# Netflix and Open Source

## April 2013

## Adrian Cockcroft

@adrianco #netflixcloud @NetflixOSS

http://www.linkedin.com/in/adriancockcroft

# Cloud Native

# NetflixOSS – Cloud Native On-Ramp

# Netflix Open Source Cloud Prize

# We are Engineers

We solve hard problems

We build amazing and complex things

We fix things when they break

# We strive for perfection

Perfect code

Perfect hardware

Perfectly operated

# But perfection takes too long…

So we compromise

Time to market vs. Quality

Utopia remains out of reach

# Where time to market wins big

Web services

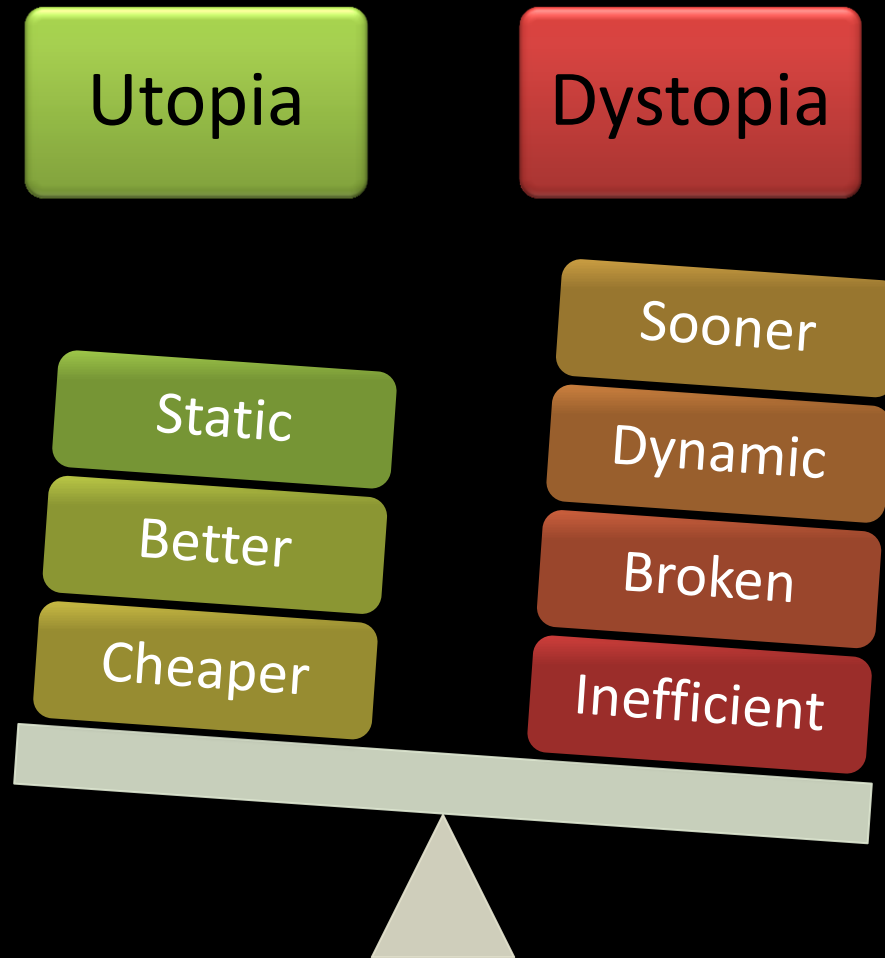Agile infrastructure - cloud

Continuous deployment

# How Soon?

Code features in days instead of months

Hardware in minutes instead of weeks

Incident response in seconds instead of hours

# Tipping the Balance

**Utopia**

**Dystopia**

Static

Sooner

Better

Dynamic

Cheaper

Broken

Inefficient

# A new engineering challenge

Construct a highly agile and highly available service from ephemeral and often broken components

# Cloud Native

How does Netflix work?

# Netflix Member Web Site Home Page
## Personalization Driven – How Does It Work?

# How Netflix Streaming Works

# Content Delivery Service
## Open Source Hardware Design + FreeBSD, bird, nginx

# November 2012 Traffic

| Rank | Upstream | | Downstream | | Aggregate | |
|---|---|---|---|---|---|---|
| | Application | Share | Application | Share | Application | Share |
| 1 | BitTorrent | 36.8% | Netflix | 33.0% | Netflix | 28.8% |
| 2 | HTTP | 9.83% | YouTube | 14.8% | YouTube | 13.1% |
| 3 | Skype | 4.76% | HTTP | 12.0% | HTTP | 11.7% |
| 4 | Netflix | 4.51% | BitTorrent | 5.89% | BitTorrent | 10.3% |
| 5 | SSL | 3.73% | iTunes | 3.92% | iTunes | 3.43% |
| 6 | YouTube | 2.70% | MPEG | 2.22% | SSL | 2.23% |
| 7 | PPStream | 1.65% | Flash Video | 2.21% | MPEG | 2.05% |
| 8 | Facebook | 1.62% | SSL | 1.97% | Flash Video | 2.01% |
| 9 | Apple PhotoStream | 1.46% | Amazon Video | 1.75% | Facebook | 1.50% |
| 10 | Dropbox | 1.17% | Facebook | 1.48% | RTMP | 1.41% |
| | Top 10 | 68.24% | Top 10 | 79.01% | Top 10 | 76.54% |

sandvine

Table 3 - Top 10 Peak Period Applications (North America, Fixed Access)

# Real Web Server Dependencies Flow

(Netflix Home page business transaction as seen by AppDynamics)

Each icon is three to a few hundred instances across three AWS zones

Cassandra

memcached

Web service

S3 bucket

Start Here

Three Personalization movie group choosers (for US, Canada and Latam)

# Cloud Native Architecture

| | | |
|---|---|---|
| Clients | Things | |

| | | |
|---|---|---|
| Autoscaled Micro Services | JVM | JVM | JVM |

| | | |
|---|---|---|
| Autoscaled Micro Services | JVM | JVM | Memcached |

| | | |
|---|---|---|
| Distributed Quorum NoSQL Datastores | Cassandra | Cassandra | Cassandra |

Zone A          Zone B          Zone C

# Non-Native Cloud Architecture

**Agile Mobile Mammals**

iOS/Android

**Cloudy Buffer**

Server

**Datacenter Dinosaurs**

M... ...pps

# New Anti-Fragile Patterns

Micro-services

Chaos engines

Highly available systems composed
from ephemeral components

# Stateless Micro-Service Architecture

## Linux Base AMI (CentOS or Ubuntu)

Optional Apache frontend, memcached, non-java apps

Monitoring
Log rotation to S3 AppDynamics machineagent Epic/Atlas

### Java (JDK 6 or 7)

AppDynamics appagent monitoring

GC and thread dump logging

#### Tomcat

Application war file, base servlet, platform, client interface jars, Astyanax

Healthcheck, status servlets, JMX interface, Servo autoscale

# Cassandra Instance Architecture

## Linux Base AMI (CentOS or Ubuntu)

Tomcat and Priam on JDK

Healthcheck, Status

Monitoring

AppDynamics machineagent Epic/Atlas

### Java (JDK 7)

AppDynamics appagent monitoring

GC and thread dump logging

#### Cassandra Server

Local Ephemeral Disk Space – 2TB of SSD or 1.6TB disk holding Commit log and SSTables

# Cloud Native

Master copies of data are cloud resident

Everything is dynamically provisioned

All services are ephemeral

# Dynamic Scalability

# Asgard

# Cloud Deployment Scalability

New Autoscaled AMI – zero to 500 instances from 21:38:52 - 21:46:32, 7m40s
Scaled up and down over a few days, total 2176 instance launches, m2.2xlarge (4 core 34GB)

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|---------|--------|------|---------|------|
| 41.0 | 104.2 | 149.0 | 171.8 | 215.8 | 562.0 |



Instance Launch Duration for a large Autoscale Group

# Ephemeral Instances

- Largest services are autoscaled
- Average lifetime of an instance is 36 hours

# Leveraging Public Scale

1,000 Instances          100,000 Instances

**Public**     **Grey Area**     **Private**

Startups          Netflix          Google

# How big is Public?

AWS Maximum Possible Instance Count 3.7 Million
Growth >10x in Three Years,  >2x Per Annum



AWS upper bound estimate based on the number of public IP Addresses
Every provisioned instance gets a public IP by default

# Availability

Is it running yet?

How many places is it running in?

How far apart are those places?

# Antifragile API Patterns
## Functional Reactive with Circuit Breakers and Bulkheads

The STRANGE WORLD of the FUTURE

STRANDED without video!
No way to fill their empty hours!
They were victims of...

WAY OUT

THE CLOUD OF BROKEN STREAMS

# Outages

- Running very fast with scissors
  - Mostly self inflicted – bugs, mistakes
  - Some caused by AWS bugs and mistakes

- Next step is multi-region
  - Investigating and building in stages during 2013
  - Could have prevented some of our 2012 outages

# Managing Multi-Region Availability

NETFLIX
OSS

AWS
Route53

UltraDNS

DynECT
DNS

Denominator

Regional Load Balancers

Regional Load Balancers

Zone A
Cassandra Replicas

Zone B
Cassandra Replicas

Zone C
Cassandra Replicas

Zone A
Cassandra Replicas

Zone B
Cassandra Replicas

Zone C
Cassandra Replicas

A portable way to manage multiple DNS providers from Java

# Configuration State Management

Datacenter CMDB's woeful

Cloud native is the solution

Dependably complete

# Edda – Configuration History

Eureka Services metadata

AWS Instances, ASGs, etc.

AppDynamics Request flow

Edda

Monkeys

NETFLIX
SIMIAN ARMY

# Edda Query Examples

Find any instances that have ever had a specific public IP address
$ curl "http://edda/api/v2/view/instances;publicIpAddress=1.2.3.4;_since=0"
 ["i-0123456789","i-012345678a","i-012345678b"]

Show the most recent change to a security group
$ curl "http://edda/api/v2/aws/securityGroups/sg-0123456789;_diff;_all;_limit=2"
--- /api/v2/aws.securityGroups/sg-0123456789;_pp;_at=1351040779810
+++ /api/v2/aws.securityGroups/sg-0123456789;_pp;_at=1351044093504
@@ -1,33 +1,33 @@
 {
…
     "ipRanges" : [
       "10.10.1.1/32",
       "10.10.1.2/32",
+       "10.10.1.3/32",
-       "10.10.1.4/32"
…
 }

A Cloud Native Open Source Platform

# Inspiration

# Three Questions

Why is Netflix doing this?

How does it all fit together?

What is coming next?

# Beware of Geeks Bearing Gifts: Strategies for an Increasingly Open Economy

*Simon Wardley - Researcher at the Leading Edge Forum*

# How did Netflix get ahead?

**Netflix Business + Developer Org**
- Doing it right now
- SaaS Applications
- PaaS for agility
- Public IaaS for AWS features
- Big data in the cloud
- Integrating many APIs
- FOSS from github
- Renting hardware for 1hr
- Coding in Java/Groovy/Scala

**Traditional IT Operations**
- Taking their time
- Pilot private cloud projects
- Beta quality installations
- Small scale
- Integrating several vendors
- Paying big $ for software
- Paying big $ for consulting
- Buying hardware for 3yrs
- Hacking at scripts

NETFLIX
OSS

# Netflix Platform Evolution

2009-2010        2011-2012        2013-2014

Bleeding Edge Innovation → Common Pattern → Shared Pattern

Netflix ended up several years ahead of the industry, but it's not a sustainable position

# Making it easy to follow

**Exploring the wild west each time**   **vs. laying down a shared route**

# How does it all fit together?

# NetflixOSS Continuous Build and Deployment

# NetflixOSS Services Scope

## AWS Account

- Asgard Console
- Archaius Config Service
- Cross region Priam C*
- Explorers Dashboards
- Atlas Monitoring
- Genie Hadoop Services

### Multiple AWS Regions

- Eureka Registry
- Exhibitor ZK
- Edda History
- Simian Army

#### 3 AWS Zones

**Application Clusters**
Autoscale Groups
Instances

**Priam**
Cassandra
Persistent Storage

**Evcache**
Memcached
Ephemeral Storage

# NetflixOSS Instance Libraries

## Initialization
- Baked AMI – Tomcat, Apache, your code
- Governator – Guice based dependency injection
- Archaius – dynamic configuration properties client
- Eureka - service registration client

## Service Requests
- Karyon - Base Server for inbound requests
- RxJava – Reactive pattern
- Hystrix/Turbine – dependencies and real-time status
- Ribbon - REST Client for outbound calls

## Data Access
- Astyanax – Cassandra client and pattern library
- Evcache – Zone aware Memcached client
- Curator – Zookeeper patterns
- Denominator – DNS routing abstraction

## Logging
- Blitz4j – non-blocking logging
- Servo – metrics export for autoscaling
- Atlas – high volume instrumentation

# NetflixOSS Testing and Automation

**Test Tools**
- CassJmeter – Load testing for Cassandra
- Circus Monkey – Test account reservation rebalancing

**Maintenance**
- Janitor Monkey – Cleans up unused resources
- Efficiency Monkey
- Doctor Monkey
- Howler Monkey – Complains about expiring certs

**Availability**
- Chaos Monkey – Kills Instances
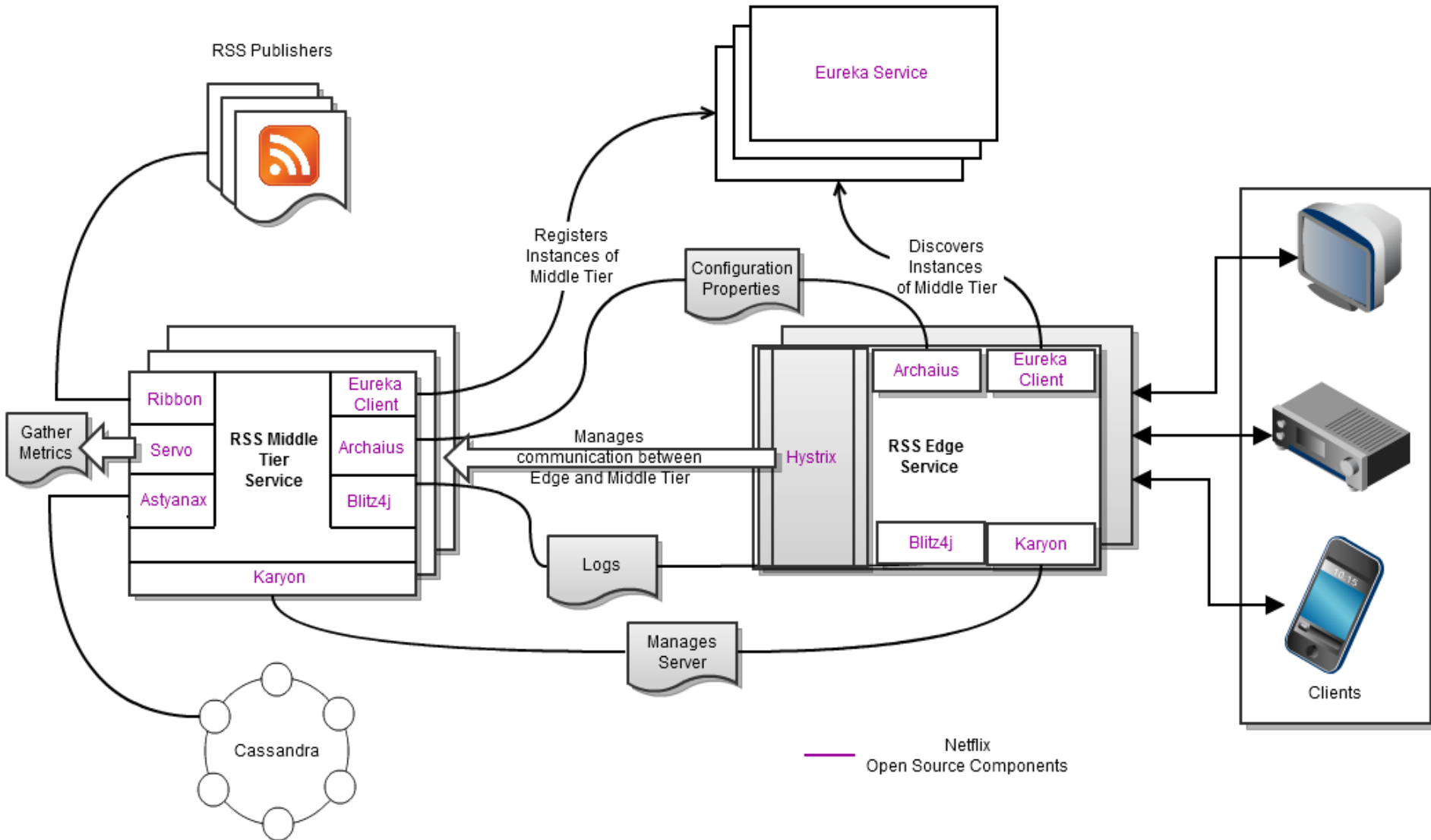- Chaos Gorilla –  Kills Availability Zones
- Chaos Kong – Kills Regions
- Latency Monkey – Latency and error injection

**Security**
- Security Monkey
- Conformity Monkey

# Example Application – RSS Reader

# What's Coming Next?

Better portability

Higher availability

Easier to deploy

Contributions from end users

Contributions from vendors

More Features

More Use Cases

NETFLIX OSS

# Vendor Driven Portability
## Interest in using NetflixOSS for Enterprise Private Clouds



"It's done when it runs Asgard"
Functionally complete
Demonstrated March
Release 3.3 in 2Q13

Some vendor interest
Many missing features
Bait and switch AWS API strategy

Some vendor interest
Needs AWS compatible Autoscaler

# AWS 2009
## Baseline features needed to support NetflixOSS

# Netflix Cloud Prize

## Boosting the @NetflixOSS Ecosystem

# In 2012 Netflix Engineering won this..

# We'd like to give out prizes too

But what for?
Contributions to NetflixOSS!
Shared under Apache license
Located on github

# Judges choice award

## Best example application mash-up

## Best usability enhancement

## Best portability enhancement

## Best new monkey

## Best new feature

## Best datastore integration

## Best contribution to code quality

## Best contribution to operational tools

## Best contribution to performance

# How long do you have?

Entries open March 13<sup>th</sup>

Entries close September 15<sup>th</sup>

Six months...

# Who can win?

Almost anyone, anywhere...

Except current or former Netflix or
AWS employees

# Who decides who wins?

Nominating Committee

Panel of Judges

Aino Corry
Program Chair for Qcon/GOTO

Simon Wardley
Strategist

Martin Fowler
Chief Scientist Thoughtworks

Werner Vogels
CTO Amazon

Joe Weinman
SVP Telx, Author "Cloudonomics"

Yury Izrailevsky
VP Cloud Netflix

# What are Judges Looking For?

Eligible, Apache 2.0 licensed

Original and useful contribution to NetflixOSS

Code that successfully builds and passes a test suite

A large number of watchers, stars and forks on github

NetflixOSS project pull requests

Good code quality and structure

Documentation on how to build and run it

Evidence that code is in use by other projects, or is running in production

# What do you win?

One winner in each of the 10 categories

Ticket and expenses to attend AWS Re:Invent 2013 in Las Vegas

A Trophy

**$10,000 cash and $5,000 in AWS Credits**
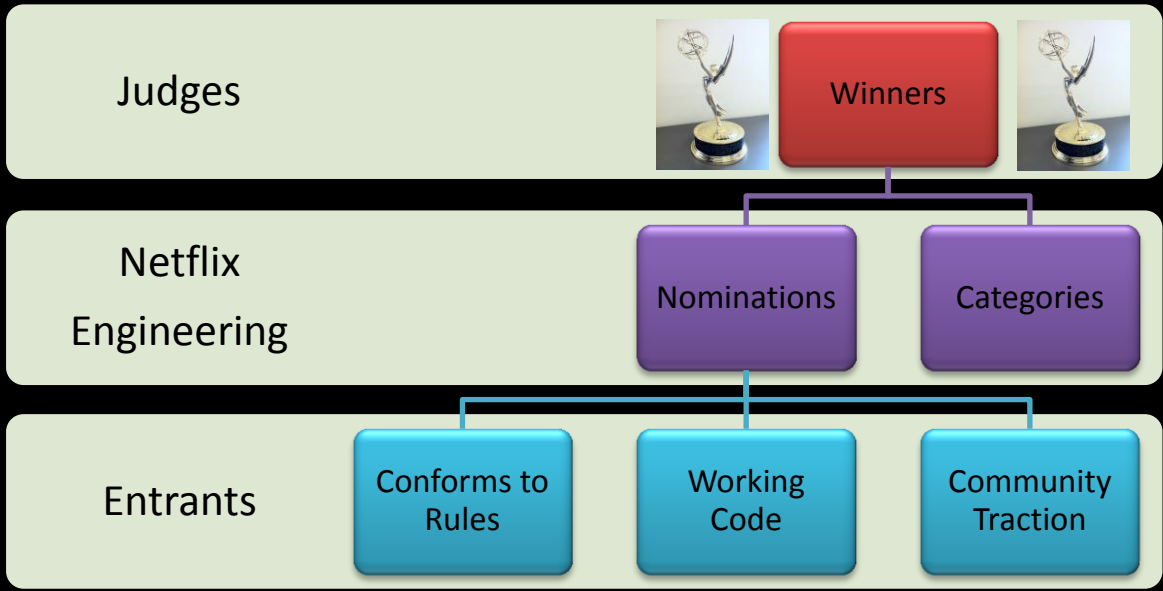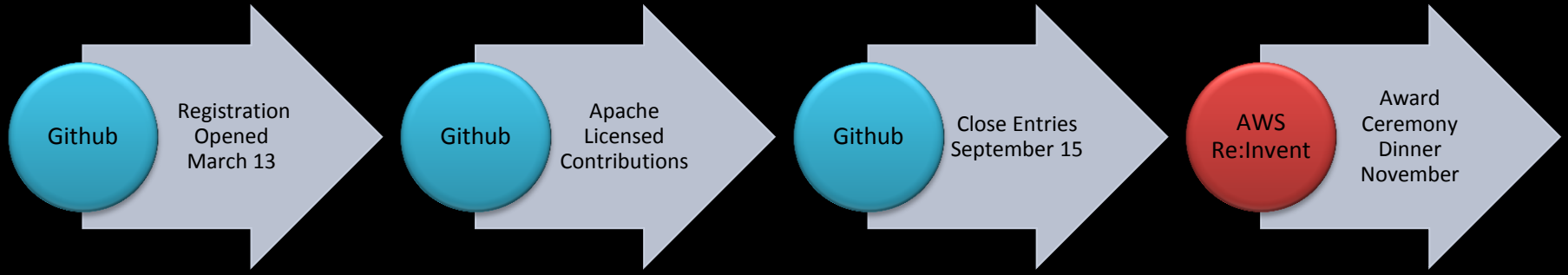
# How do you enter?

Get a (free) github account
Fork github.com/netflix/cloud-prize
Send us your email address
Describe and build your entry

Twitter #cloudprize

Functionality and scale now, portability coming

Moving from parts to a platform in 2013

Netflix is fostering an ecosystem

Rapid Evolution - Low MTBIAMSH

(Mean Time Between Idea And Making Stuff Happen)

# Takeaway

*Netflix is making it easy for everyone to adopt Cloud Native patterns.*

*Open Source is not just the default, it's a strategic weapon.*

http://netflix.github.com

http://techblog.netflix.com

http://slideshare.net/Netflix

http://www.linkedin.com/in/adriancockcroft

@adrianco #netflixcloud @NetflixOSS