer

jsprague@redhat.com
@caramelcode



Adam Minter
Associate Project Manager

aminter@redhat.com