# Monitoring Solr at Scale

**Bloomberg** Engineering

**Berlin Buzzwords**
**June 18, 2019**

**Ken LaPorte**
**Team Lead, Search Infrastructure**

**TechAtBloomberg.com**

# Has this happened to you?

# What's a Bloomberg?

- **Largest provider of financial news and information**

- **Our strength is quickly and accurately delivering data, news and analytics**

- **Creating highly-performant and accurate information retrieval systems is core to our strength**
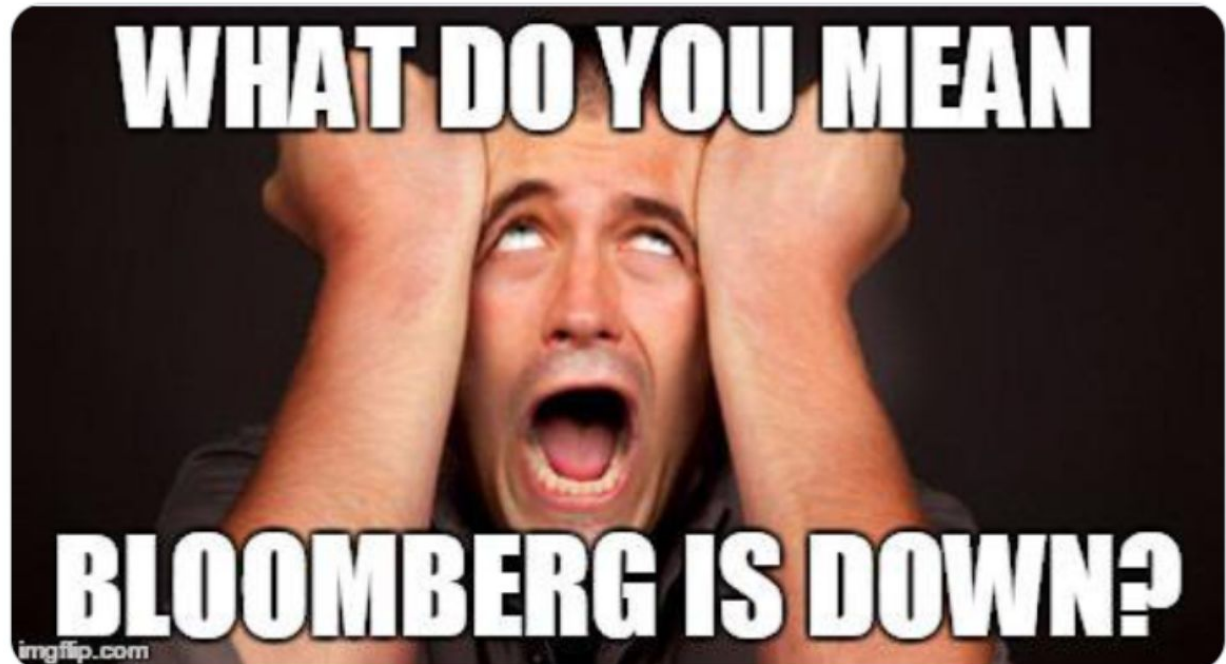
**Bloomberg**

Engineering

**Advisory Desk** @advdesk · Apr 17, 2015
Maybe it could be a Blockbuster #BloombergDown

I SURVIVED THE
BLOOMBERG CRASH
FRIDAY APR 17, 2015

**Ian Shepherdson** @IanShepherdson · Apr 17, 2015
Oxygen! I need oxygen! #bloombergdown

WHAT DO YOU MEAN

BLOOMBERG IS DOWN?
imgflip.com

# Search Infrastructure at Bloomberg

- **Hundreds of search applications & ZK Ensembles**
  - Diverse use cases and scale
  - Displaced other technologies
- **10s of billions documents**
- **100s of millions new documents daily**
- **1,000s of servers**
- **10,000s of Solr instances**
- **10,000s of queries per second**
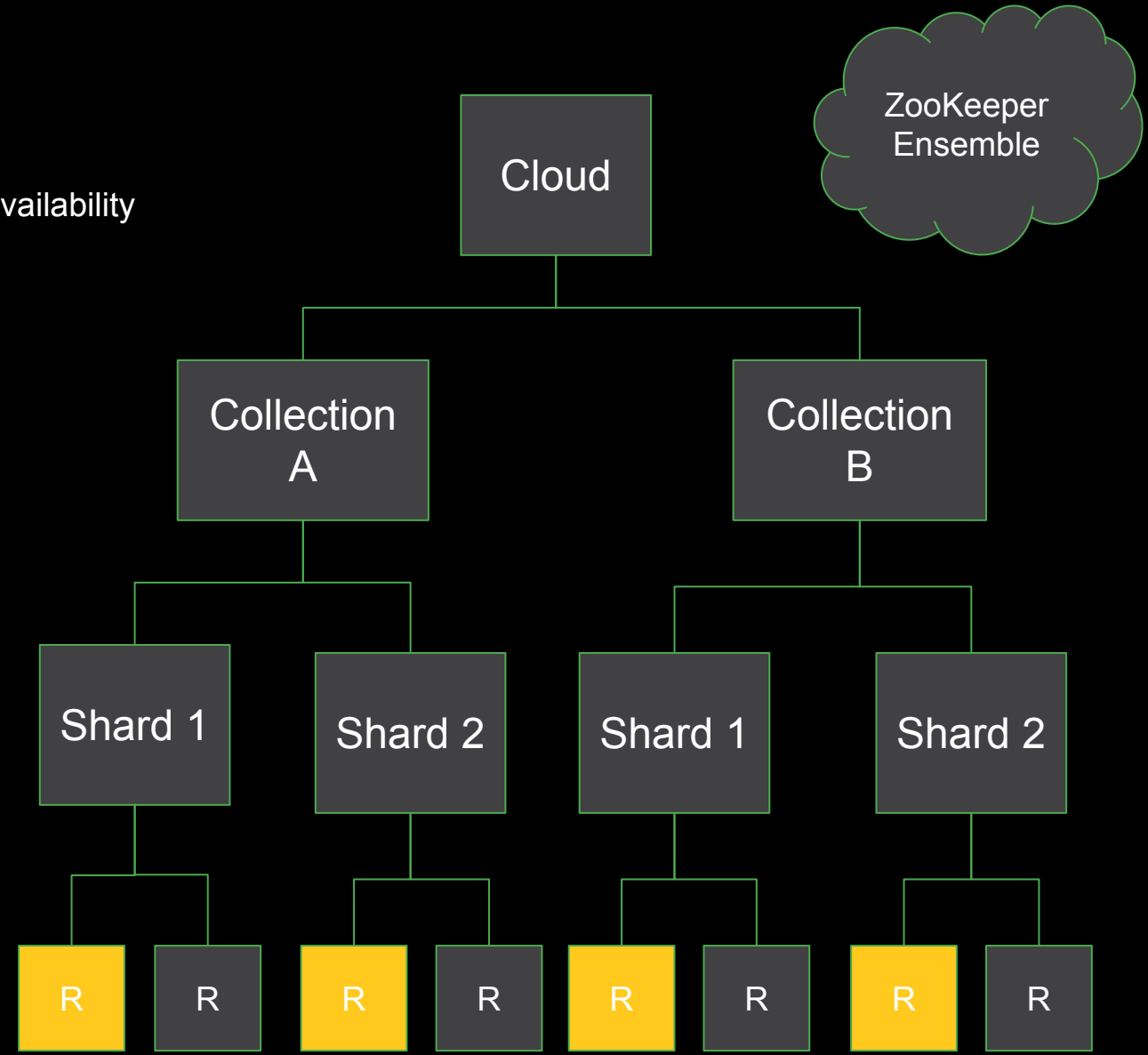- **Critical to Bloomberg and the global financial markets**

# Solr Open Source Contributions at Bloomberg

- **Analytics component**
- **Streaming expressions**
- **Learning to rank**
- **Bug fixes**
- **3 committers and 2 PMC members**
    - **On teams throughout Bloomberg**

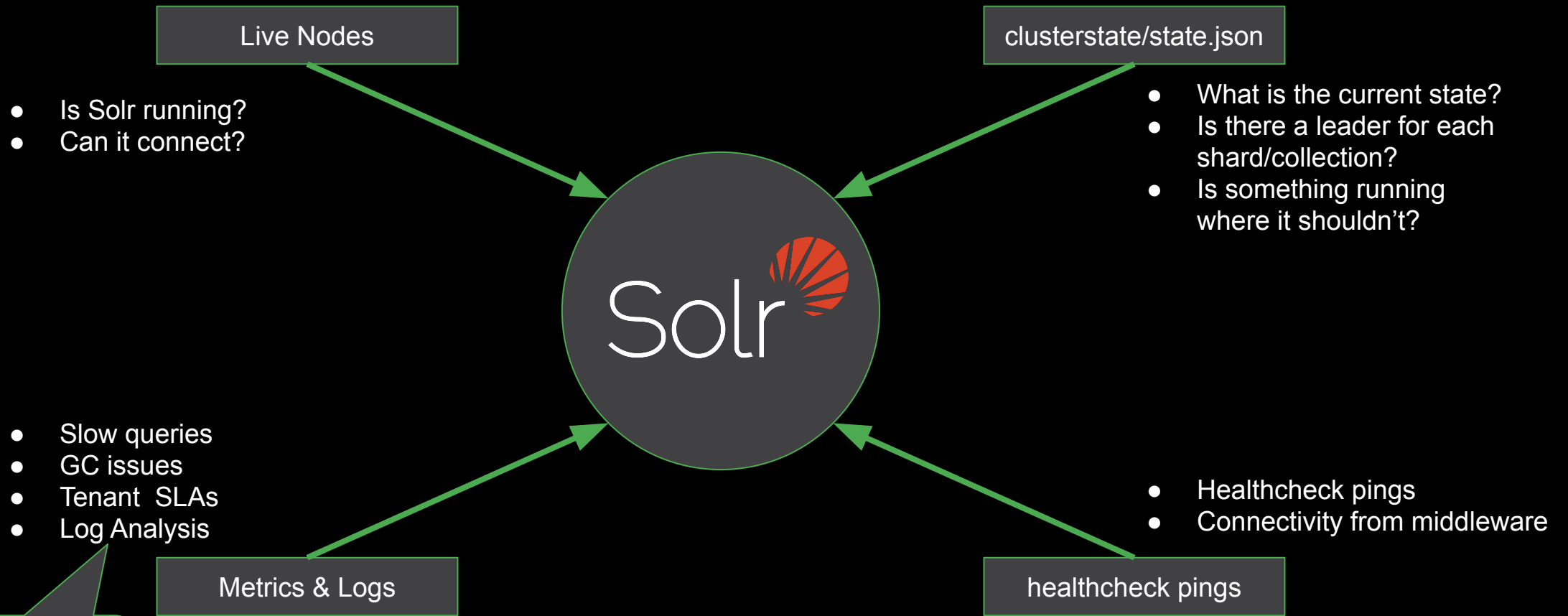Logos courtesy Apache Software Foundation

6

# Some Terms

- Solr Cloud
  - A cluster setup that allows for fault tolerance and high availability

- Collection
  - A complete logical index

- Shard
  - A logical piece (or slice) of a collection

- Replica
  - Simply a copy of a shard
- Leader
  - The shard replica responsible for indexing
  - Represented in yellow

- ZooKeeper
  - A consensus management system used by Solr
- Service Provider
  - The service management team (Me)
- Tenant
  - The application team using a service
  - Hostile users <- Credit Anshum
- Alarming
  - Active notification of relevant events
- Monitoring
  - Next slide please



*Leader
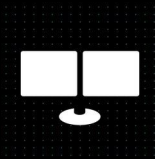
# What does it mean to monitor Solr?

Live Nodes

clusterstate/state.json

- Is Solr running?
- Can it connect?

- What is the current state?
- Is there a leader for each shard/collection?
- Is something running where it shouldn't?

- Slow queries
- GC issues
- Tenant  SLAs
- Log Analysis

- Healthcheck pings
- Connectivity from middleware

Metrics & Logs

healthcheck pings

Anshum Gupta's
Self defence talk

# What about ZooKeeper?

ZooKeeper 3.5 5

Dyn Reconfig.
Jetty/Rest

ruok

- IMOK
- Basic status/connectivity

mntr

- ZooKeeper stats
  - Size
  - # Watches
  - # ZNodes
  - Latency
- Version
- Leadership

- # Connections
- # ZNodes
- Latency
- Log Analysis

Metrics & Logs

- Connection details

cons

# Can you see me now?

- Initial version

- Single Collection View

- Cross Environments

- Point in time

# Service Provider View

- Display state of all clouds

- Switch between environments

- Bulk commands

  - Restart clouds by rack

- Still point in time

- Scale still a problem

| Solr | Zookeeper | | | |
|---|---|---|---|---|
| ☑Active ☑Recovering ☑Down ☑Gone | | | | ☑Open Collection in New Window |
| | | | | |
| cloud500 | | myServer | myServer | |
| cloud501 | | myServer | myServer | |
| cloud502 | | myServer | myServer | |
| cloud503 | | myServer | myServer | |
| cloud504 | | myServer | myServer | |
| cloud505 | | myServer | myServer | |
| cloud506 | | myServer | myServer | |
| cloud507 | | myServer | myServer | |
| cloud508 | | myServer | myServer | |
| cloud509 | | myServer | myServer | |
| cloud51 | | myServer | myServer | |
| cloud510 | | myServer | myServer | |
| cloud511 | | myServer | myServer | |
| cloud512 | | myServer | myServer | |
| cloud513 | | myServer | myServer | |
| cloud514 | | myServer | myServer | |
| Legend | Good | Recovering Down | Offline | |
| Bulk | | | Start All | Stop All | Env. DEV |

Names changed for display purposes

# Hello Niteowl!

- Active monitor of state changes

- Backend

  - Intelligent alarms

- Frontend

  - More scalable

  - Service provider only



Data for display purposes only.
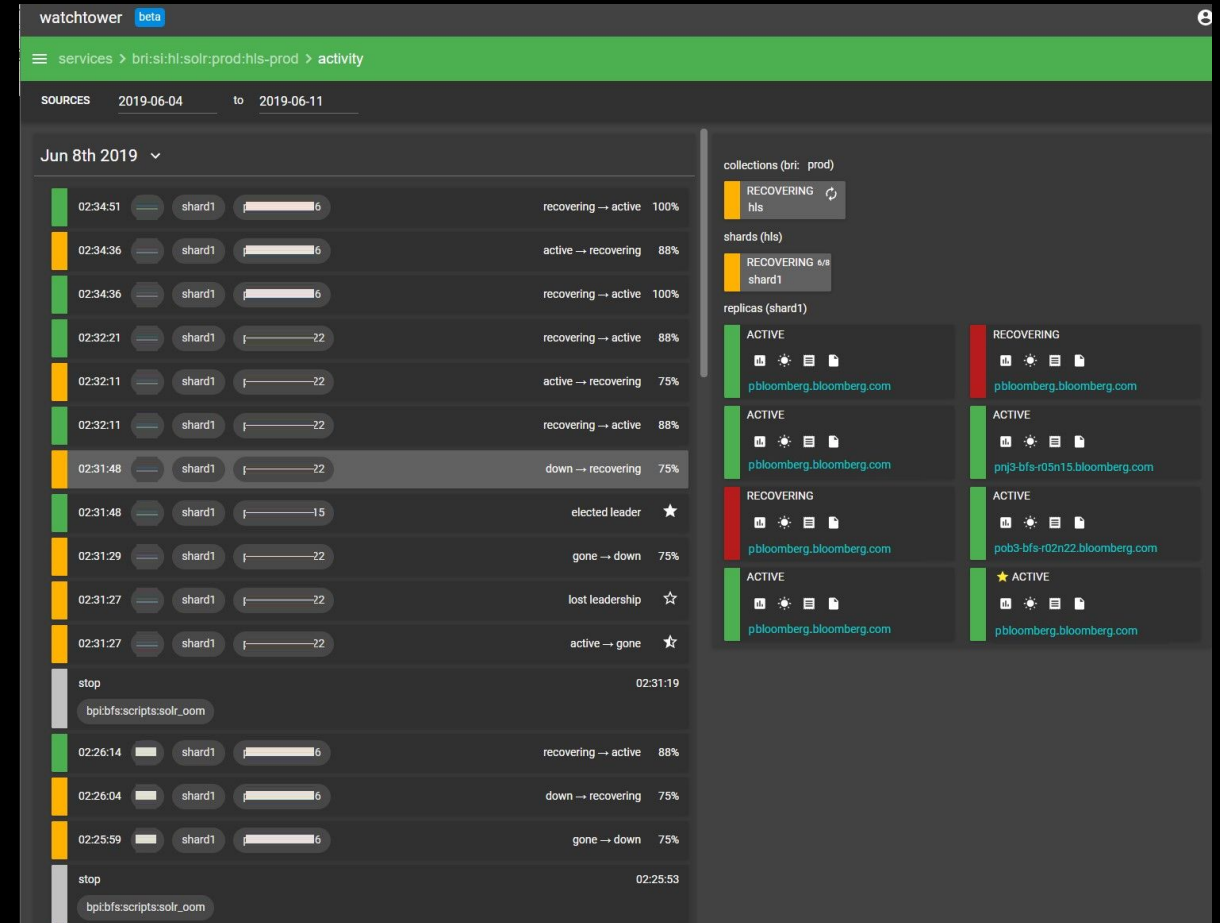
# All along the WatchTower….

- Allows Svc provider & tenants to share a UI

- Leverages same Niteowl back-end

- Day 2 actions
  - Process management (start, stop, etc.)
  - Schema changes
  - Query interface
  - Metrics
  - Security



Credit Jimi Hendrix & the Noun Project
Names blurred/changed for IP

# Alarming, Activity & Analysis

- Alarming
  - Cloud aware
  - Consolidation across signals
  - Delay Alarms to prevent flutter

- Activity
  - Auditability is critical
  - Every state change is recorded
  - Each user action is logged
  - Additional sources (OOMs)

- Analysis
  - Events indexed into Solr
  - Analytics helped address multitude of issues
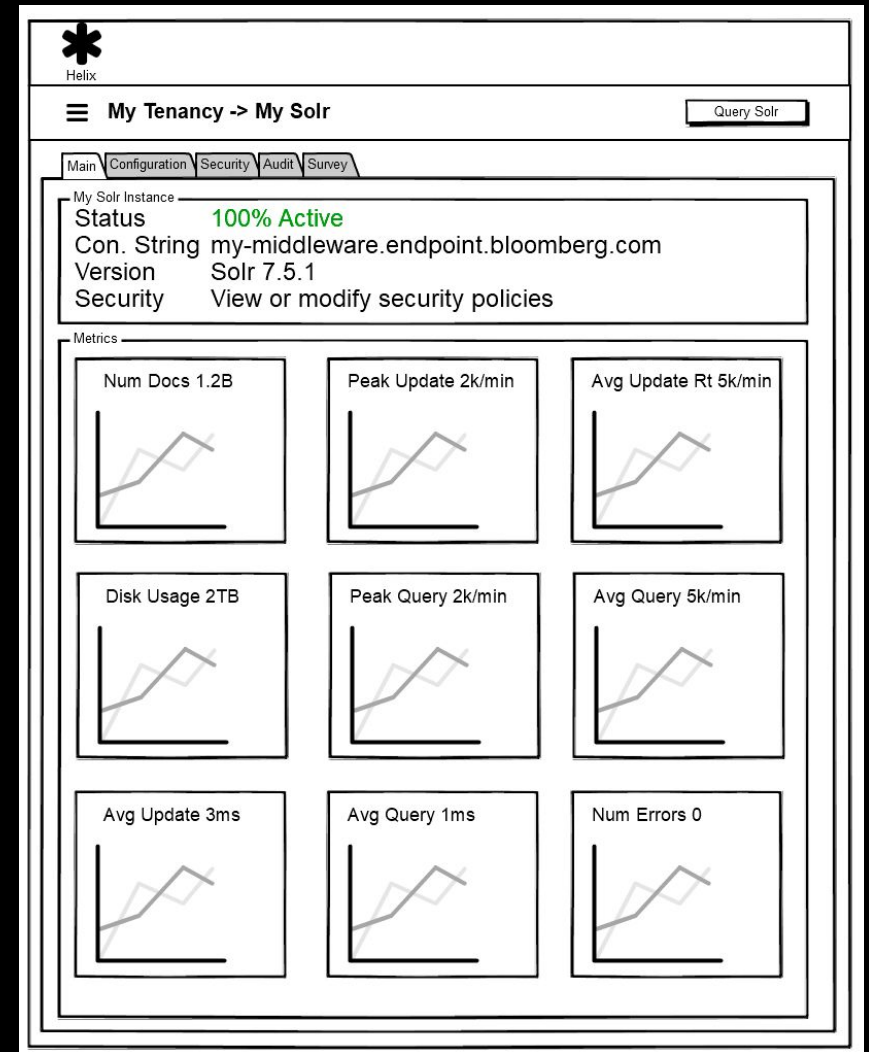
# Slice & Dice

- Slice by:
  - Host list
  - State
  - and more

- Dice by:
  - Technology
  - Tenancy
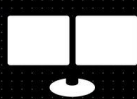  - Network Segment
  - and more



Names blurred/changed for IP

# nth time's a charm

- Enter Helix

- New UI
  - Shared among service providers
  - Think AWS console

- New Purpose
  - Client centric
  - Administration through automation

- Obviously still in wireframe/planning stages

- Our final UI 😁

# + = A Good Night's Sleep

- Chat Ops Goals:
  - Laptops off
  - Only wake up when necessary
  - Read-only

- Commands
  - What's wrong (sup)?
  - Detailed views (how is *foo*)
  - Host by host breakdown (hosts)
  - More to come

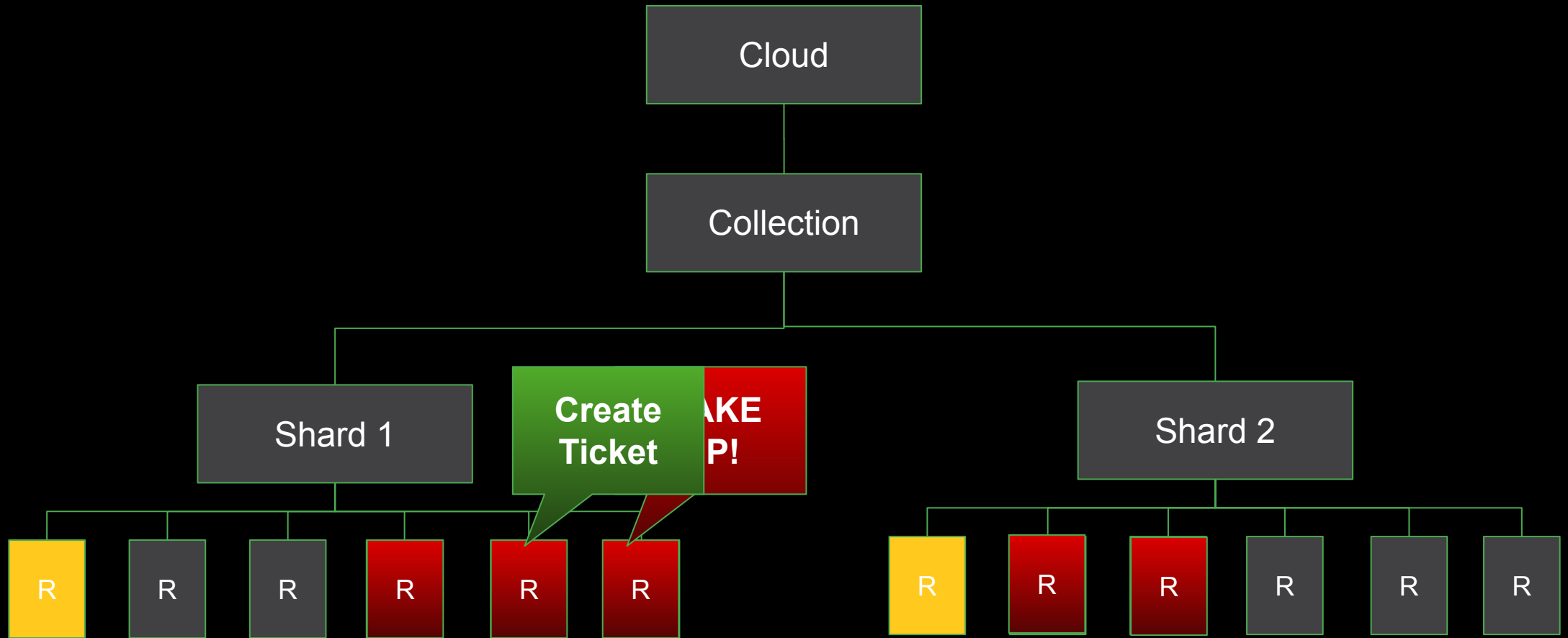- Publication of events



```
Ken Laporte   Sup
BFS Bot  APP  zk:
    bri:si:bfs:zk:dev:dan1-qa-dev
Ken Laporte   how is dan1-qa
BFS Bot  APP  bri:si:bfs:zk:dev:dan1-qa-dev [zk] has 1 node in the following state:
    1 gone
Ken Laporte   hosts
BFS Bot  APP  Issues found in the following hosts:
    d.bloomberg.com [ X / Y ]

    1   zk                [dan1-qa]
```



```
UPDATE

Service: bri:si:m:solr:dev:m-dev
Alarm Type: com.bloomberg.bfs.banshee.service.alarms.ReplicasDownAlarmRule@39a77334
Severity: BF
Associated DRQS: NONE
Underlying Conditions:
bri:si:msg:solr:dev:m-dev#m-p001#shard5 has only 1/2 replicas healthy
bri:si:msg:solr:dev:m-dev#m-p001#shard4 has only 1/2 replicas healthy
bri:si:msg:solr:dev:m-dev#m-p001#shard12 has only 1/2 replicas healthy
```
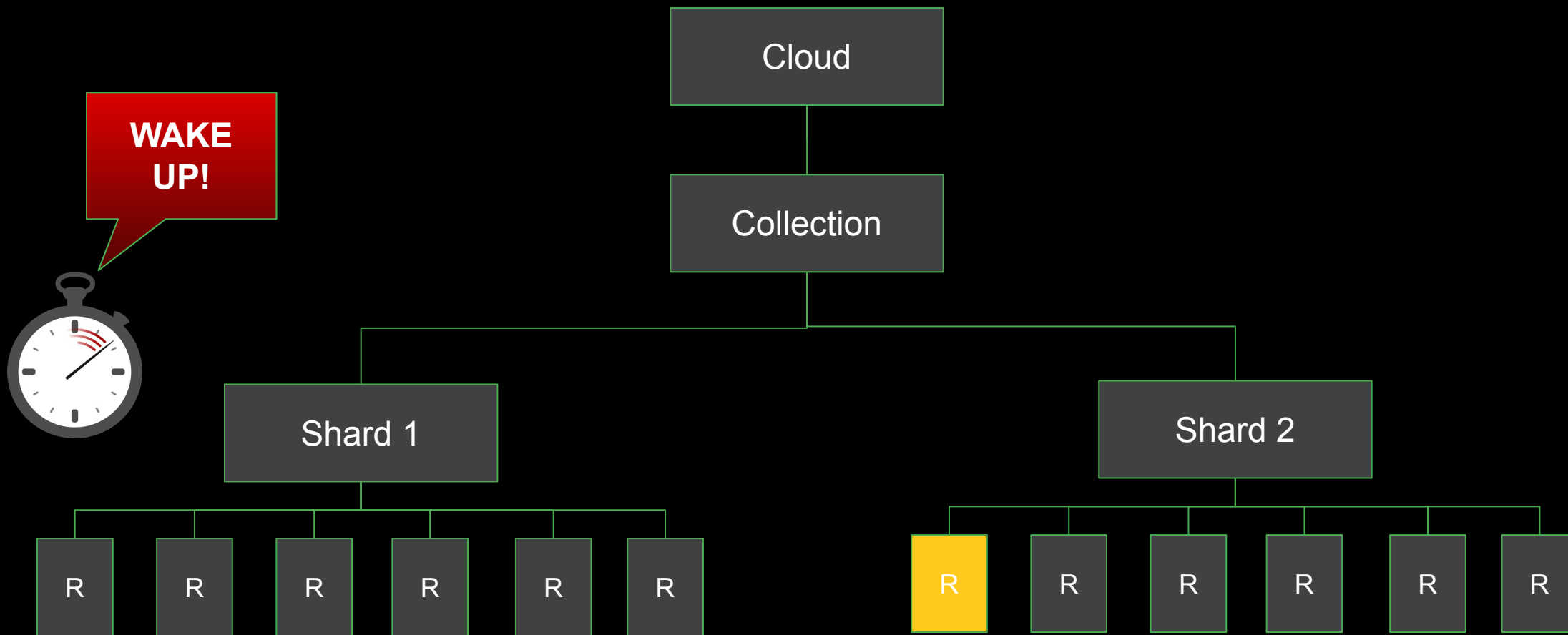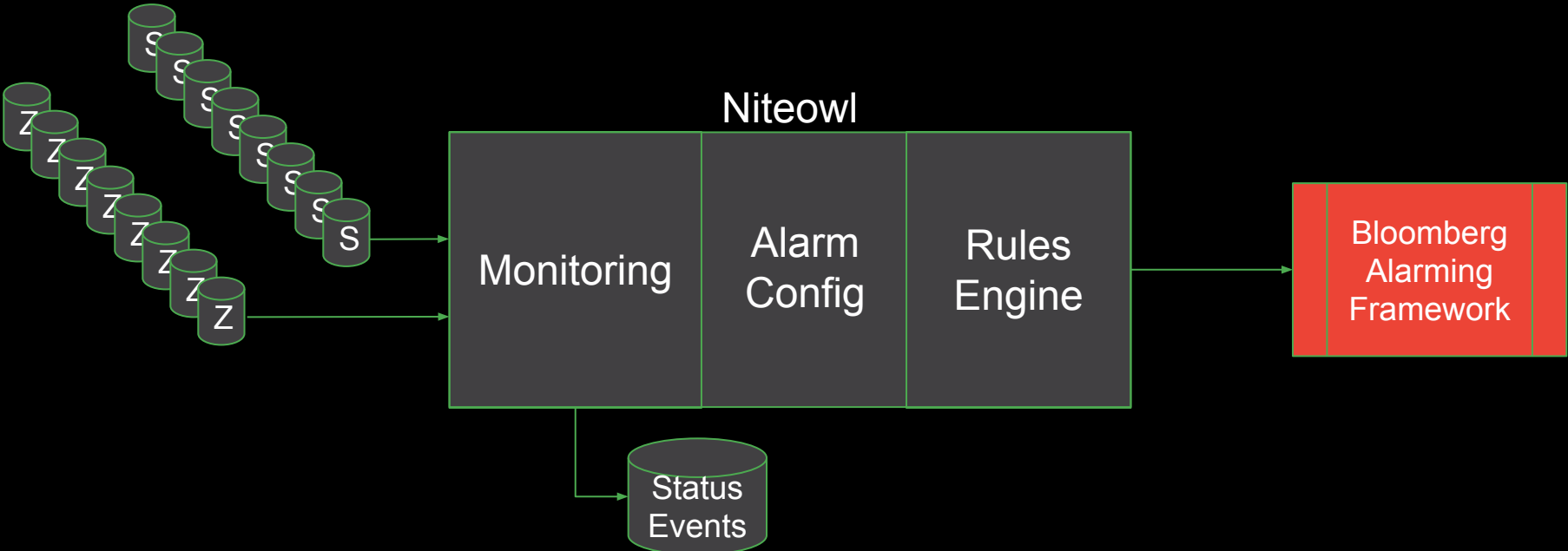
# So what does all this look like?

# So what does all this look like?

Niteowl

| Monitoring | Alarm Config | Rules Engine |
|---|---|---|

Bloomberg Alarming Framework

# Enter Kafka & Banshee

Niteowl | Kafka | Banshee

Monitoring

Event
Publishers

Rules
Engine

Alarm
Config

Event
Consumers

Status
Events

Bloomberg
Alarming
Framework
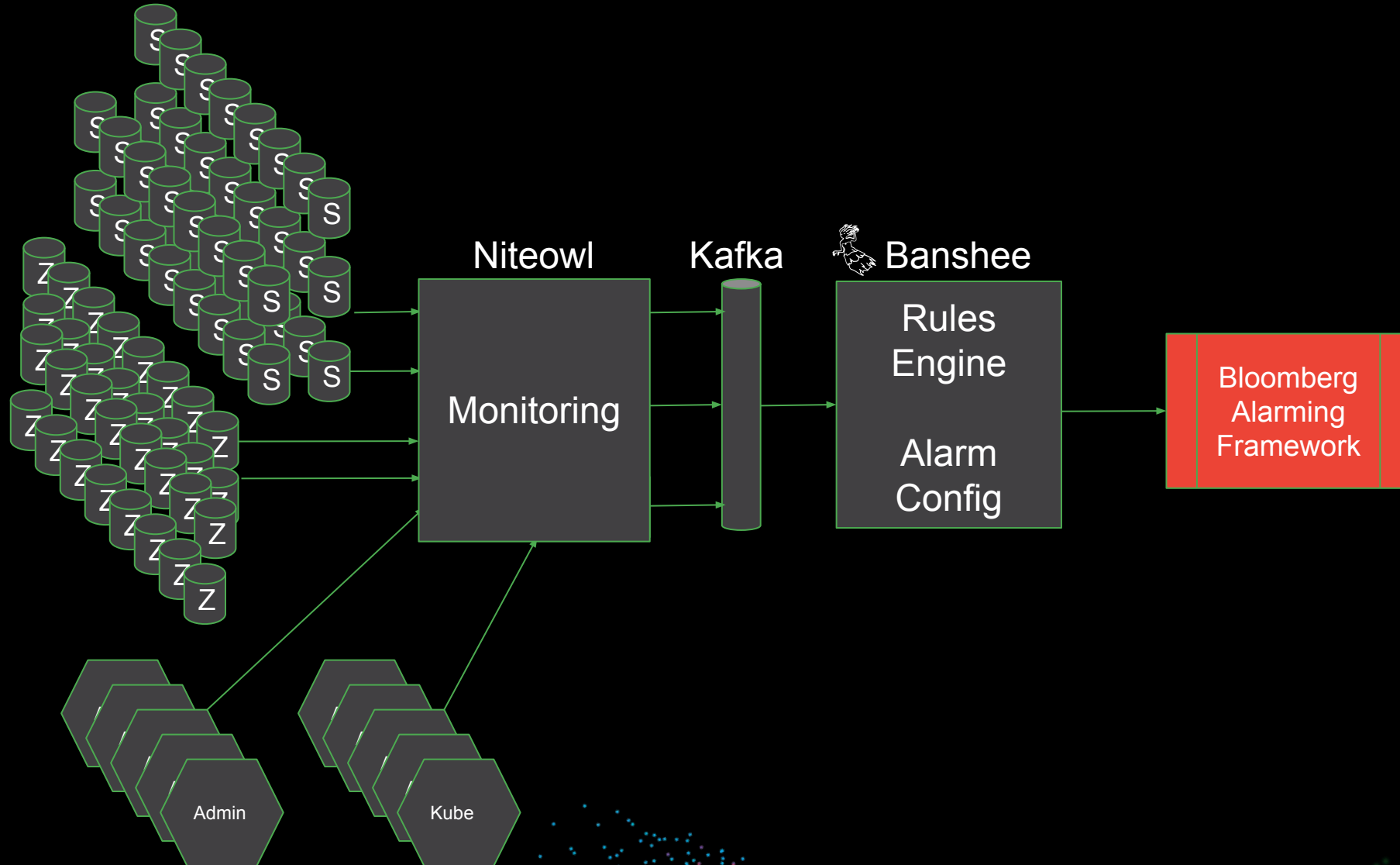
# Run Niteowl via Kubernetes and spread the load



23

# What's Next?

- 2019 Goal: Rules Engine
  - Leverage Kafka Streams to create robust alarms engine

- 2019 Goal: Open Sourcing Niteowl
  - Reality: Tightly integrated into the Bloomberg Infrastructure

- 2020 Goal: Detangle Niteowl and hopefully publish as open source

# Thank You, Sleep Well!

**https://www.bloomberg.com/careers**

# Questions?

Ken LaPorte
klaporte@bloomberg.net

**Engineering**

**Bloomberg**

**TechAtBloomberg.com**