# Winston:  Helping Netflix Engineers Sleep at Night

Our journey… assisting engineers reduce operational load and MTTR

# On-Call !

# Sayli Karmarkar

Senior Software Engineer
Diagnostics and Remediation Engineering (DaRE)

skarmarkar@netflix.com

@HikerTechy

https://www.linkedin.com/in/saylikarmarkar

**DaRE Team's Focus**
Build platforms, tools and libraries
to help teams reduce MTTR for operational issues.

# Traditional On-Call Timeline

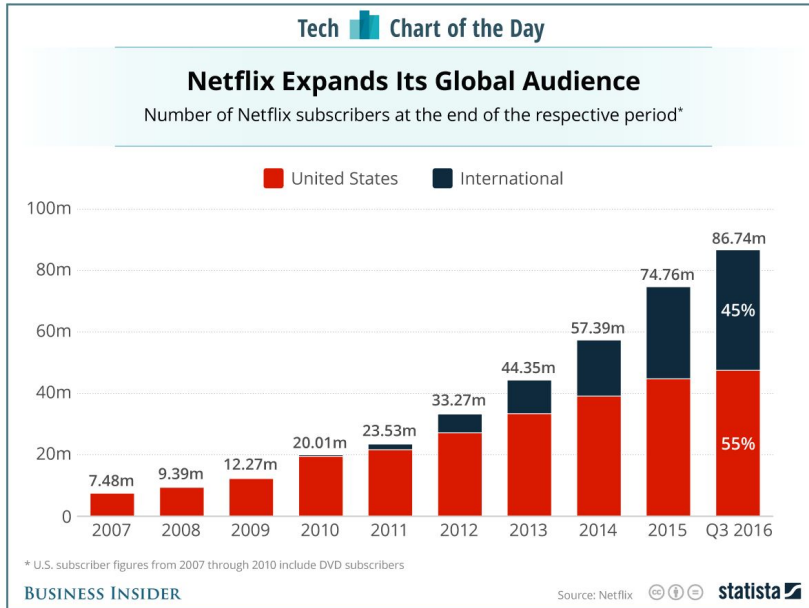| 2:00 AM | 2:02 AM | 2:07 AM | 2:10 AM | 2:15 AM | 2:20 AM | 2:30 AM |
|---------|---------|---------|---------|---------|---------|---------|
| PagerDuty Alert | Engineer Wakes up | Logs in and ACK | Studies the alert | Checks runbook | Runs diagnostics | Fixes/Mitigates the problem |

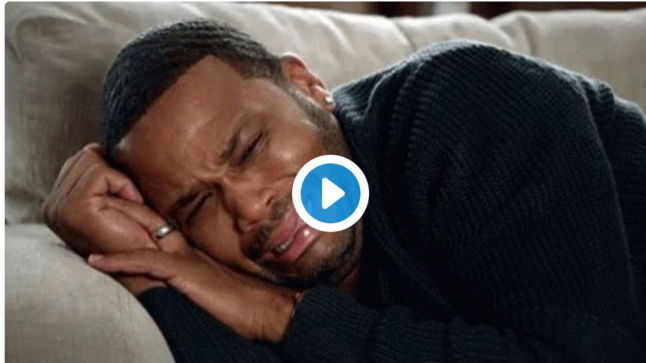# Works, but does it scale?!



**Scale and Growth**



**Availability**

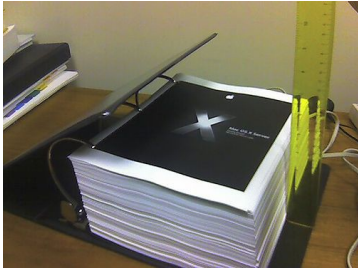# Netflix goes down. Twitter blows up



**SUNNI** ✓
@SunniAndTheCity
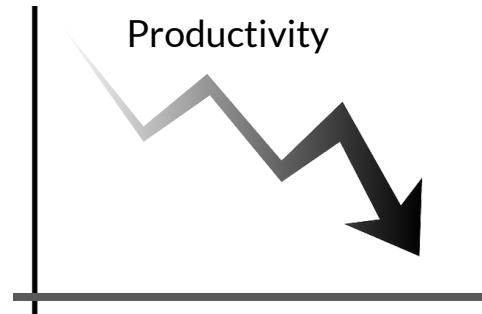
Follow

Netflix is down on a rainy Saturday afternoon.

12:42 PM - 1 Oct 2016

↩   ⟲ 105   ♥ 105

# Traditional On-Call Pain Points



**MTTR**

Productivity

# Solution?

**Automate**

- Removing False Positives
- Collecting Diagnostic Information
- Mitigating the problem to reduce impact on the customers

**Hands-free**

Feed the runbooks to an event-driven automation platform and have them executed in response to operational events

# Unique Problem? Not really ..

# Define

- **Business Goals**
- **Use-cases**
- **Customers**
- **Constraints**
- **Interactions with other services**

# Winston's Goals

- Assist engineers in **reducing MTTR and pager fatigue** by providing a platform to automate their runbooks
- Provide an **easy way to connect automated runbooks to an event**
- Let engineers **focus on the business logic** of runbooks rather than infrastructure aka PaaS.
- Provide appropriate **wrappers and libraries** to interact with other services
- Ensure **best practices** for automations and runbook lifecycle management

# What is Winston?

Winston is an event driven runbook automation platform. It is designed to host and execute runbooks in response to operational events.

# Traditional On-Call Timeline

| **2:00 AM** | 2:02 AM | 2:07 AM | 2:10 AM | 2:15 AM | 2:20 AM | **2:30 AM** |
|---|---|---|---|---|---|---|
| PagerDuty Alert | Engineer Wakes up | Logs in and ACK | Studies the alert | Checks runbook | Runs diagnostics | Fixes the problem |

# On-Call With Winston

2:00 AM

ALERT

Winston
I'M WINSTON WOLFE. I SOLVE PROBLEMS.

*False Positive*

*Mitigates the problem*

*Assisted Diagnostics*

2:05 AM

2:05 AM

2:15 AM

# Evaluation - Build / Reuse / Buy

# Stackstorm

**+**

- A generic pluggable Event-Driven Automation Platform
- Designed with availability and reliability in mind
- Open source + Code following good design practices
- Good alignment with respect to goals and future direction

**-**

- High availability and reliability not exercised a lot
- Dependency on MongoDB and RabbitMQ
- No easy way of adding and updating automation

# Good Starting Point ..



**As a Service (High Availability and Reliability )**

**Iterate and Evaluate Regularly**

# A closer look at a Winston Instance

# V1.0 Winston HA Deployment

# Challenges

- Added cognitive load resulting in less adoption

- How to help engineers choose operational efficiency over new features?

- Recommended and safe automation and lifecycle management practices are often not followed

- Simple use-cases are not trivial to on-board

# Winston Studio

- One stop portal for all things Winston
- Supports **C**reate, **R**ead, **U**pdate, **D**elete, **E**xecute and **D**iagnose functionality
- Implements best practises
  - Compliance/Auditing
  - Persistence
  - Security (Authentication/Authorization)
- Self serve & scalable

# Winston Studio

# Runbook View

# Executions

**Execution Details**

netflix_datapipeline » broker_offline_process_alert  PROD

Jul 18th 2016, 12:53 AM

Execution: 578c8ae293582e1450d327b1

▽ Parameters

Execution Region *  us-east-1

Region in which your Automation will be executed by Winston

alert_env *  prod

alert_matchset *  ▭ ✕   Enter a value and press return to add it to the array

A comma separated list of instance ids

alert_metadata  {"action":{"actionKey":"sqs","actionStatus":"done","incidentKey":["sqs","BrokerOffline","kafka","prod",[▭]],"us

alert_name  BrokerOffline

alert_region *  us-east-1

team_email  ▭@netflix.com

Comma separated list of email addresses the alert email should be sent to

timeout  7200

Copy Parameters                                                          ▷ Run in PROD

▽ Standard Output (stdout)

```
1 Instance Ids: ['▭']
2 Active instance states for app kafka : []
3 Active instance states for app kafkabroker : []
4 Active instance states for app kskafka : []
5 Active instance states for app k2 : []
6 Active Instances: {'k2': [], 'kafkabroker': [], 'kskafka': [], 'kafka': []}
7
```

▽ Log Output (includes stderr)

```
1 2016-07-18 00:53:07,482 INFO    Loading/Refreshing App kafka
2 2016-07-18 00:53:10,468 INFO    App kafka loaded/refreshed
3 2016-07-18 00:53:10,472 INFO    Loading/Refreshing App kafkabroker
4 2016-07-18 00:53:12,036 INFO    App kafkabroker loaded/refreshed
5 2016-07-18 00:53:12,037 INFO    Loading/Refreshing App kskafka
6 2016-07-18 00:53:15,576 INFO    App kskafka loaded/refreshed
7 2016-07-18 00:53:15,591 INFO    Loading/Refreshing App k2
8 2016-07-18 00:53:18,390 INFO    App k2 loaded/refreshed
9 2016-07-18 00:53:19,027 INFO    Email sent to ['▭@netflix.com'] with subject=[[prod]
  Summary: prod us-east-1 BrokerOffline]
10
```

Jul 10, 2016 09:12 PM          18          prod          ["i-de7b6a42"]          BrokerOffline          us-east-1

# Current Winston Deployment

# Use-case Patterns

# Sample Use-cases

**False Positives**
- Broker reporting offline when AWS maintenance takes down an instance
- Cassandra ring health

**Diagnostics** - Correlation could point towards the root cause
- Checking current maintenance jobs running on a cluster when an issue occurs
- Querying dependencies upstream and downstream for anomalous behavior
- Capture current system state and logs to analyze failures and reach the root cause quicker

**Mitigation**
- Restart kafka process
- Clean up disk space

# Alert: test us-east-1 BrokerOffline    alerts/Winston   x

**winston via PAE winston alerts** <pae-winston-alerts(       3:28 PM (0 minutes ago)
to data-pipeline-.

This is an alert generated by Winston - Automated Troubleshooting and Remediation Platform.
Winston execution ID:  56314b959287d930566e34a2

PROBLEM : Following kafka broker instances were reported to be offline.

Instances terminated by AWS --

    i-cd939c7c :  kafkabroker-logtrace-us-east-1d

Alert Snapshot: http://alert-history.us-east-1.test.netflix.net/history/snapshot/kafka/BrokerOffline/us-east-1/1446071188127?checkTime=1446071188127&sourceInstance=i-c5e8ce65

If you need to, you can look at the execution details of the winston workflow at http://winstoncde-useast1c.test.netflix.net:8080/#/history/56314b959287d930566e34a2/general

# Alert: test eu-west-1 cass_pay_1-disk_space_critical

alerts/Winston    x      CDE    x      Winston    x

**winston@ubuntu.netflix.com**                    Sep 23

to cde-team, pae-winston-al.

Action Needed: True

This alert is generated by Winston - Automated Troubleshooting and Remediation Platform.
Winston execution ID: 5602b71accfde20c9742919b

PROBLEM: Instance: [i-a362390e] for App: [cass_pay_1] is reporting high disk space usage

File system with high usage: md0

File system percent used: 88%

/data/cassandra070/data size in kb: 1533394580

Internal Compaction running: False

Repair/Compaction job running: False

Attempted removal of old snapshots: True

Remediation Message: Could not recover any space after running snapshot cleanup. Manual
intervention required.

Alert Snapshot: http://alert-history.eu-west-1.test.netflix.net/history/snapshot/cde-cass-disk_space_critical-test/cass_pay_1-disk_space_critical/eu-west-1/1443018180000?checkTime=1443018521473

# Alert: prod us-west-2 BrokerOffline

Inbox  x    alerts/Winston  x    **Winston  x**

**winston@ubuntu.netflix.com**                    Oct 3 (6 days ago)

to data-pipeline-., pae-winston-al.

This is an alert generated by Winston - Automated Troubleshooting and Remediation Platform.
Winston execution ID: 56101e8b8fd4d808436c840a

PROBLEM : Following kafka broker instances were reported to be offline.

kafkabroker  i-901b5067 -  Instance is in 'running' state.

   No disk failure detected. Kafka broker restarted Successfully.

Alert Snapshot: http://alert-history.us-west-2.prod.netflix.net/history/snapshot/kafka/BrokerOffline/us-west-2/1443896640000?checkTime=1443896969030

If you need to, you can look at the execution details of the winston workflow at http://winstoncde-uswest2a.prod.netflix.net:8080/#/history/56101e8b8fd4d808436c840a/general

If any of the above troubleshooting information has errors in it or if you have suggestions for how it can be improved, please file a JIRA ticket for PAE team using https://jira.netflix.com/secure/CreateIssueDetails!init.jspa?pid=17043&issuetype=4&components=23786&priority=4

# The Road Ahead

- **Adoption / Usability**
  - Find common operational use-cases and allow them to be re-used
  - Improve discoverability of Winston by integrating into existing alerting systems
  - Polyglot support (Groovy based runbooks)

- **Safety**
  - Resource isolation using containers
  - Rate limiting capability

- **Stronger analytics**
  - Provide aggregate visualization of runbook executions

# Key Takeaways

- **Don't re-invent** the wheel
- **Start simple and iterate**. Have some room for experimentation.
- Start with use-cases where there is **more pain and less control** over the source of the problem
- Pay special attention to **usability** of your product
- Push for **changing the culture** -- Usage will follow
- Find **sponsors** for features that are much more involved
- **Implement** best practices through carefully designed user interface instead of documentation.
- **Discourage anti-patterns** that focus on long term mitigation rather than fixing the root-cause
- Talk to us and others to share insights and learnings

# **Resources**

- Winston Tech Blog link
  http://techblog.netflix.com/2016/08/introducing-winston-event-driven.html

- Stackstorm documentation
  https://docs.stackstorm.com/

- Reach out
  skarmarkar@netflix.com

  @HikerTechy

**We are hiring**
Senior Software Engineer - https://jobs.netflix.com/jobs/860752