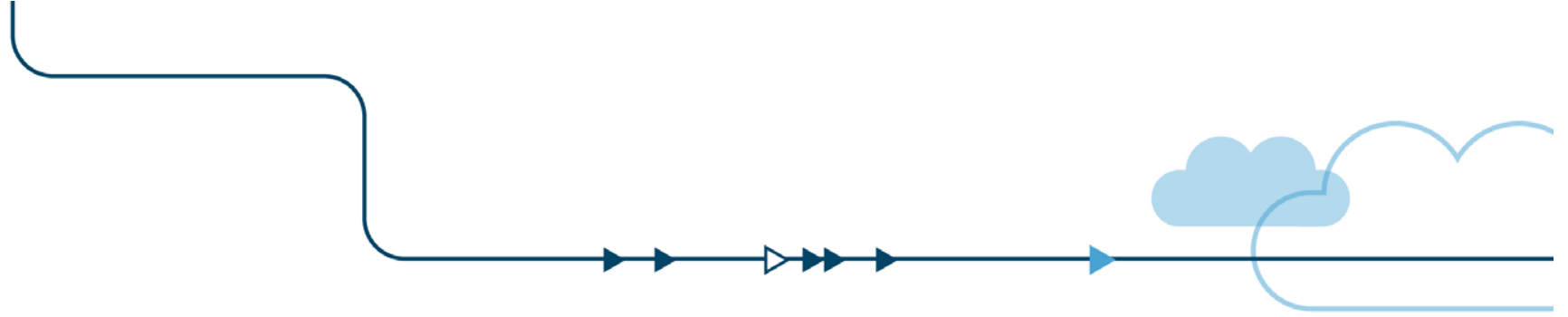




Cloud Native Cost Optimization

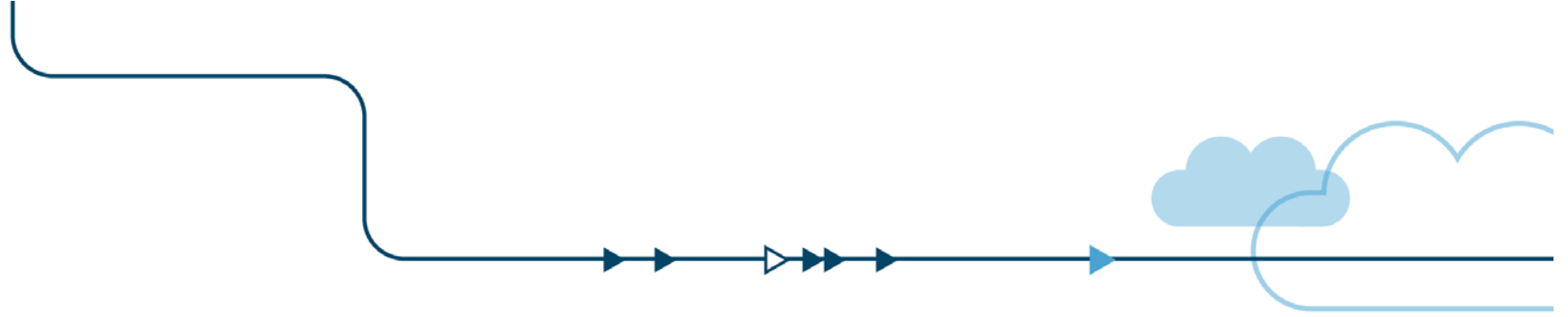
Adrian Cockcroft @adrianco
Technology Fellow - Battery Ventures
ICPE - Austin, February 2015



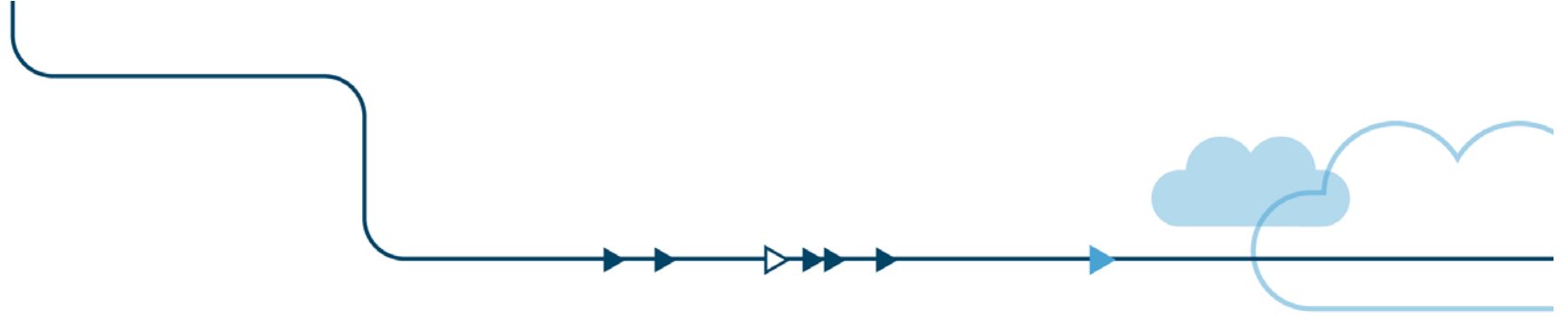


Why Does Performance Matter?





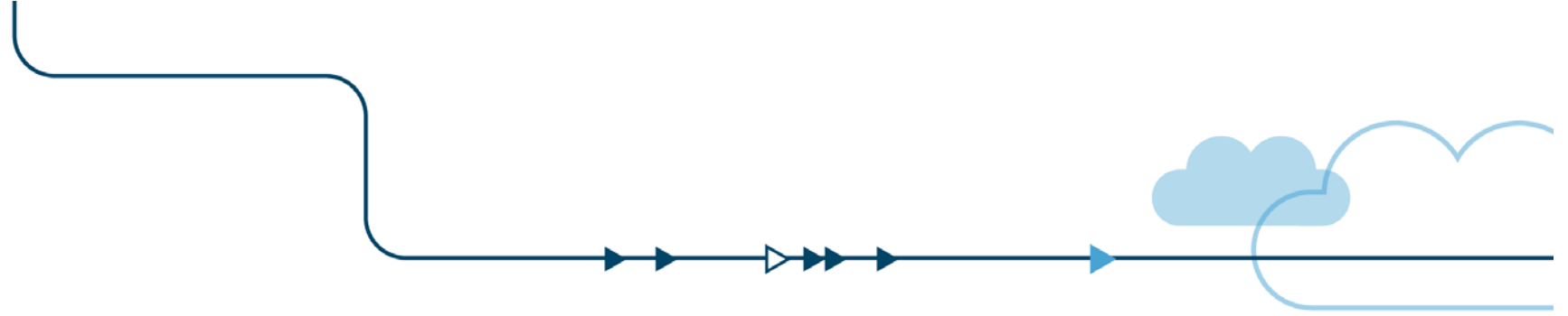
Latency Efficiency



Users: Response Latency
Developers: Release Latency
Operators: Efficiency

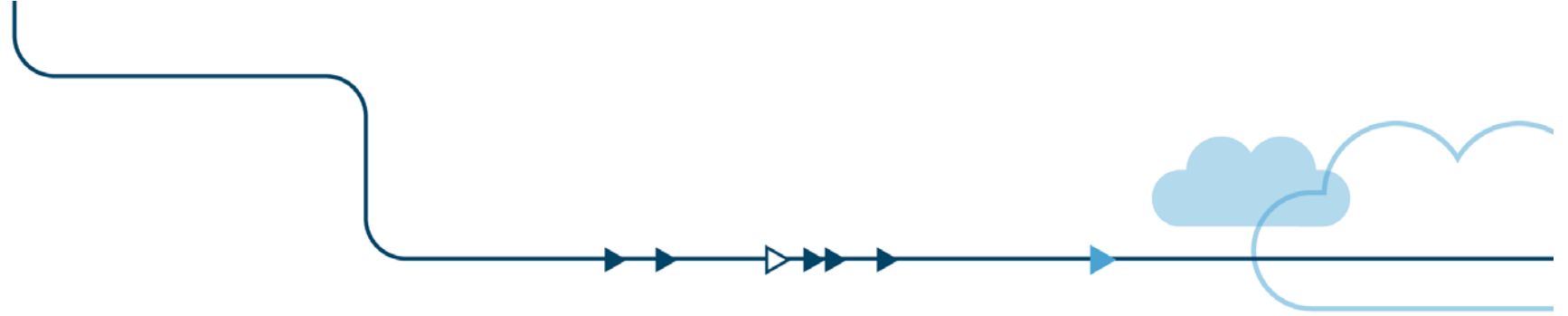
@adrianco

BV
Battery Ventures



Less Time
Less Cost





Faster Delivery

See talks by @adrianco
Speed and Scale - QCon New York
Fast Delivery - GOTO Copenhagen

A decorative line art graphic at the bottom left of the slide. It features a dark blue line that starts from the left, goes down, then right, and then continues as a horizontal line with several small arrows pointing right. This line ends with a blue arrow pointing towards a light blue cloud. From the cloud, a wavy line extends to the right, and another line loops back from the right side towards the cloud.

@adrianco

BV
Battery Ventures



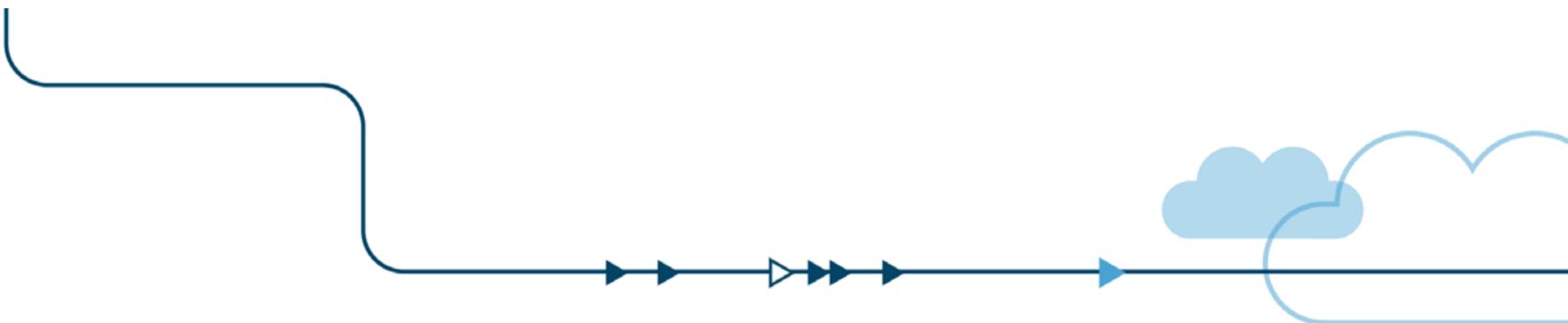
Cheaper

*This talk:
How to use Cloud Native architecture to
reduce cost without slowing down releases*



@adrianco

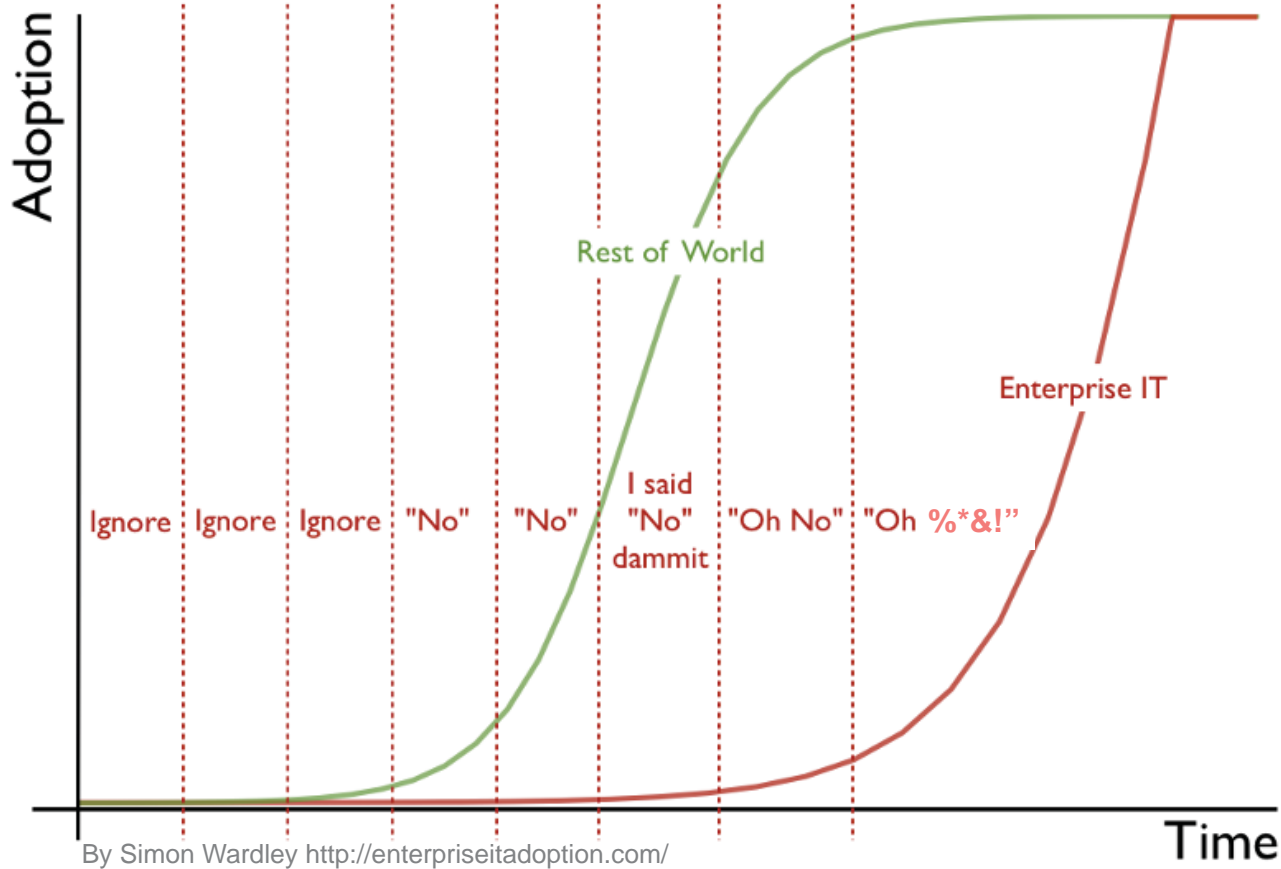
BV
Battery Ventures



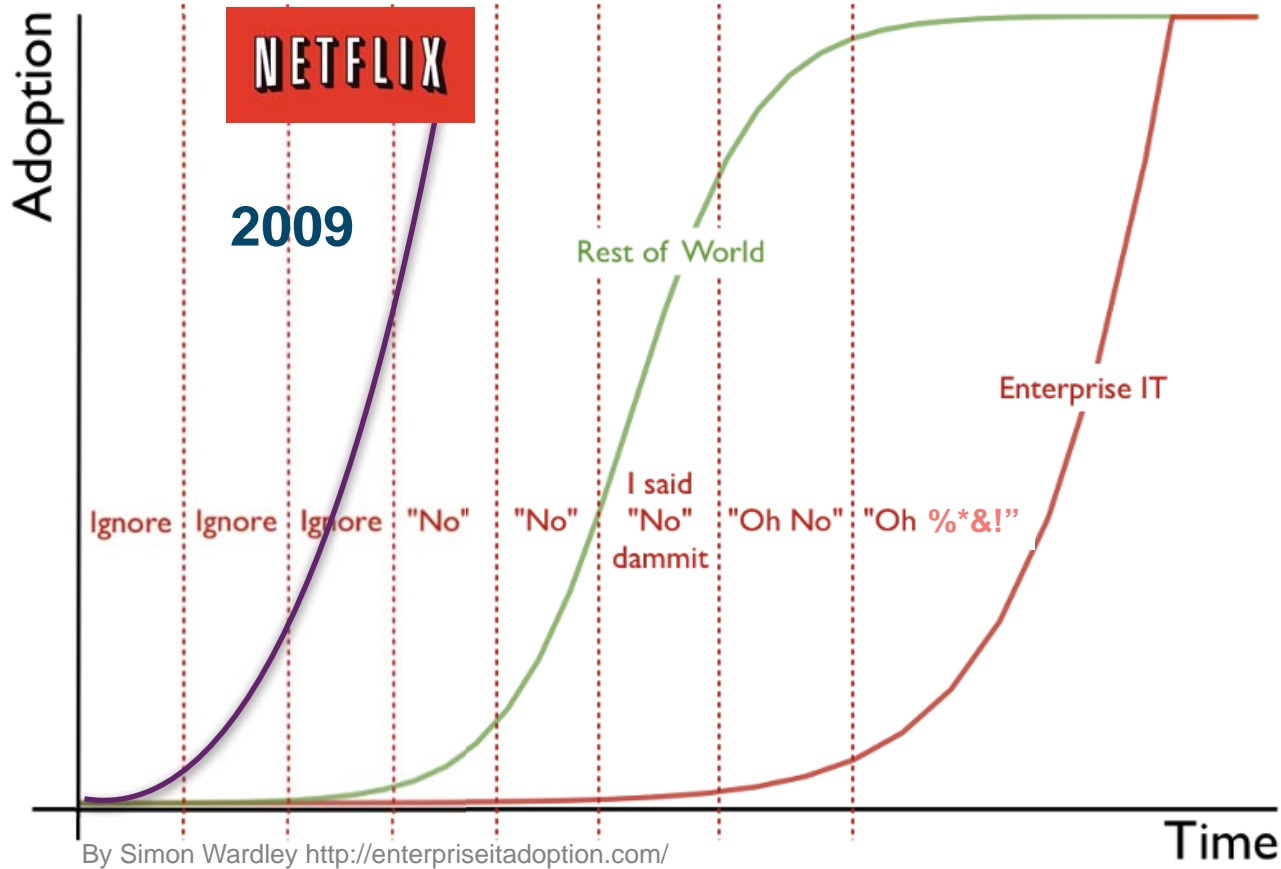
Speeding up Development
Cloud Native Applications
Cost Optimization



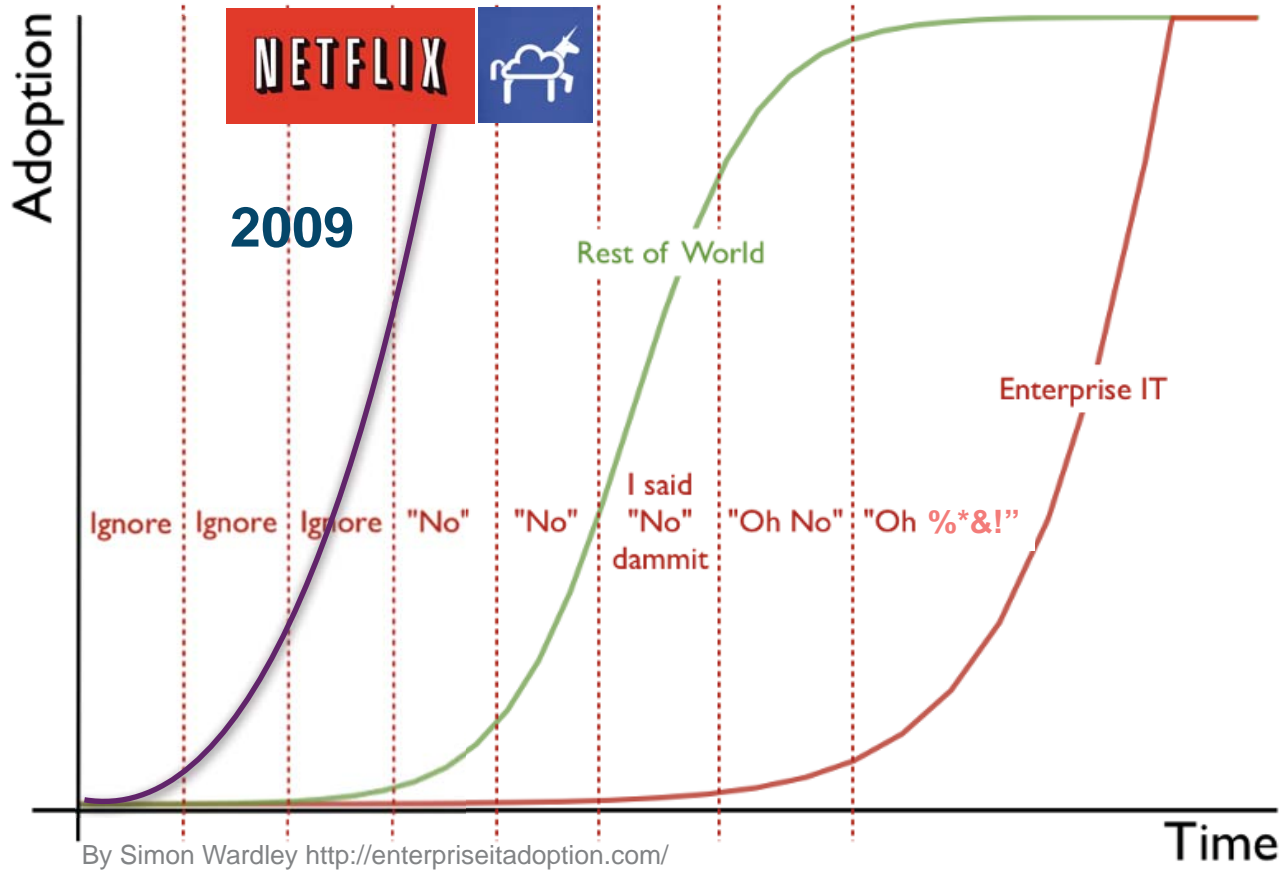
Why am I here?



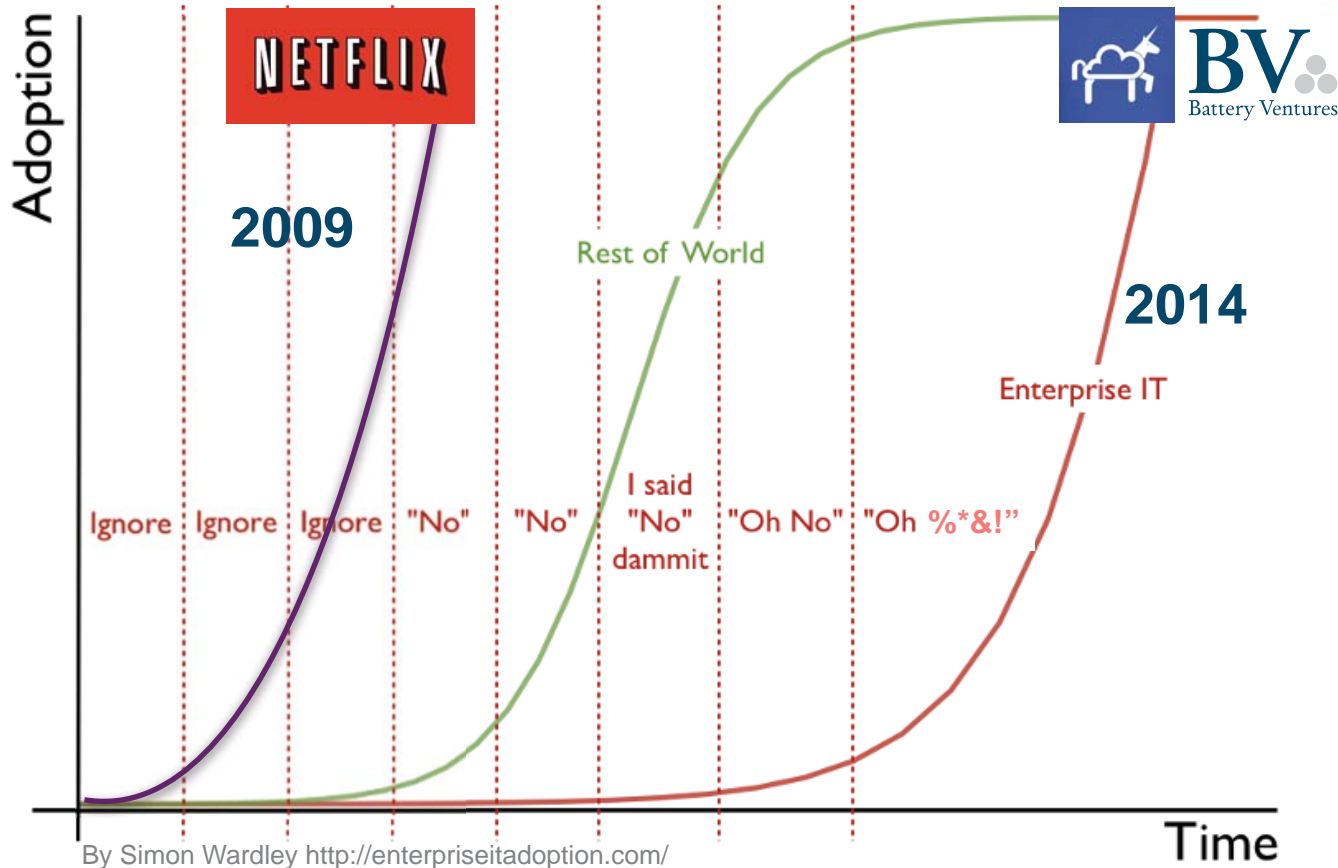
Why am I here?



Why am I here?

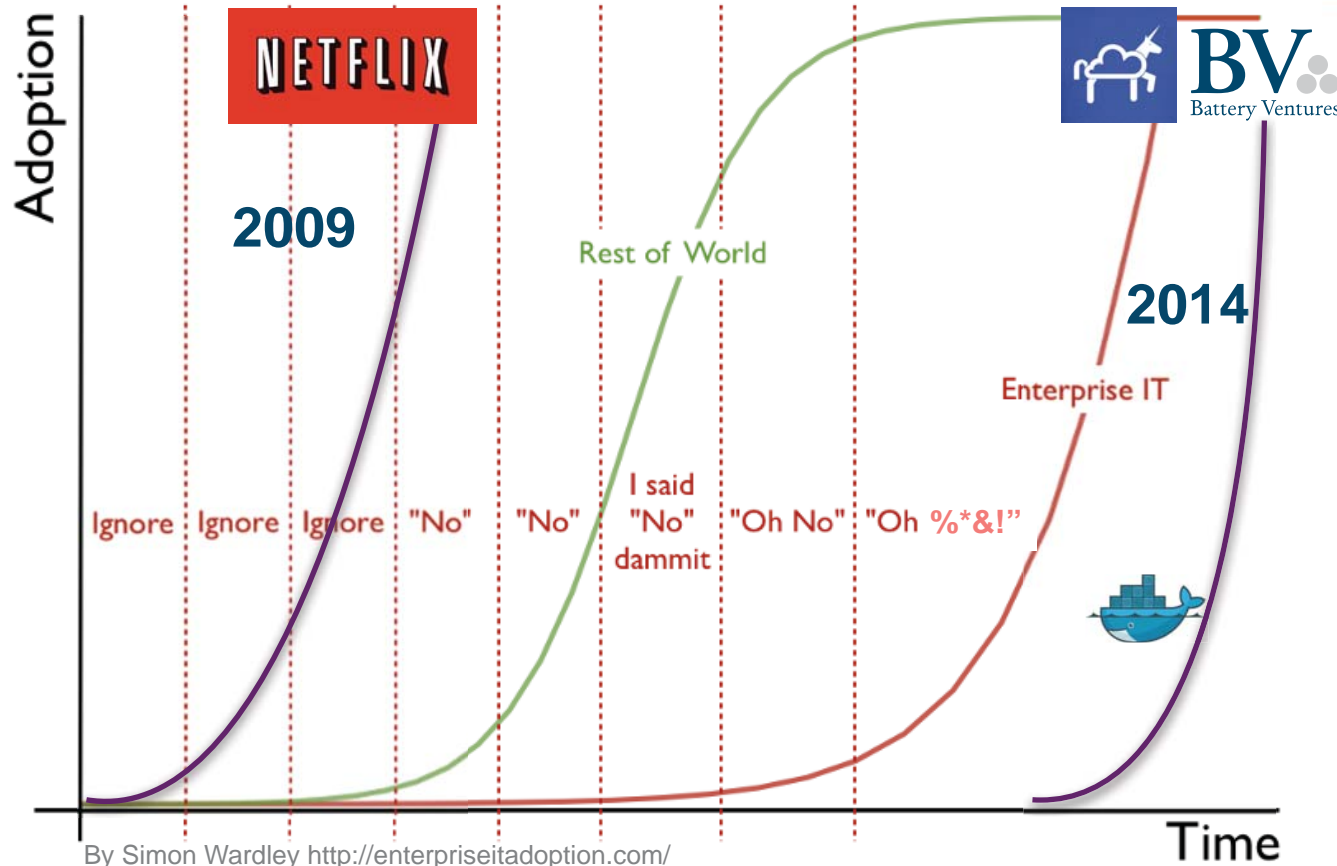


Why am I here?



@adrianco's job at the intersection of cloud and Enterprise IT, looking for disruption and opportunities.

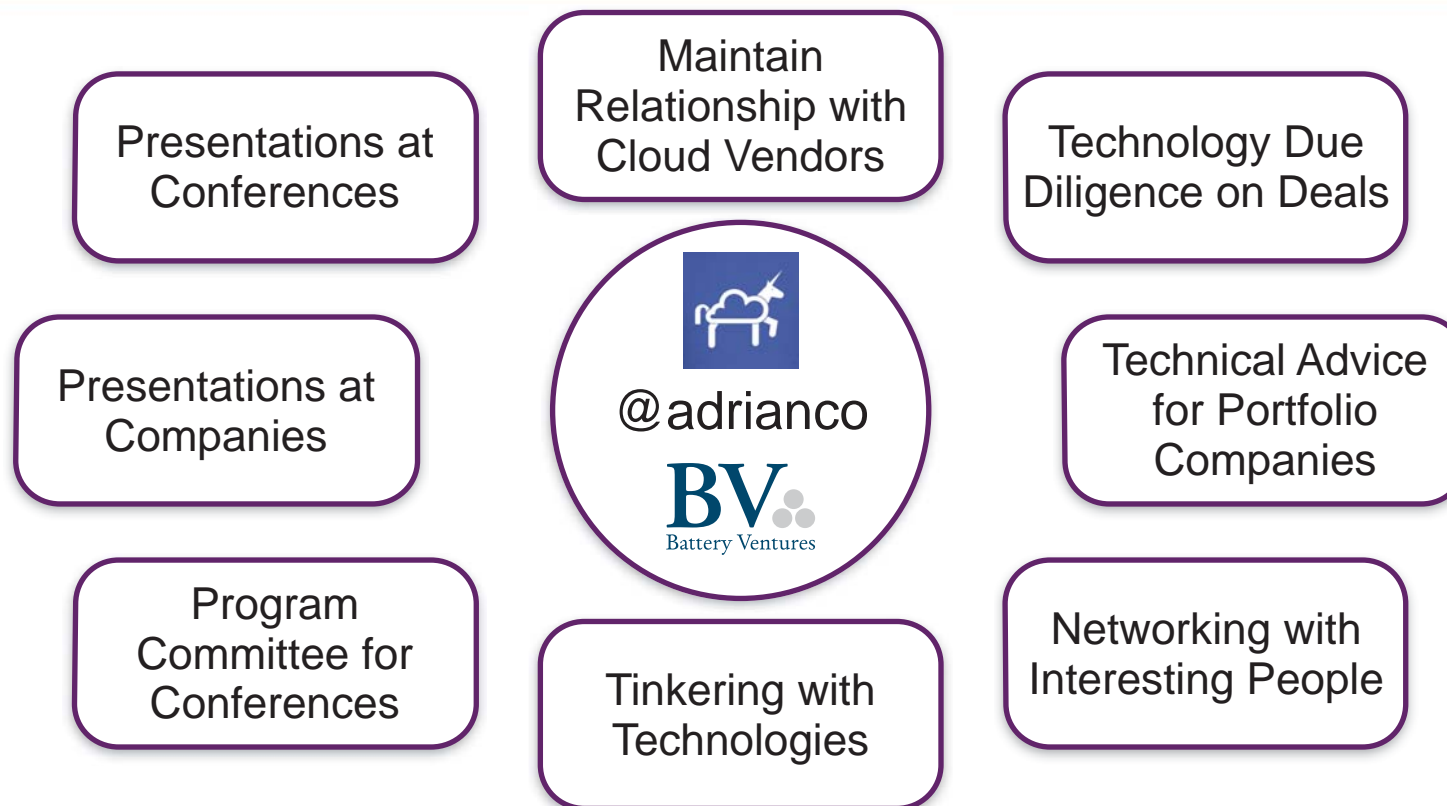
Why am I here?



@adrianco's job at the intersection of cloud and Enterprise IT, looking for disruption and opportunities.

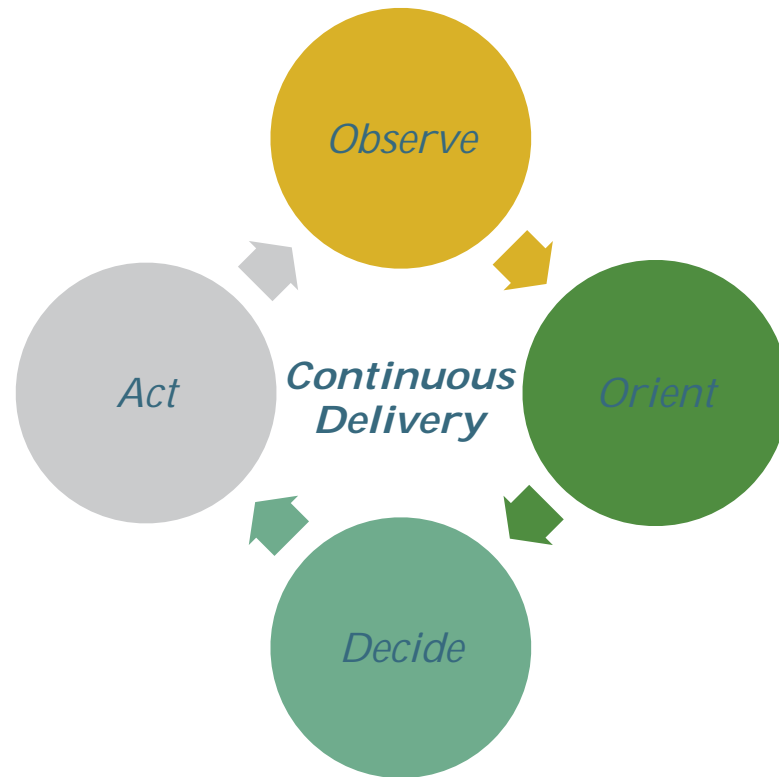
Example: Docker wasn't on anyone's roadmap for 2014. It's on everyone's roadmap for 2015.

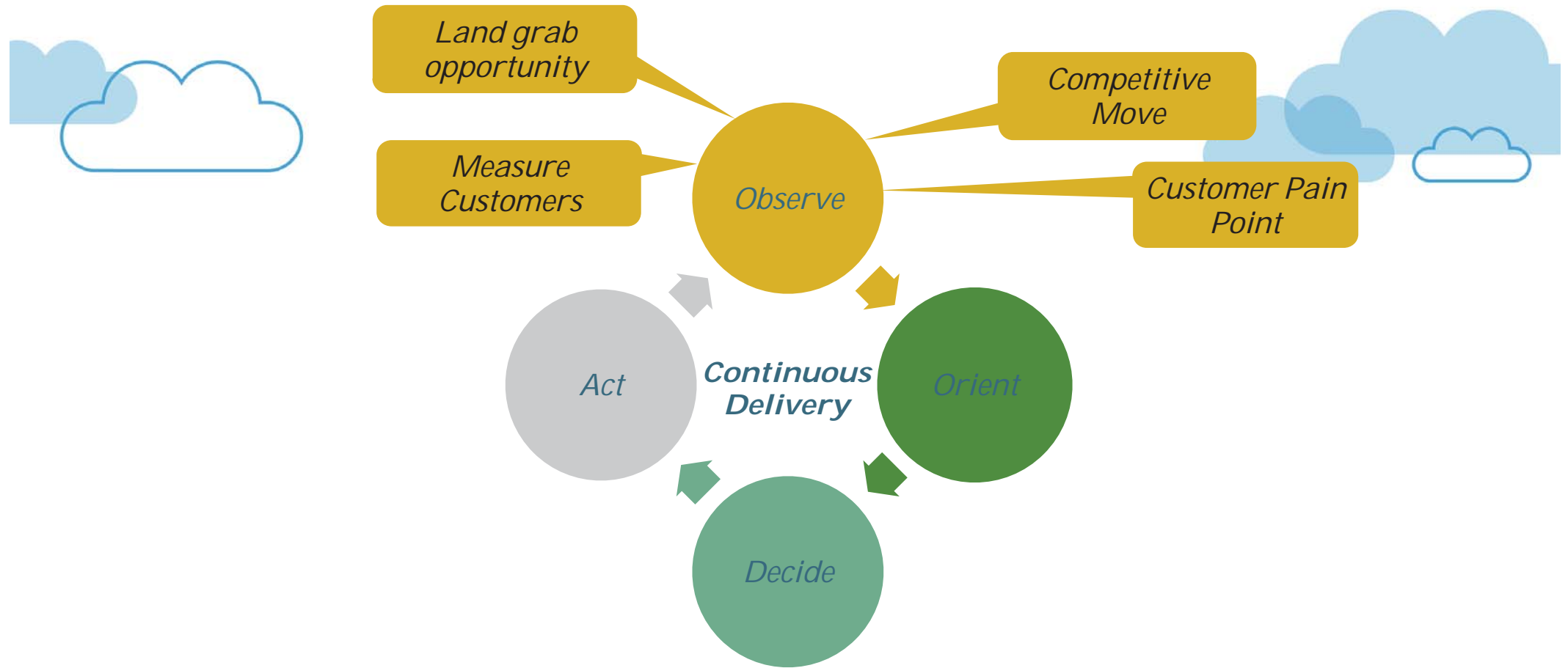
What does @adrianco do?



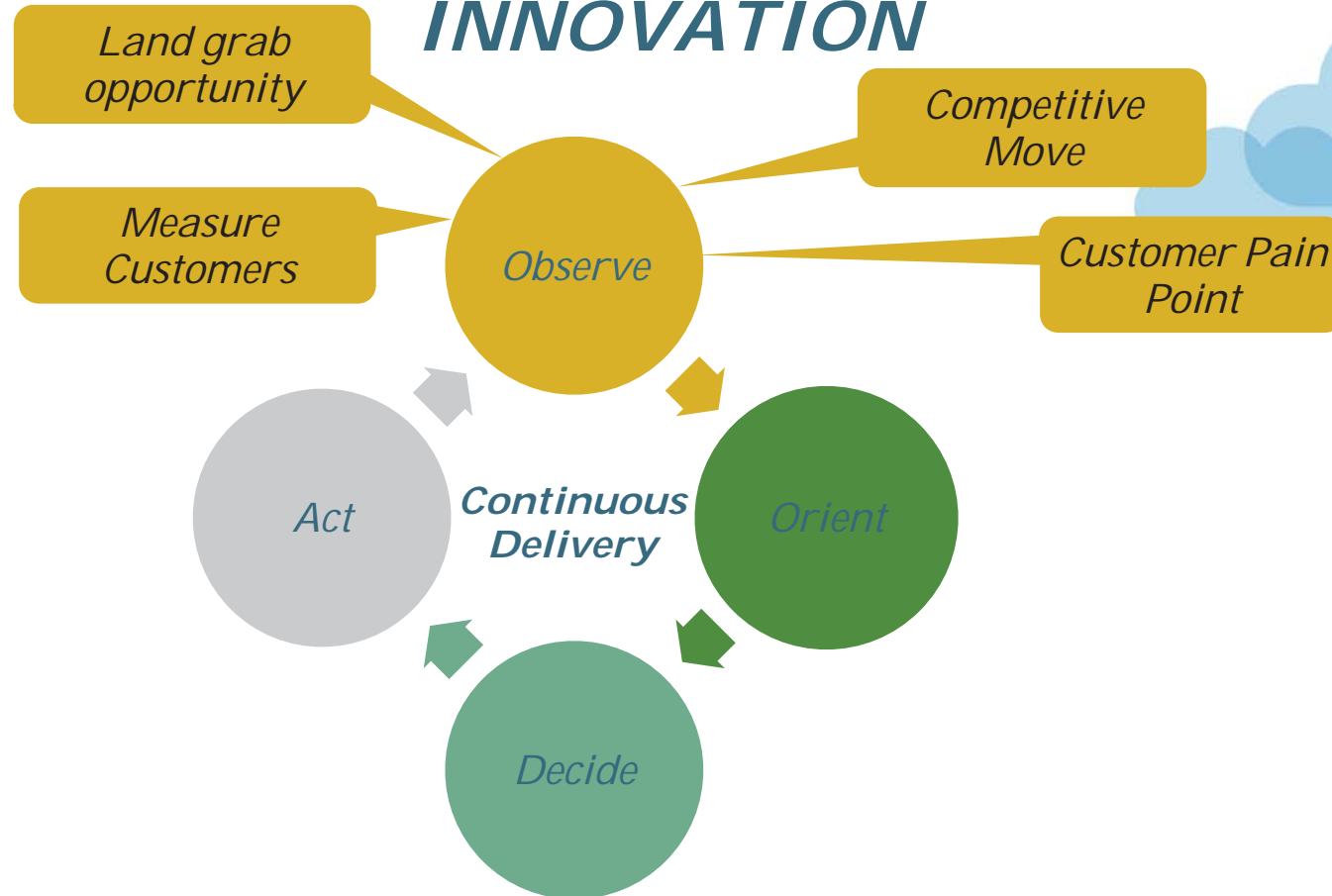


Speeding Up Development

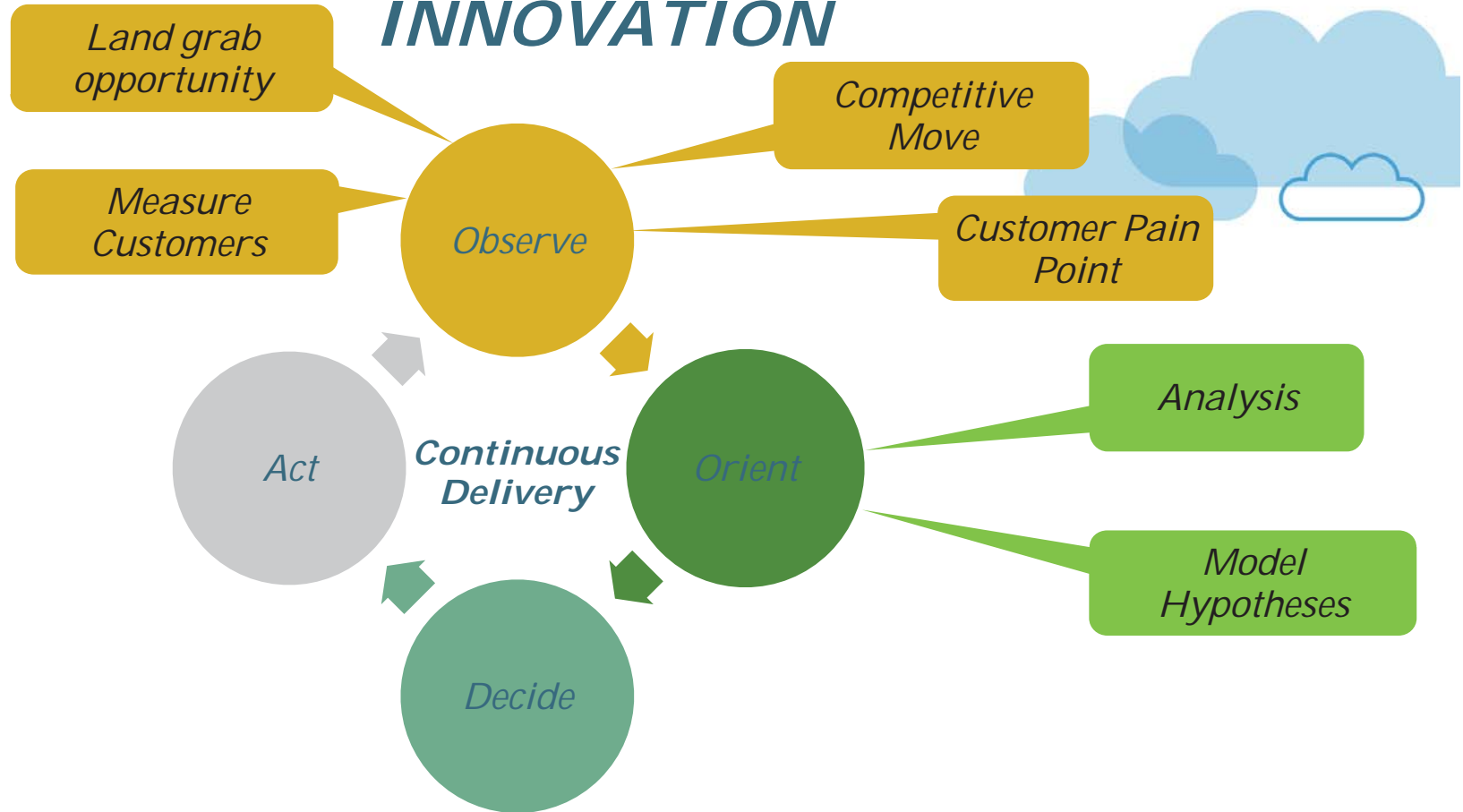


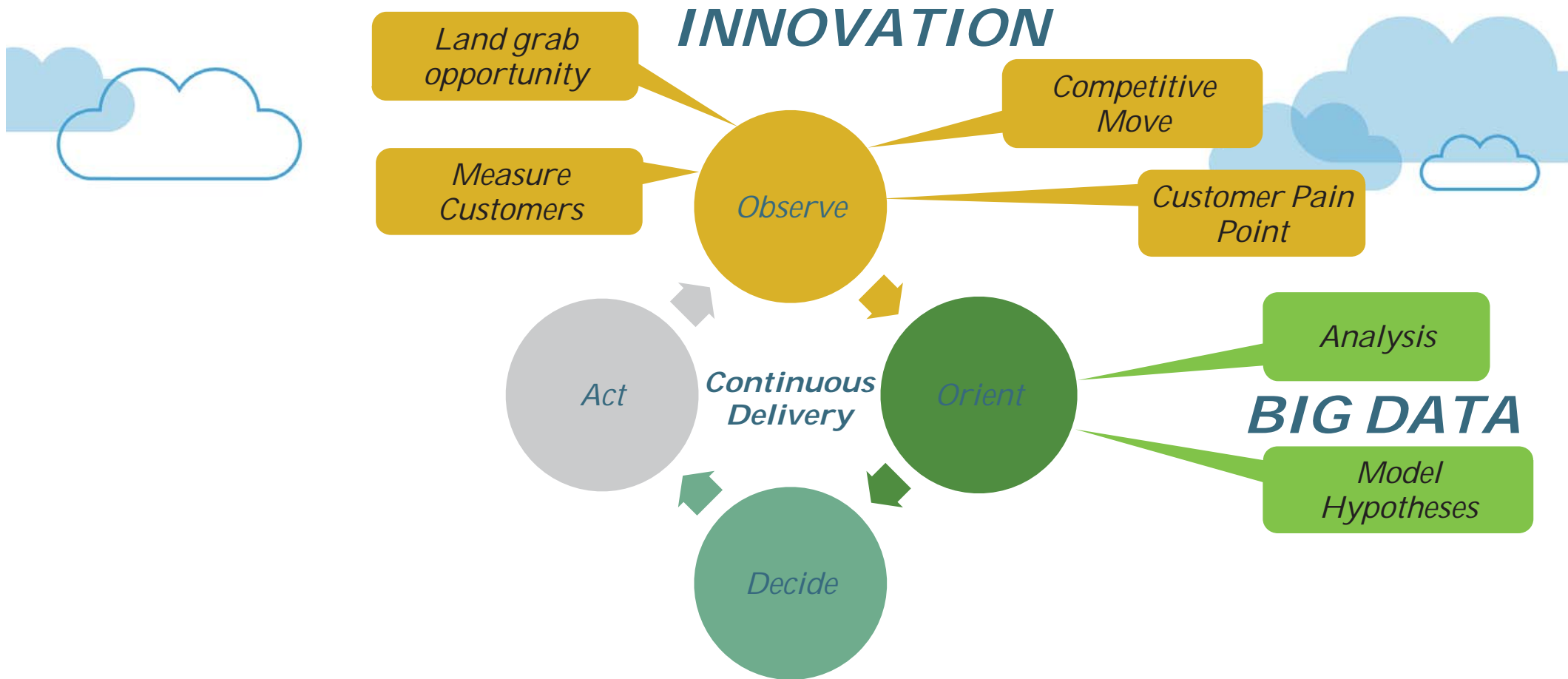


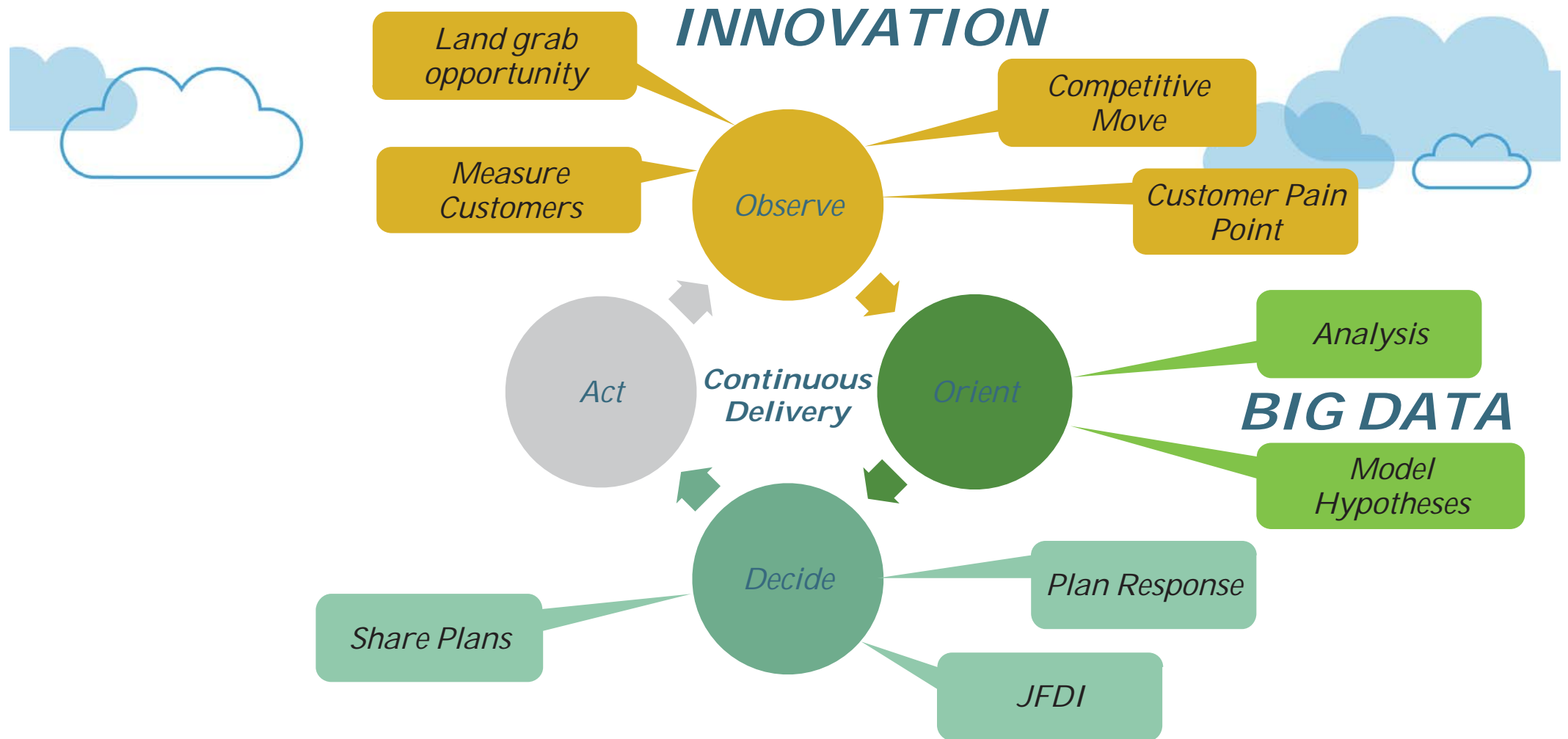
INNOVATION

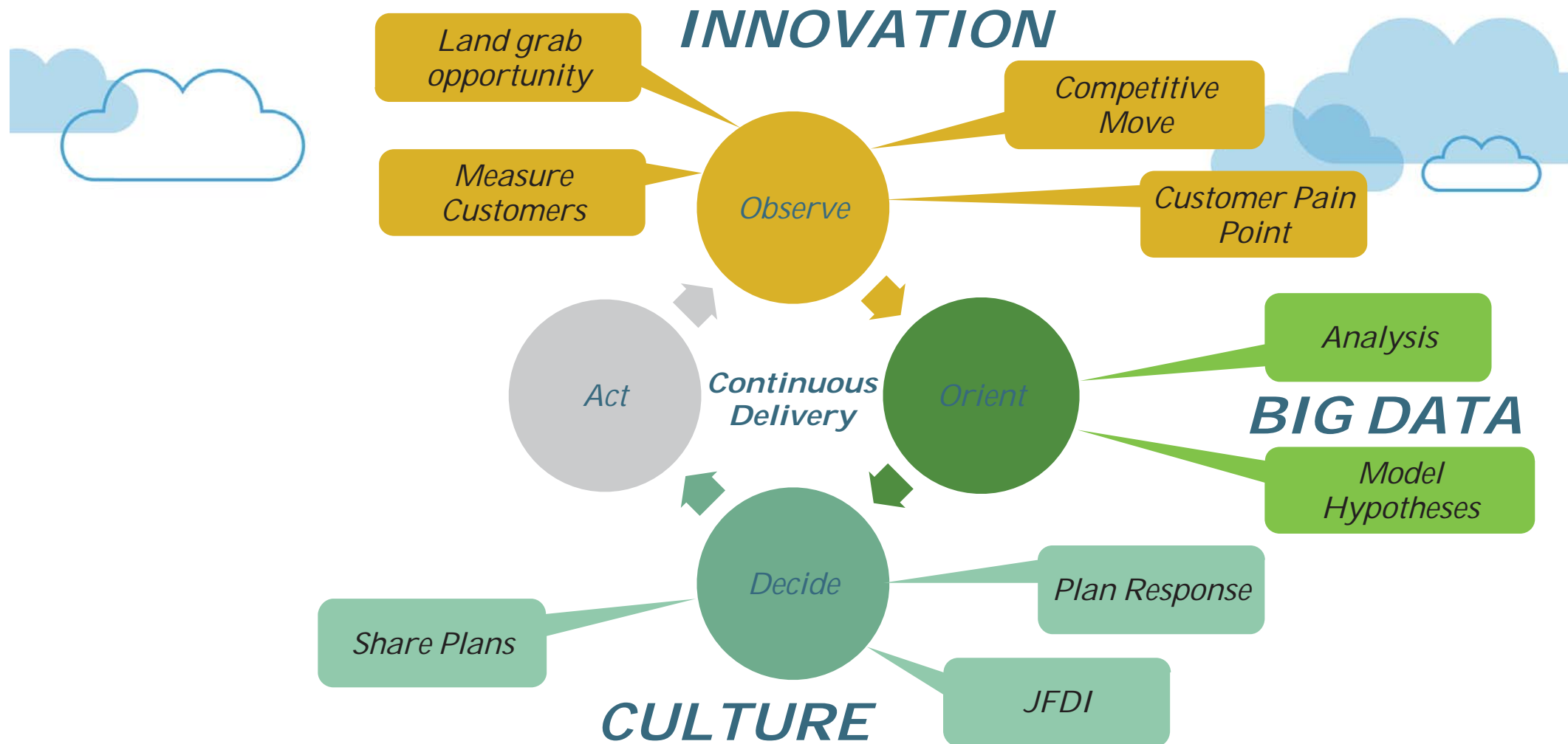


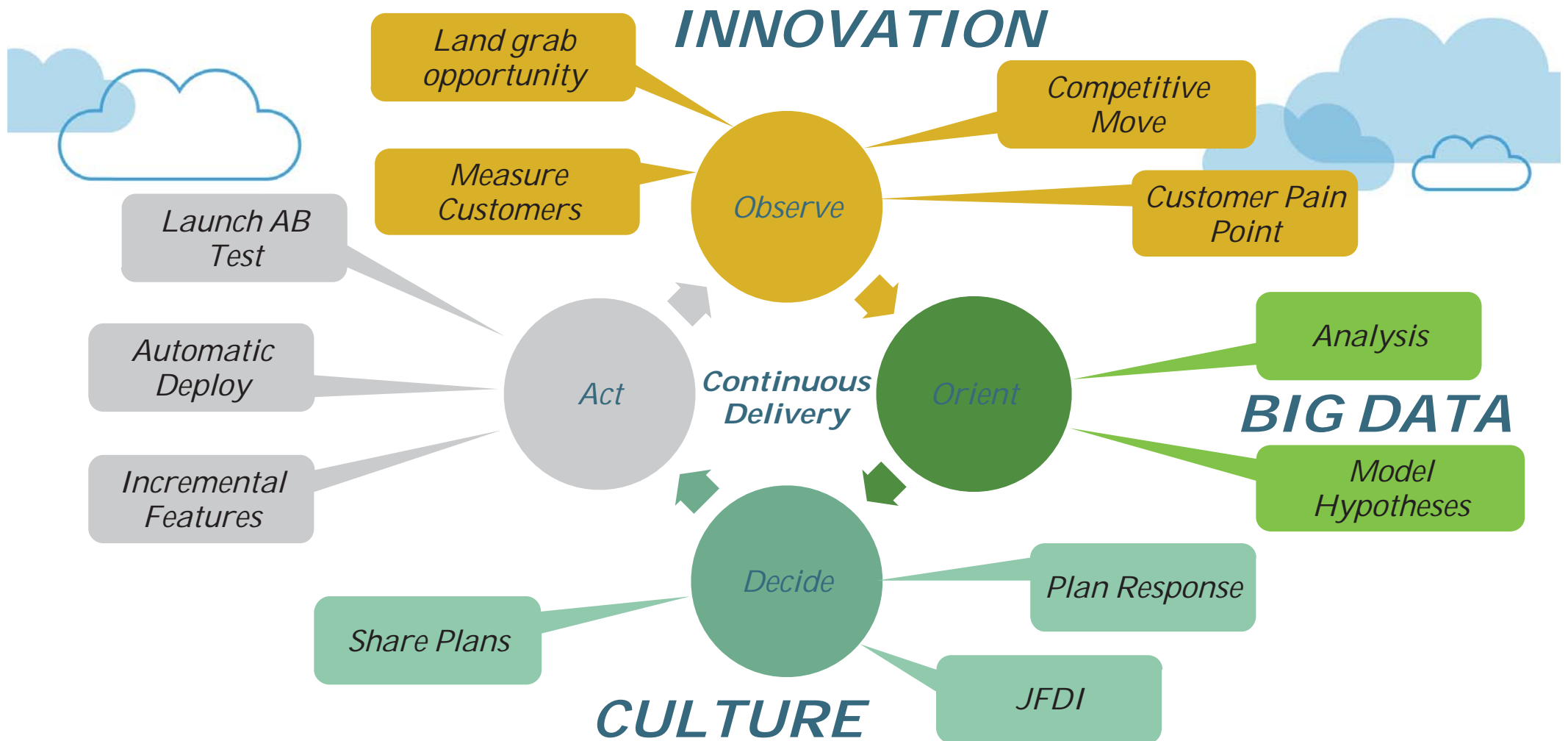
INNOVATION

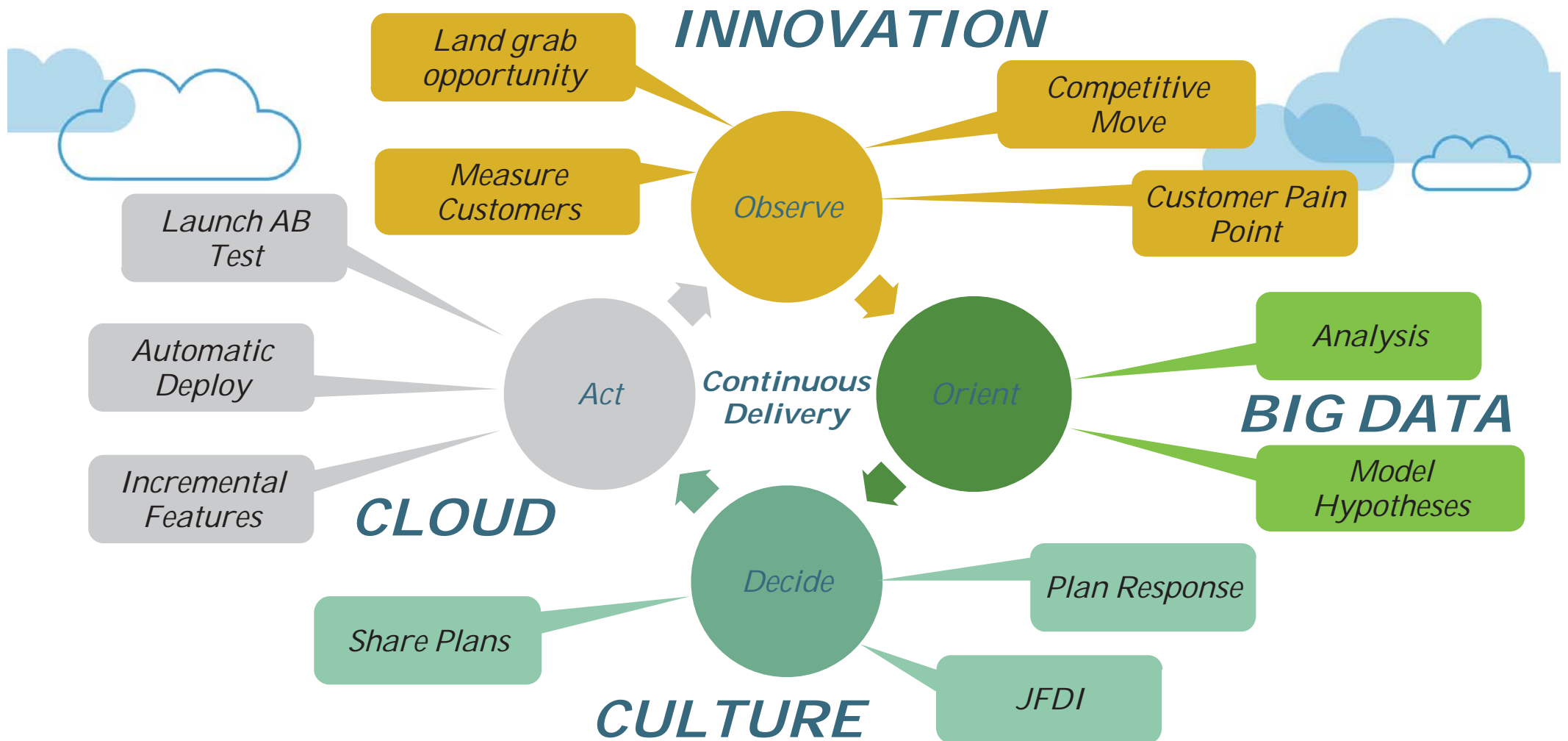


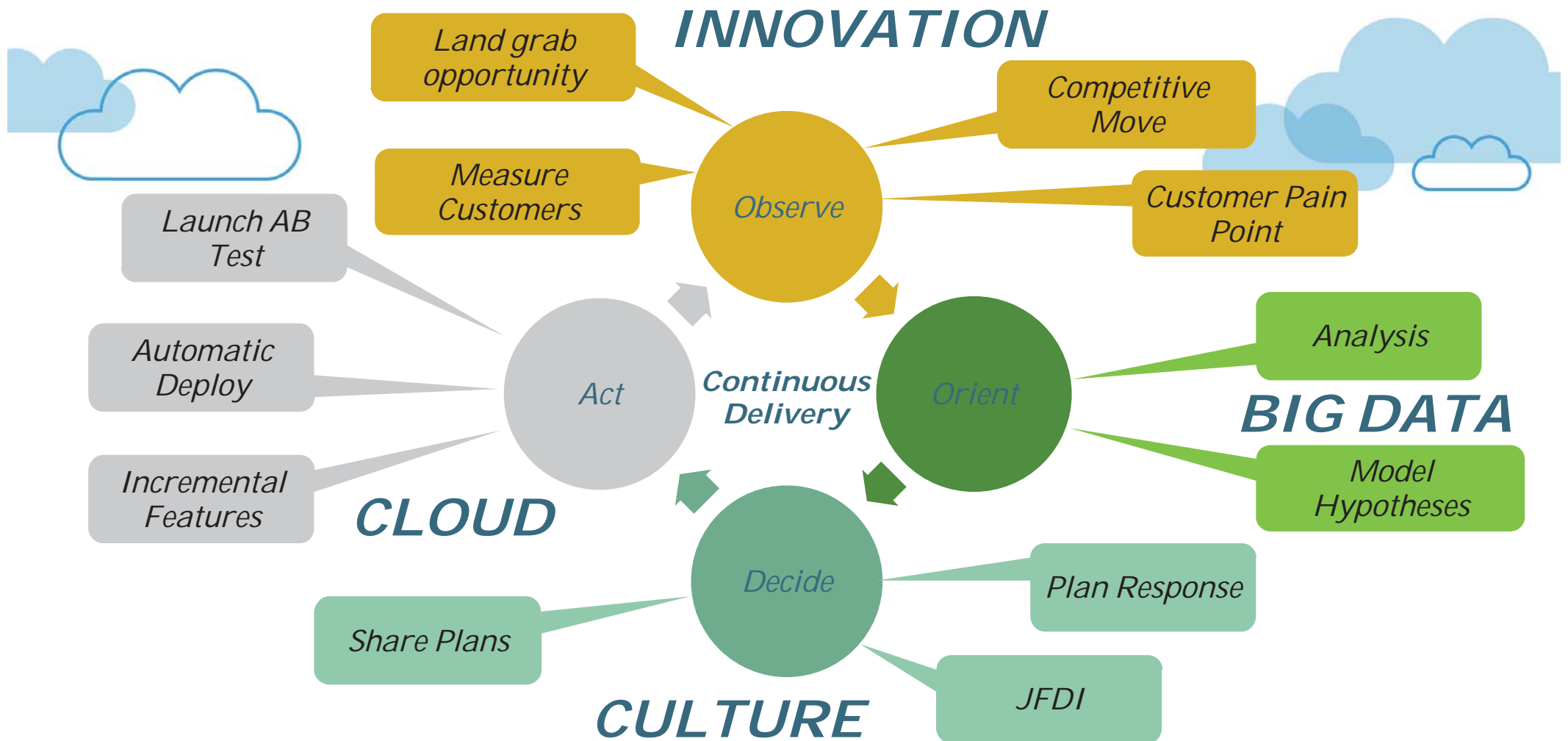


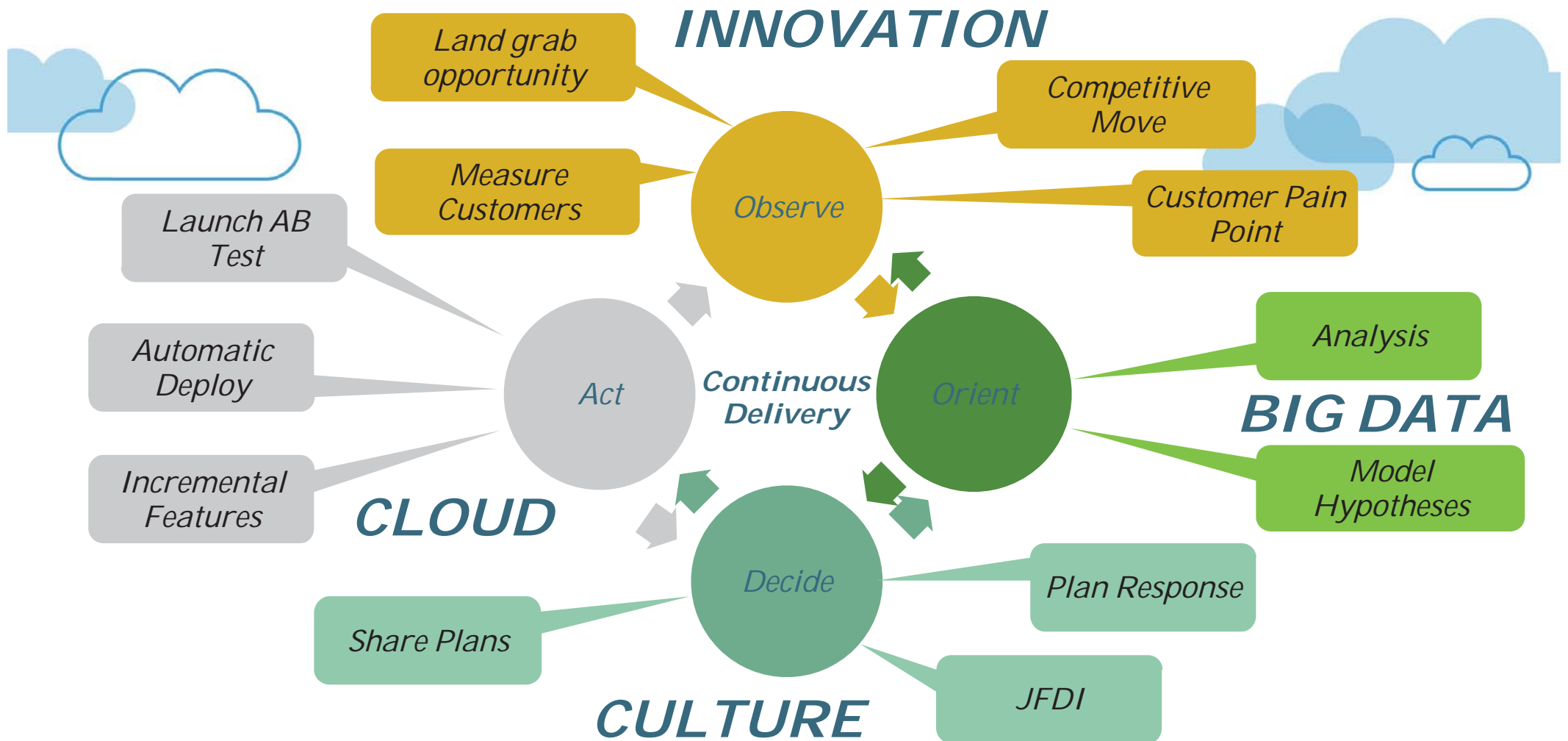


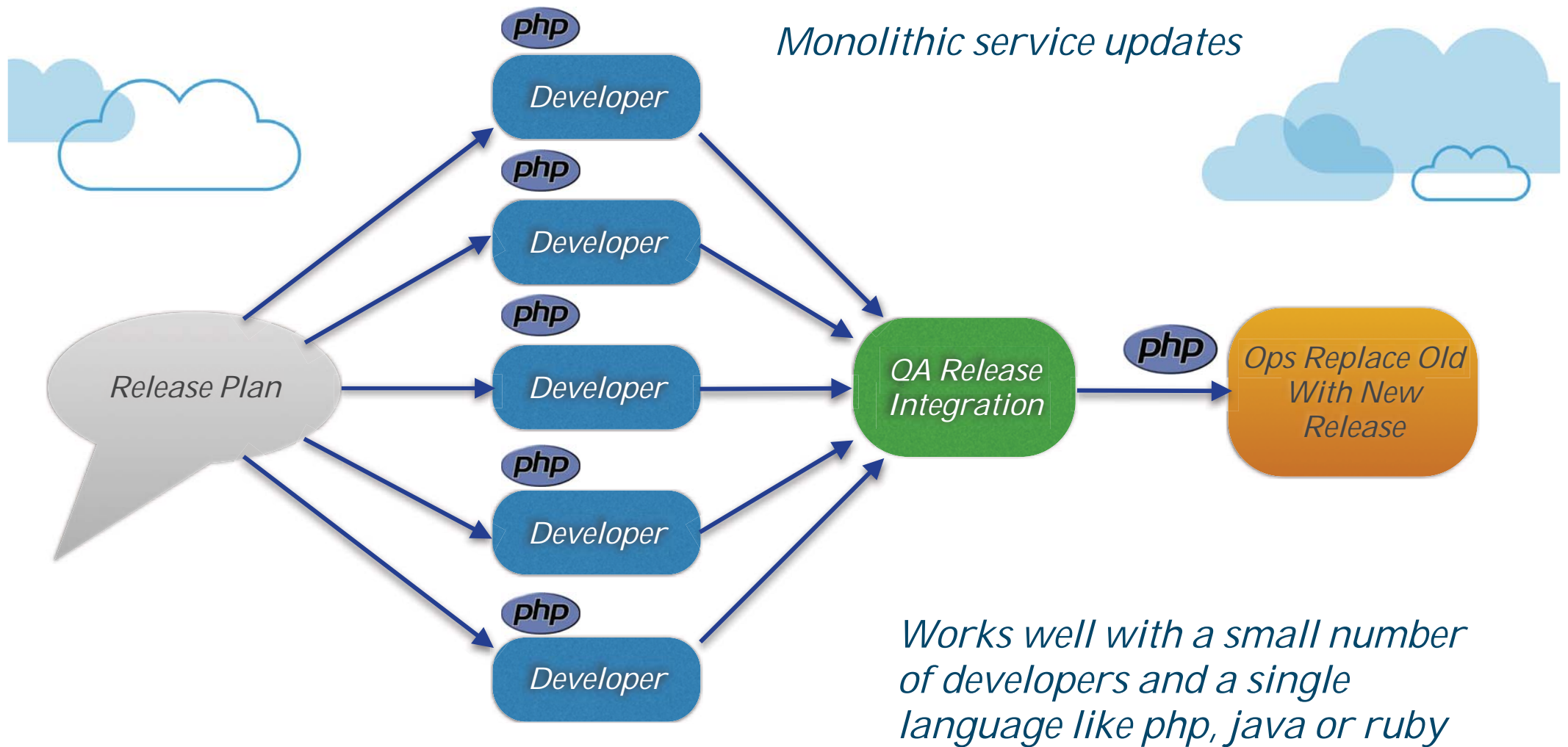


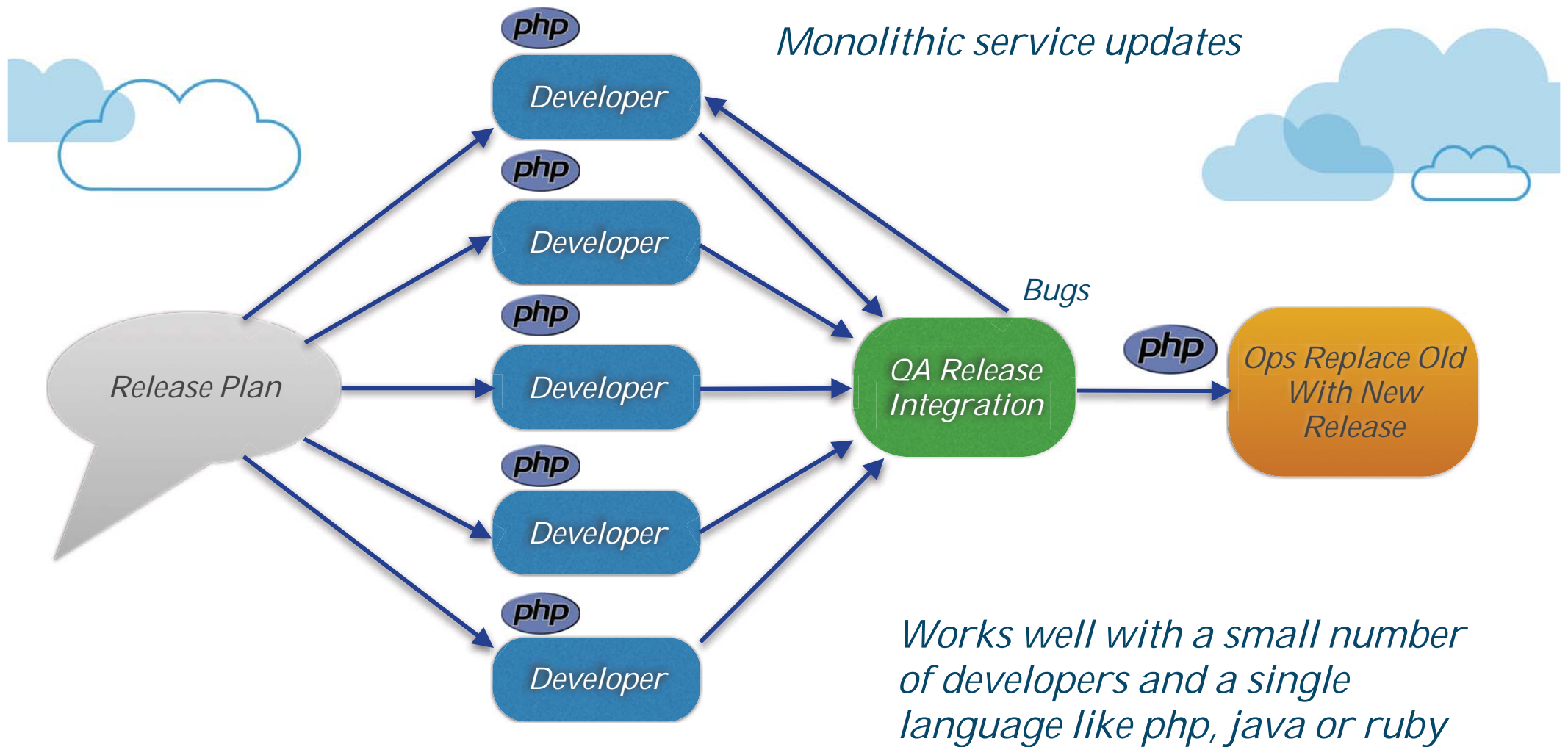


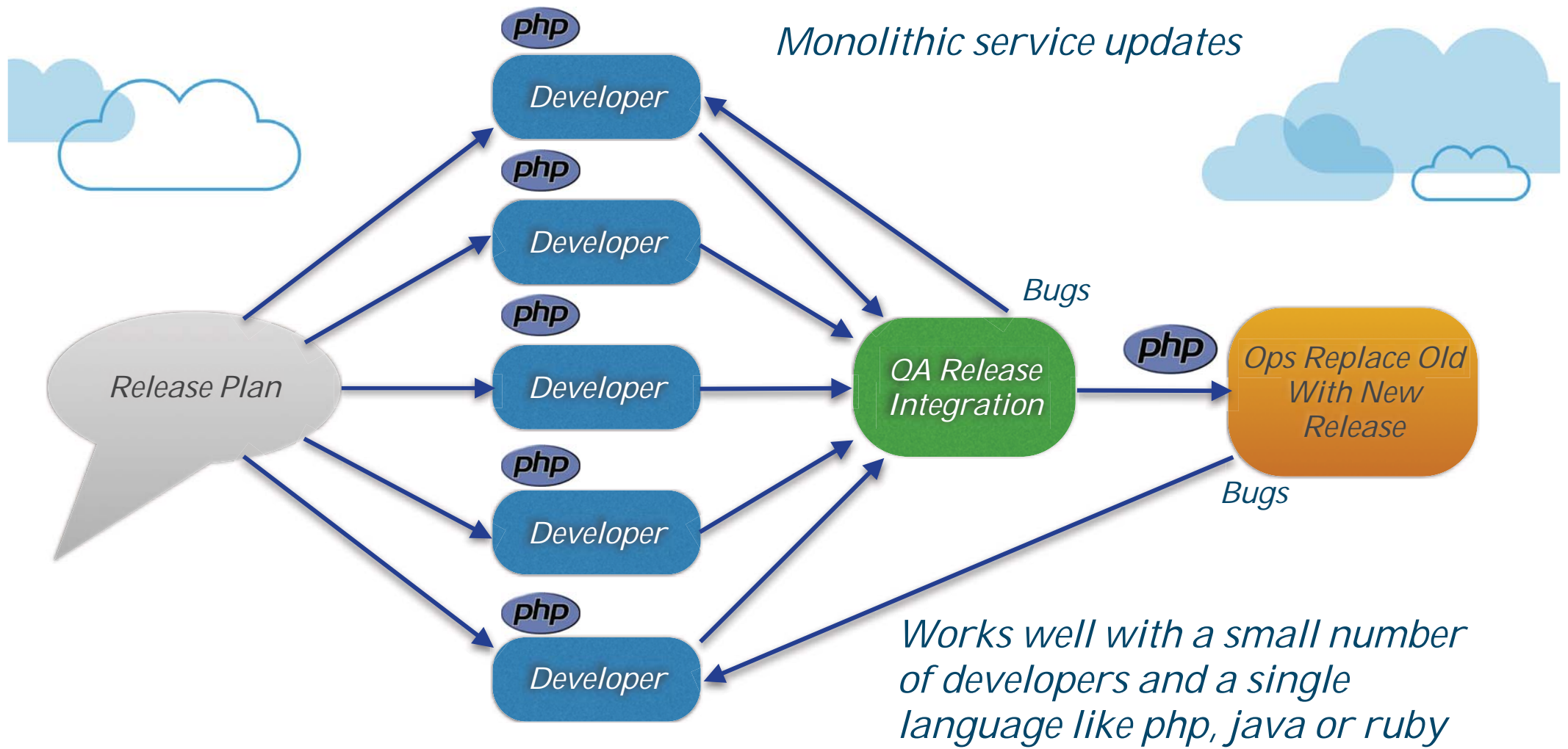












Breaking Down the SILOs



Breaking Down the SILOs



*Prod
Mgr*

UX

Dev

QA

DBA

*Sys
Adm*

*Net
Adm*

*SAN
Adm*

@adrianco

BV
Battery Ventures

Breaking Down the SILOs



Product Team Using Monolithic Delivery

Product Team Using Monolithic Delivery

*Prod
Mgr*

UX

Dev

QA

DBA

*Sys
Adm*

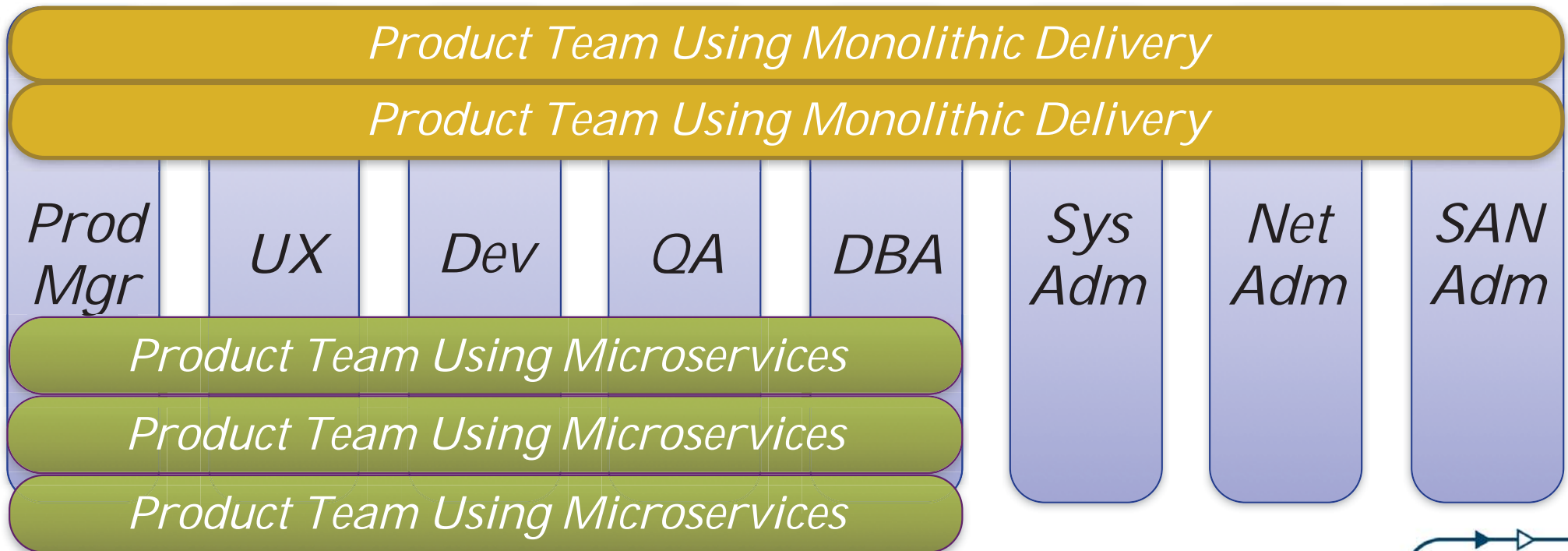
*Net
Adm*

*SAN
Adm*

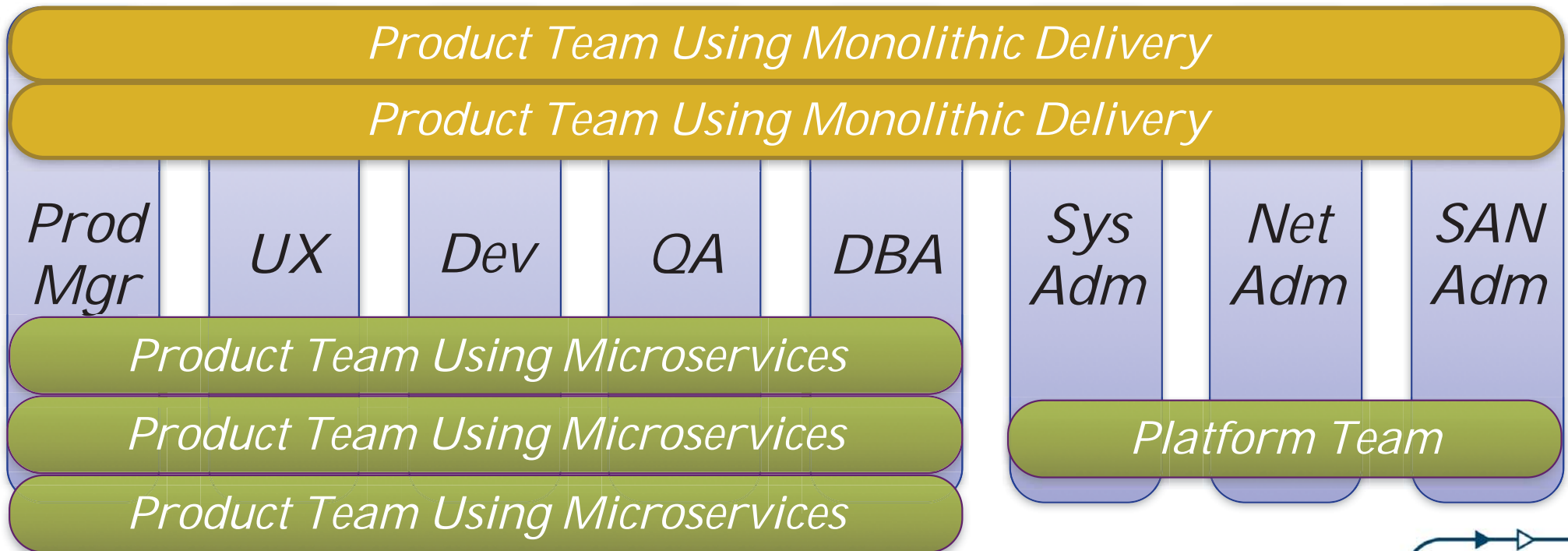
@adrianco

BV
Battery Ventures

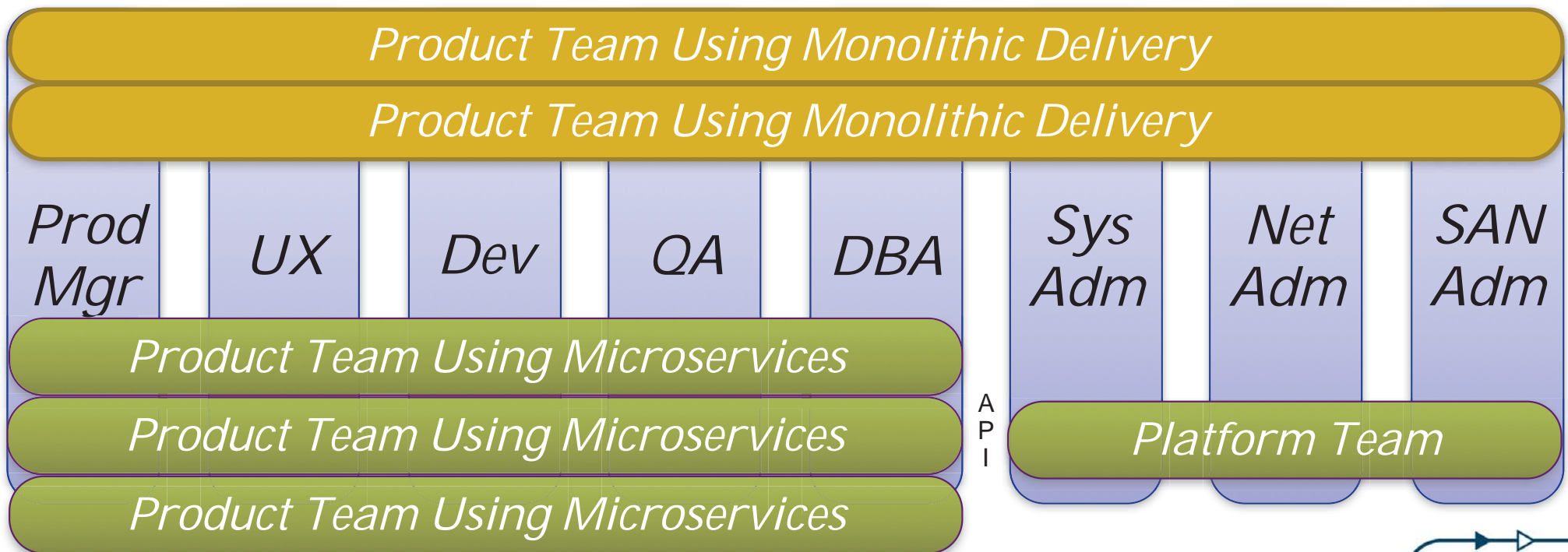
Breaking Down the SILOs



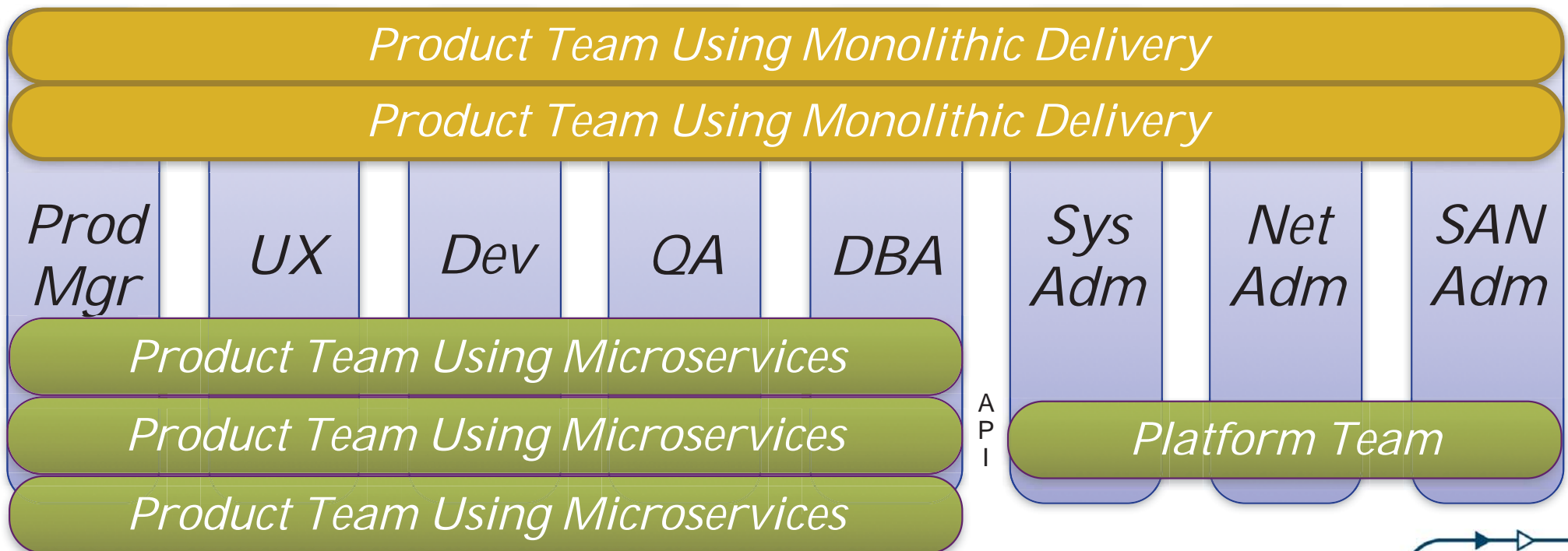
Breaking Down the SILOs



Breaking Down the SILOs

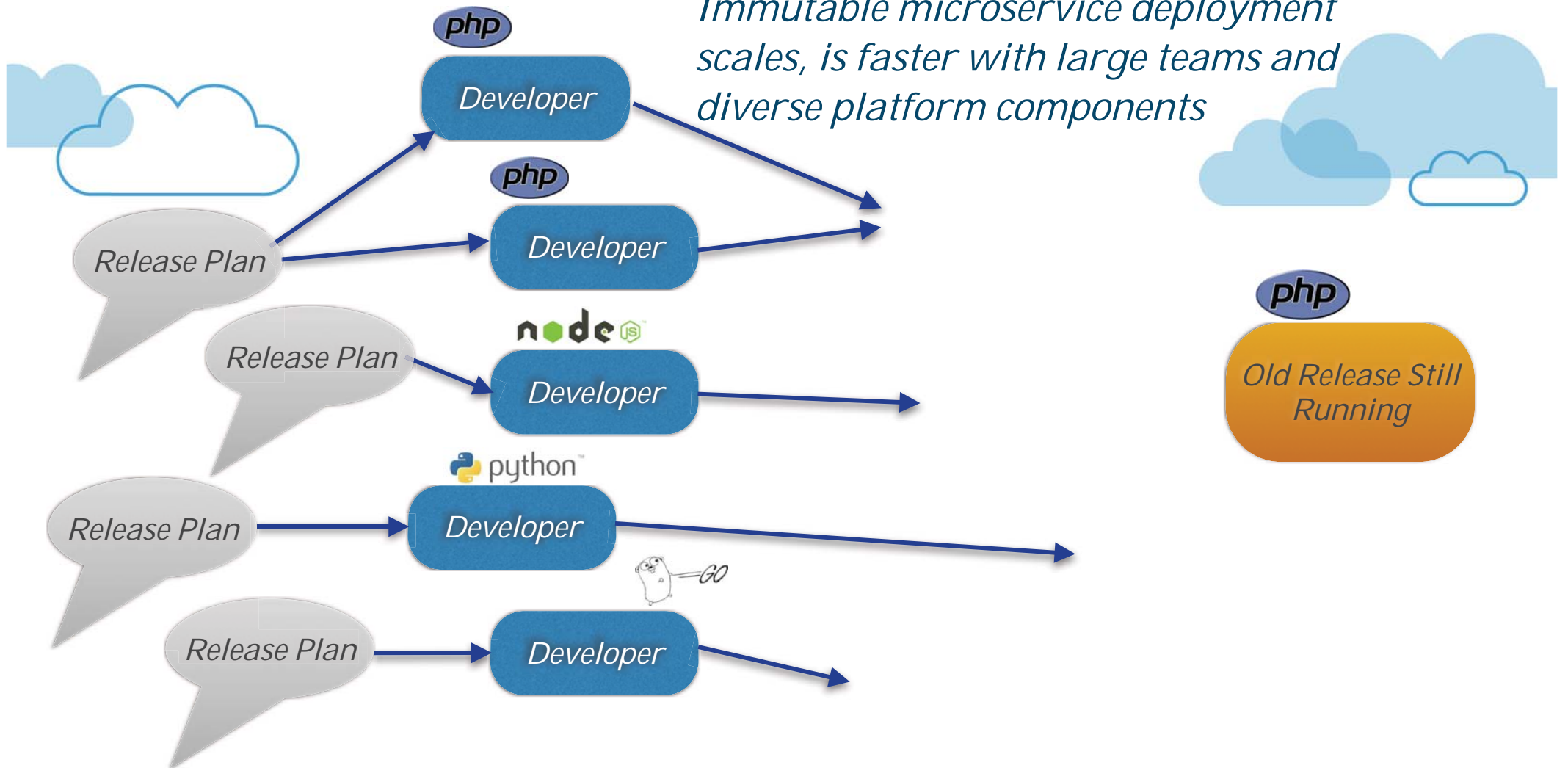


Breaking Down the SILOs

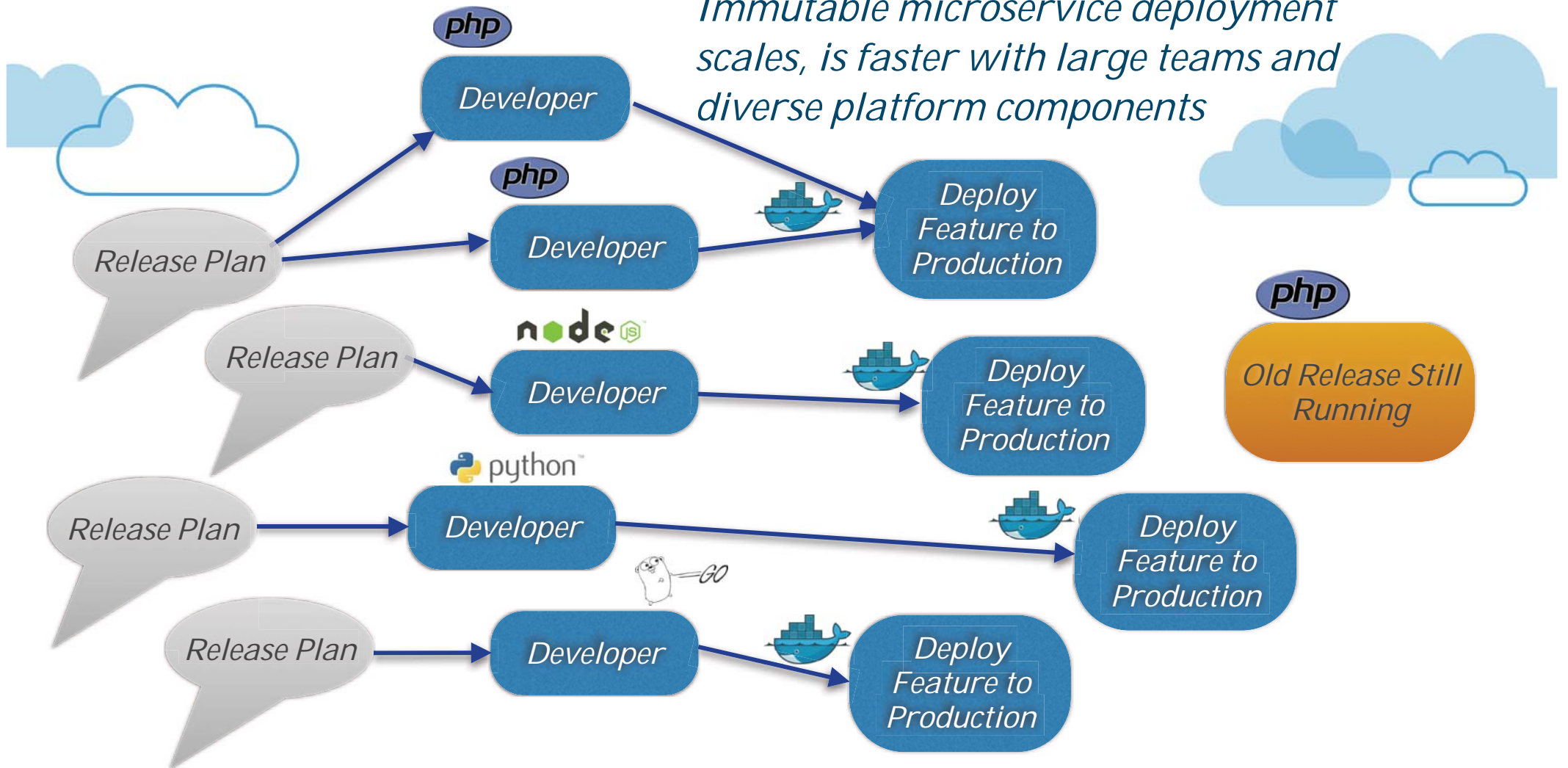


DevOps is a Re-Org!

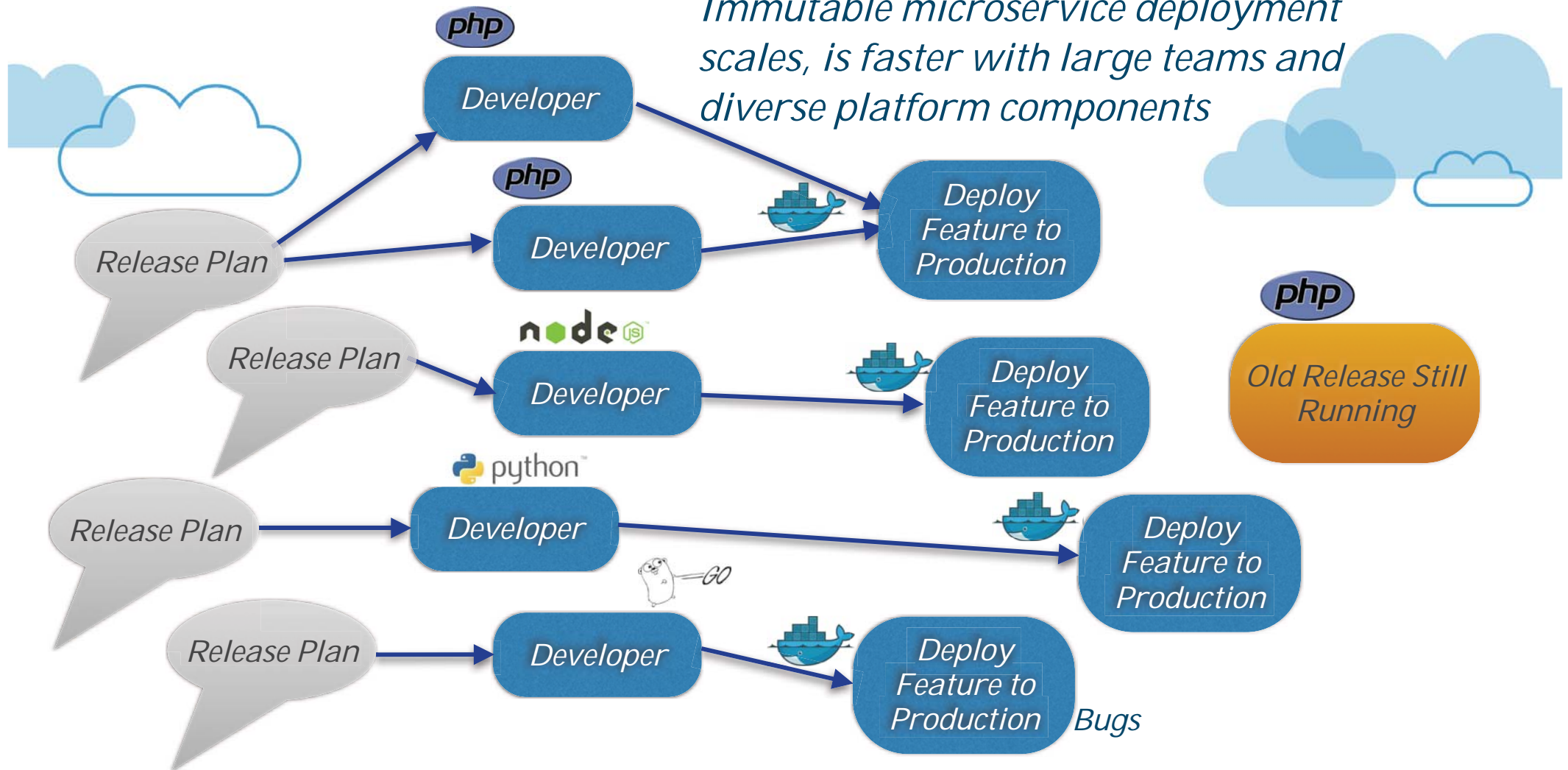
Immutable microservice deployment scales, is faster with large teams and diverse platform components



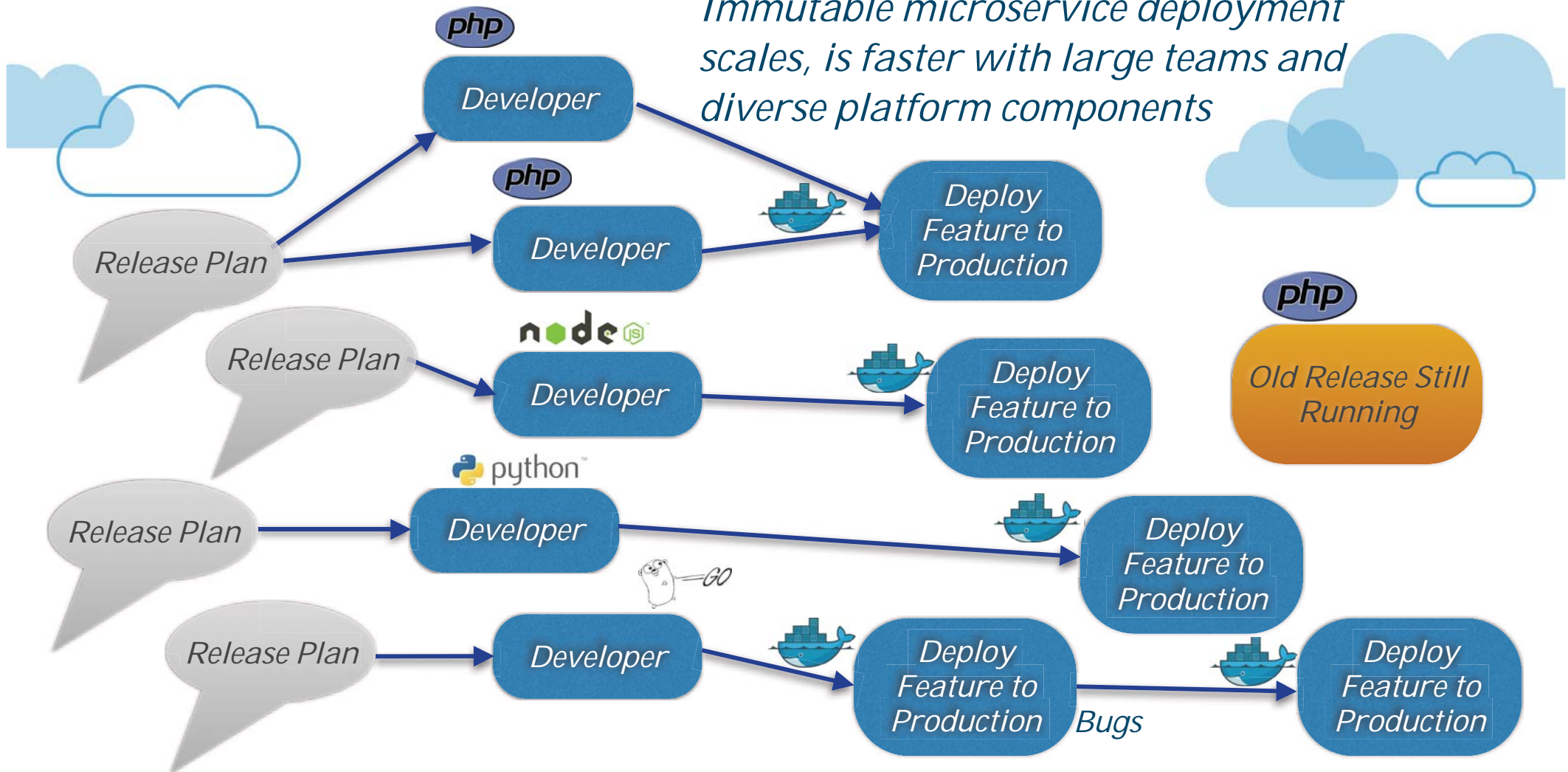
Immutable microservice deployment scales, is faster with large teams and diverse platform components

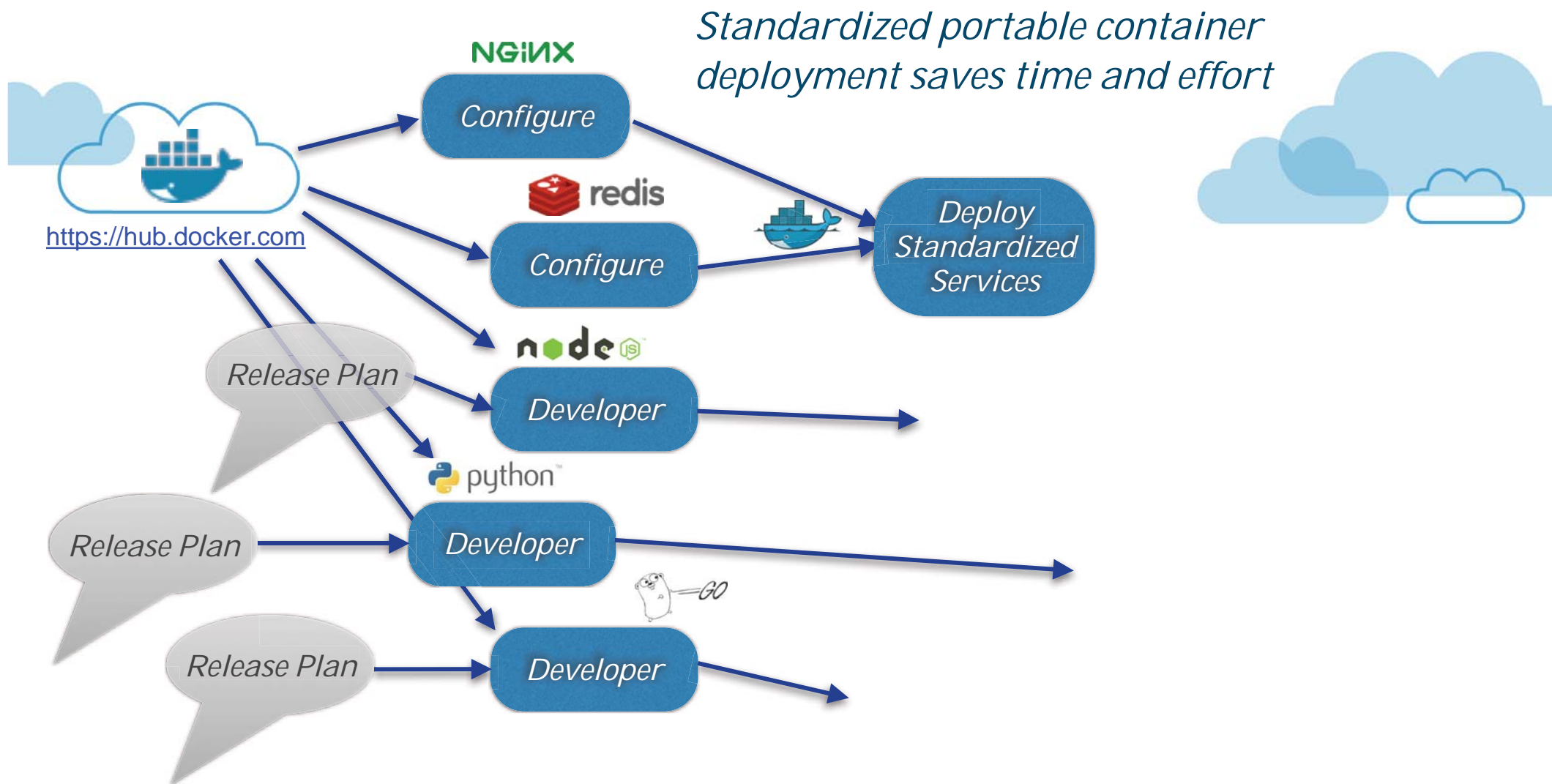


Immutable microservice deployment scales, is faster with large teams and diverse platform components

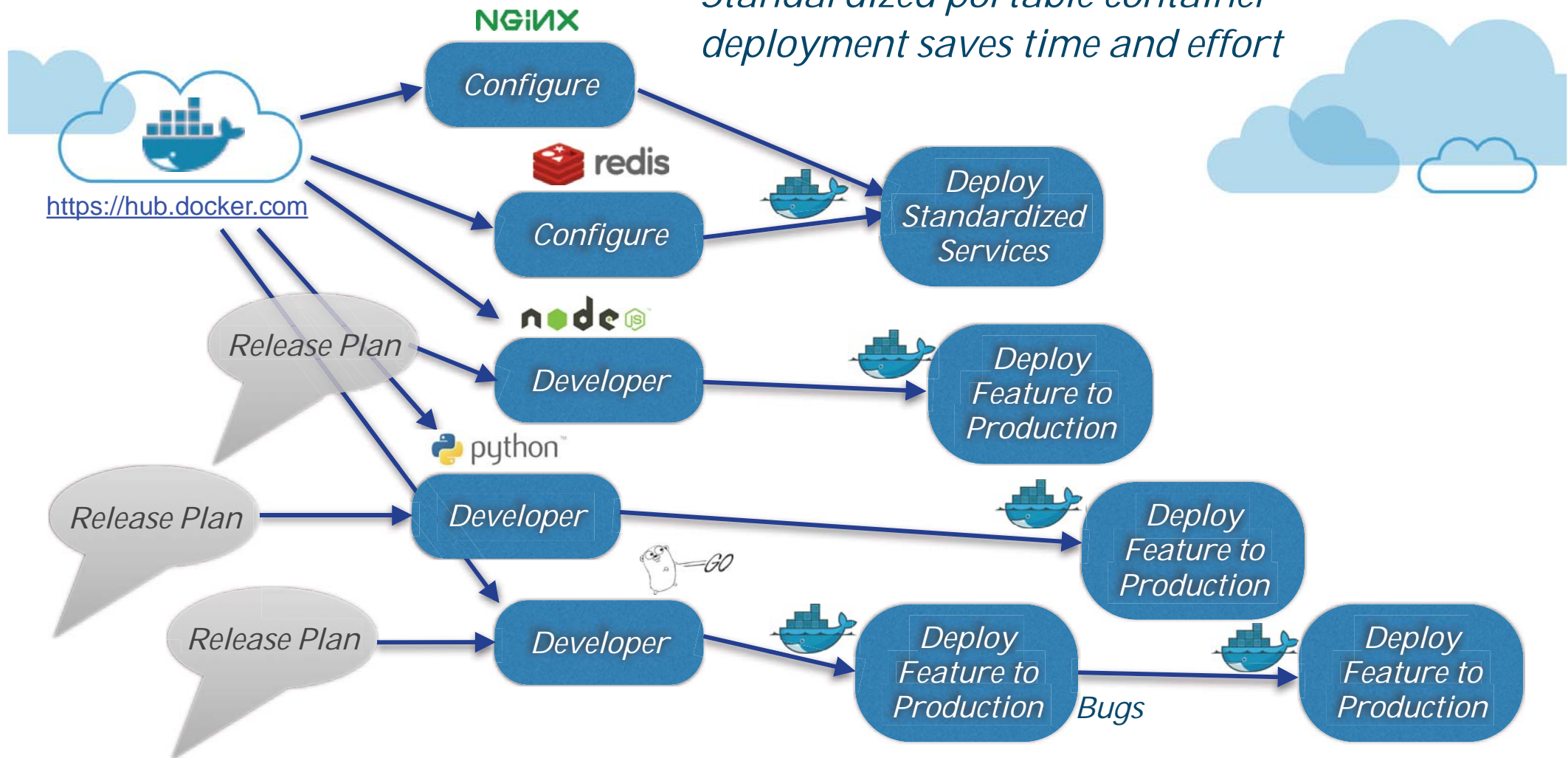


Immutable microservice deployment scales, is faster with large teams and diverse platform components





Standardized portable container deployment saves time and effort



Developing at the Speed of Docker

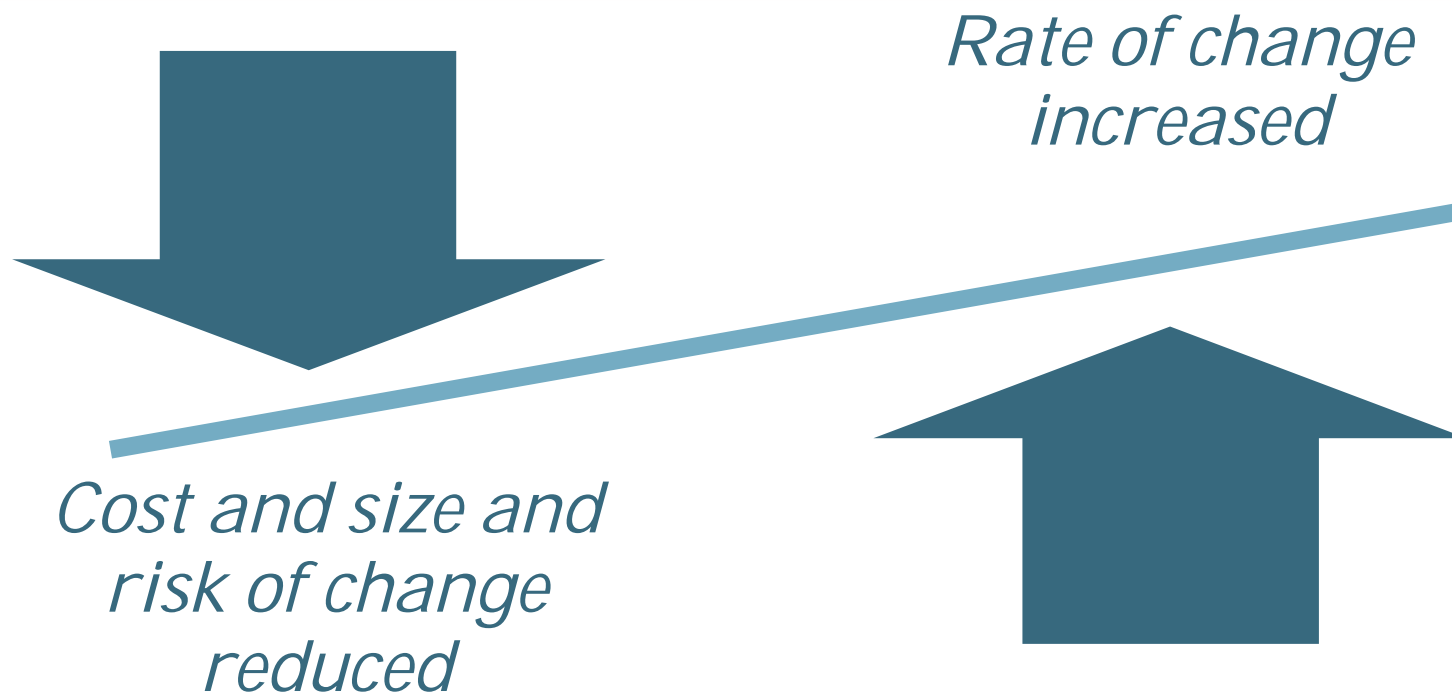


Developing at the Speed of Docker



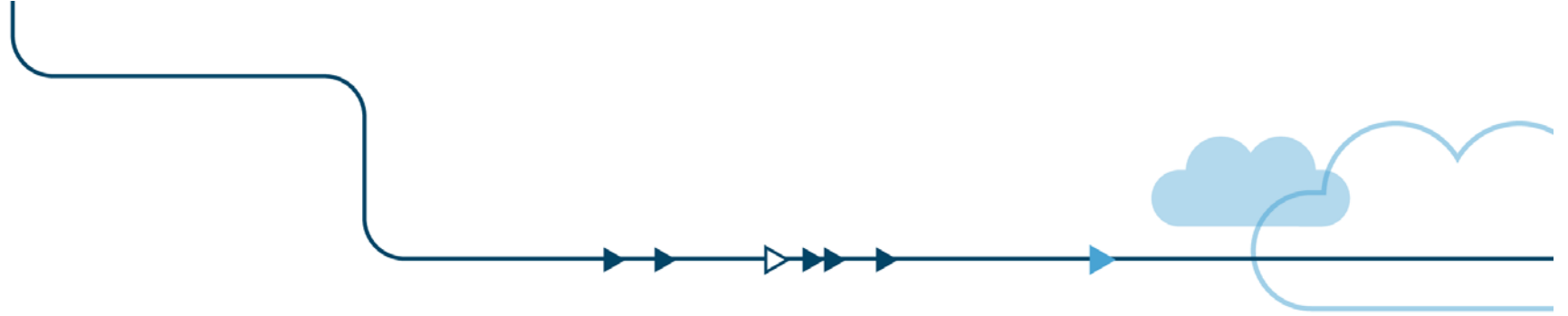
► *Speed is addictive, hard to go back to taking much longer to get things done*

What Happened?





Cloud Native Applications



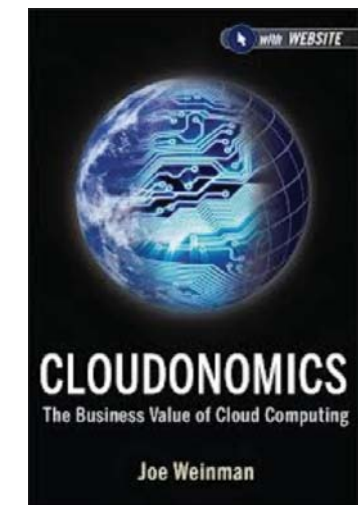
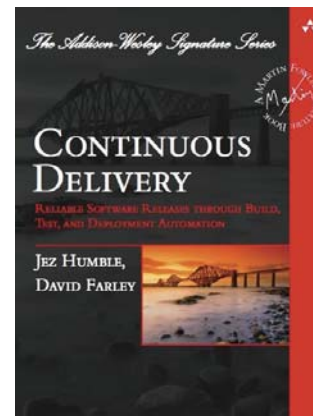
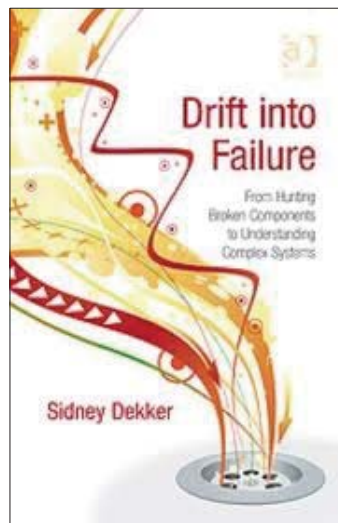
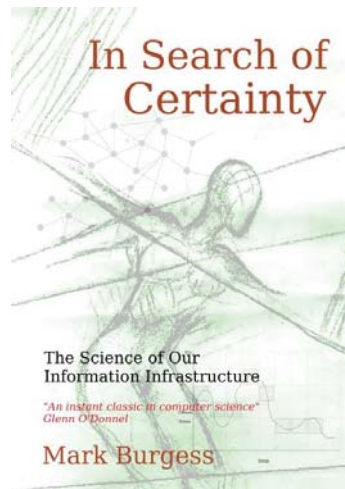
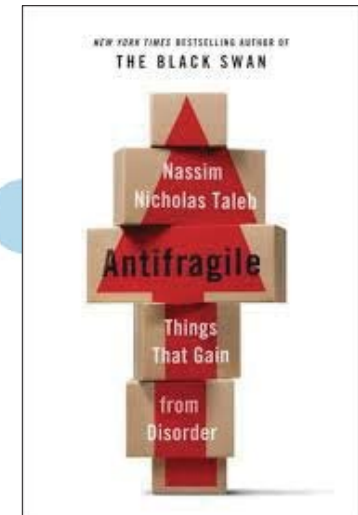
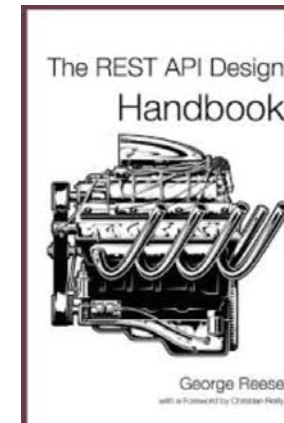
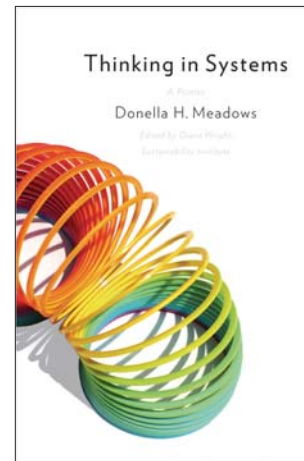
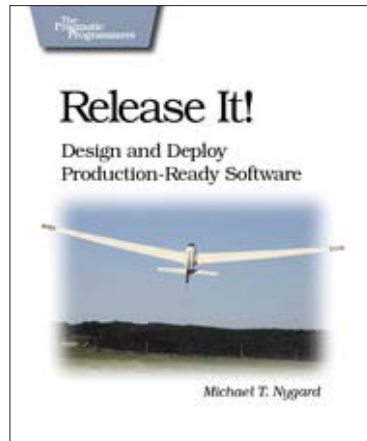
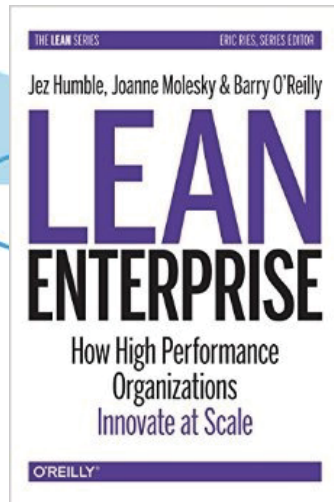
Cloud Native

A new engineering challenge

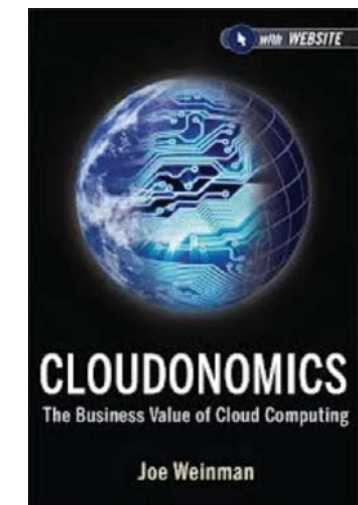
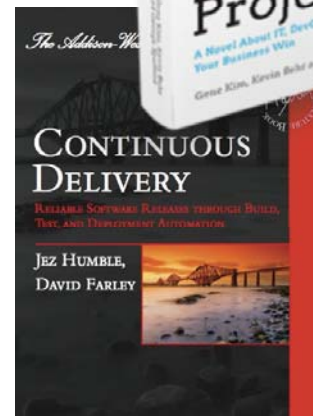
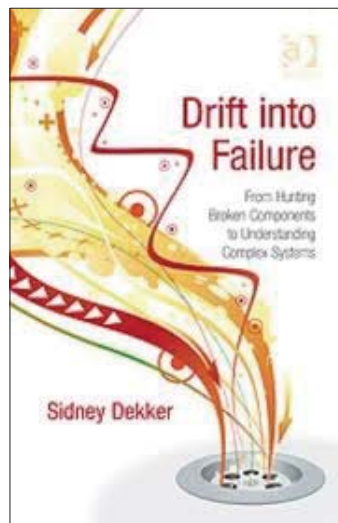
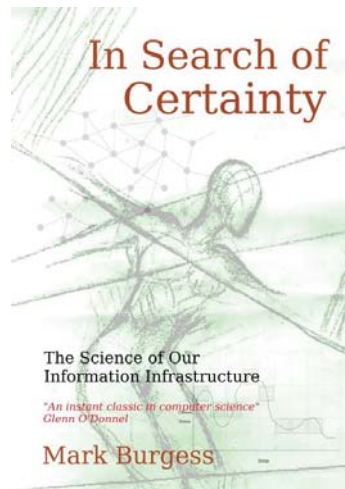
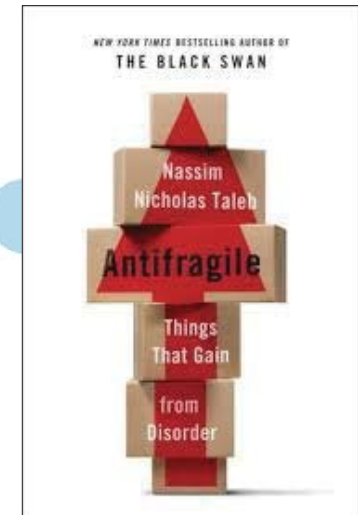
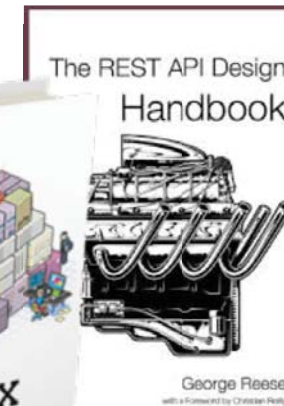
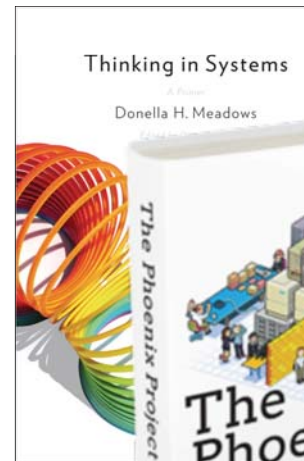
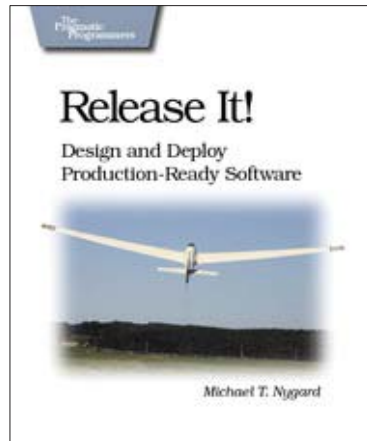
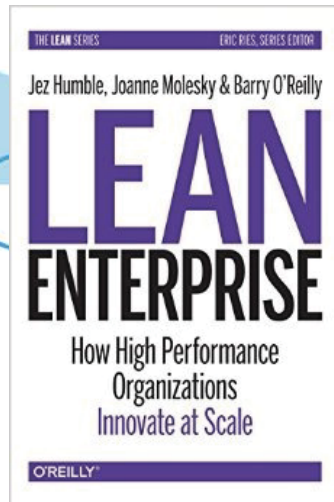
*Construct a highly agile and highly
available service from ephemeral and
assumed broken components*



Inspiration



Inspiration



State of the Art in Cloud Native Microservice Architectures



GILT

<http://www.infoq.com/presentations/scale-gilt>



NETFLIX | OSS

AWS Re:Invent : Asgard to Zuul <https://www.youtube.com/watch?v=p7ysHhs5hl0>
Resiliency at Massive Scale https://www.youtube.com/watch?v=ZfYJHtVL1_w
Microservice Architecture <https://www.youtube.com/watch?v=CriDUYtfrjs>



GROUPON

<http://www.slideshare.net/mcculloughsean/itier-breaking-up-the-monolith-philly-ete>



<http://www.infoq.com/presentations/Twitter-Timeline-Scalability>
<http://www.infoq.com/presentations/twitter-soa>
<http://www.infoq.com/presentations/Zipkin>



<https://speakerdeck.com/mattheath/scaling-micro-services-in-go-highload-plus-plus-2014>

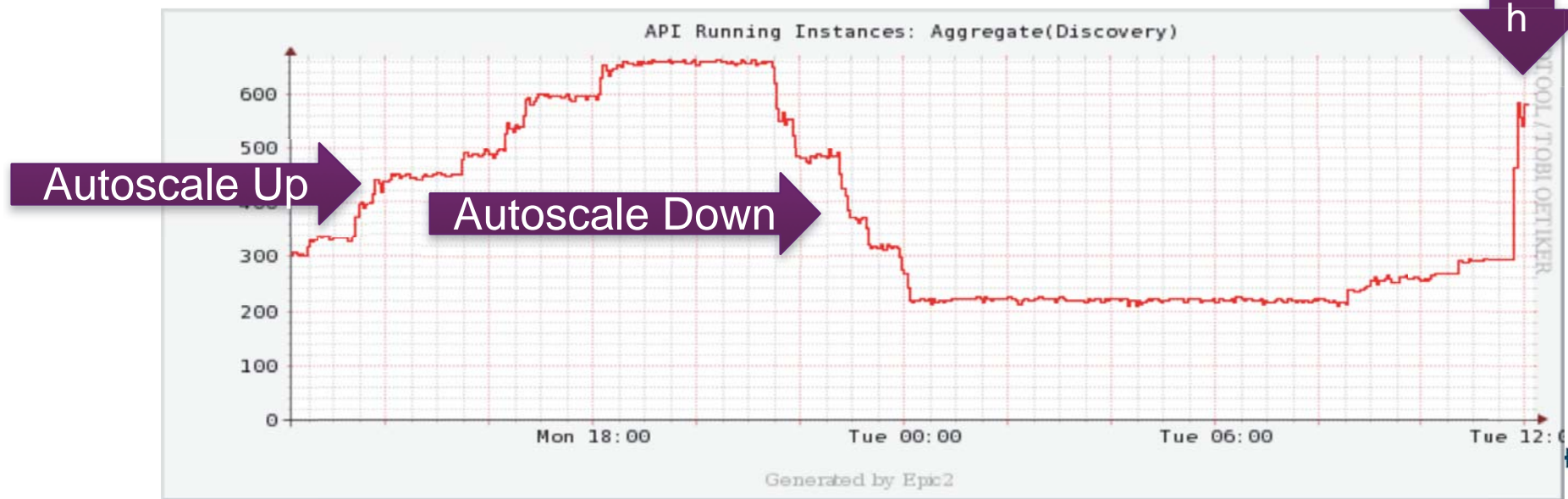
NETFLIX | OSS Trust with Verification

- *Edda - the “black box flight recorder” for configuration state*
- *Chaos Monkey - enforcing stateless business logic*
- *Chaos Gorilla - enforcing zone isolation/replication*
- *Chaos Kong - enforcing region isolation/replication*
- *Security Monkey - watching for insecure configuration settings*
- *See over 40 NetflixOSS projects at netflix.github.com*
- *Get “Technical Indigestion” trying to keep up with techblog.netflix.com*



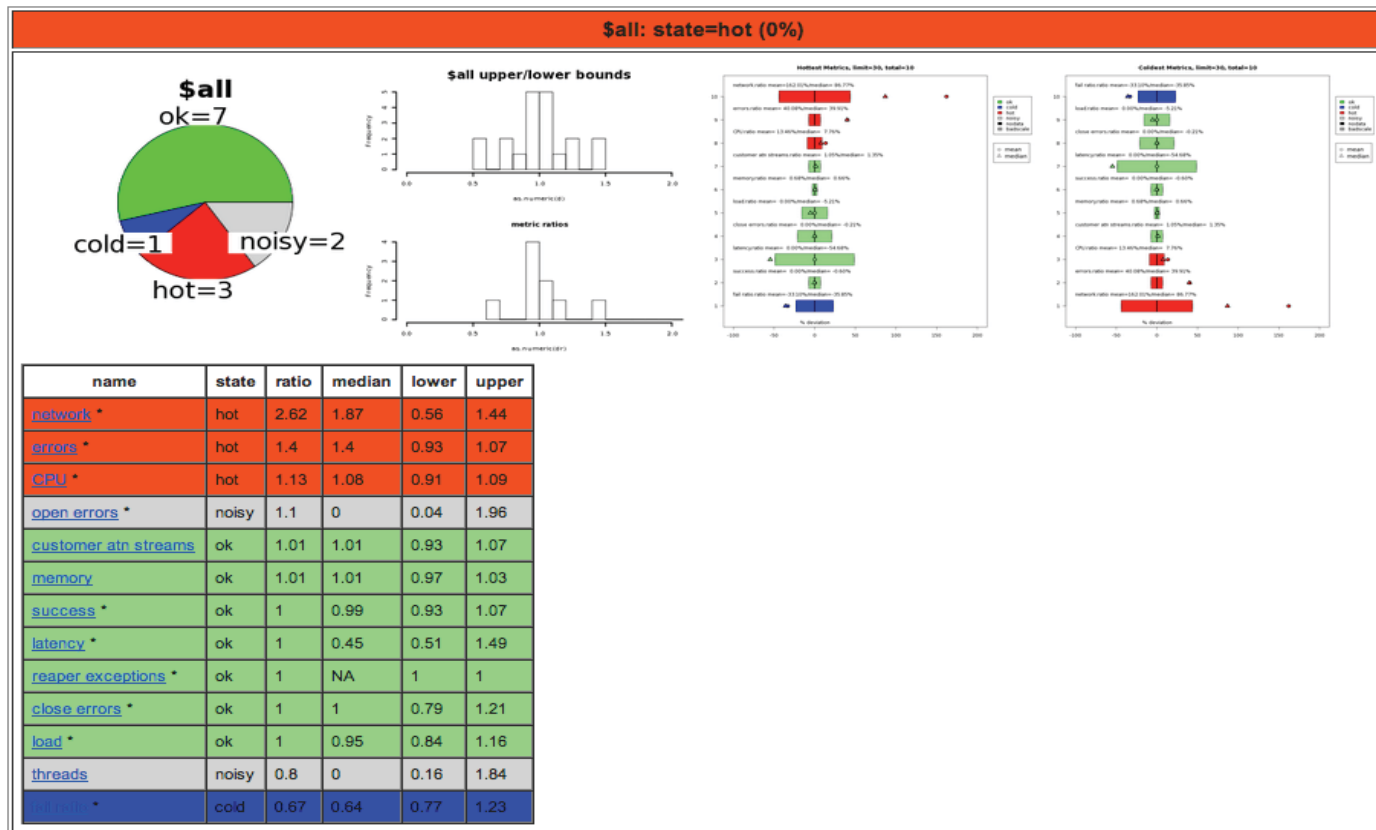
Autoscaled Ephemeral Instances at Netflix

Largest services use autoscaled red/black code pushes
Average lifetime of an instance is 36 hours



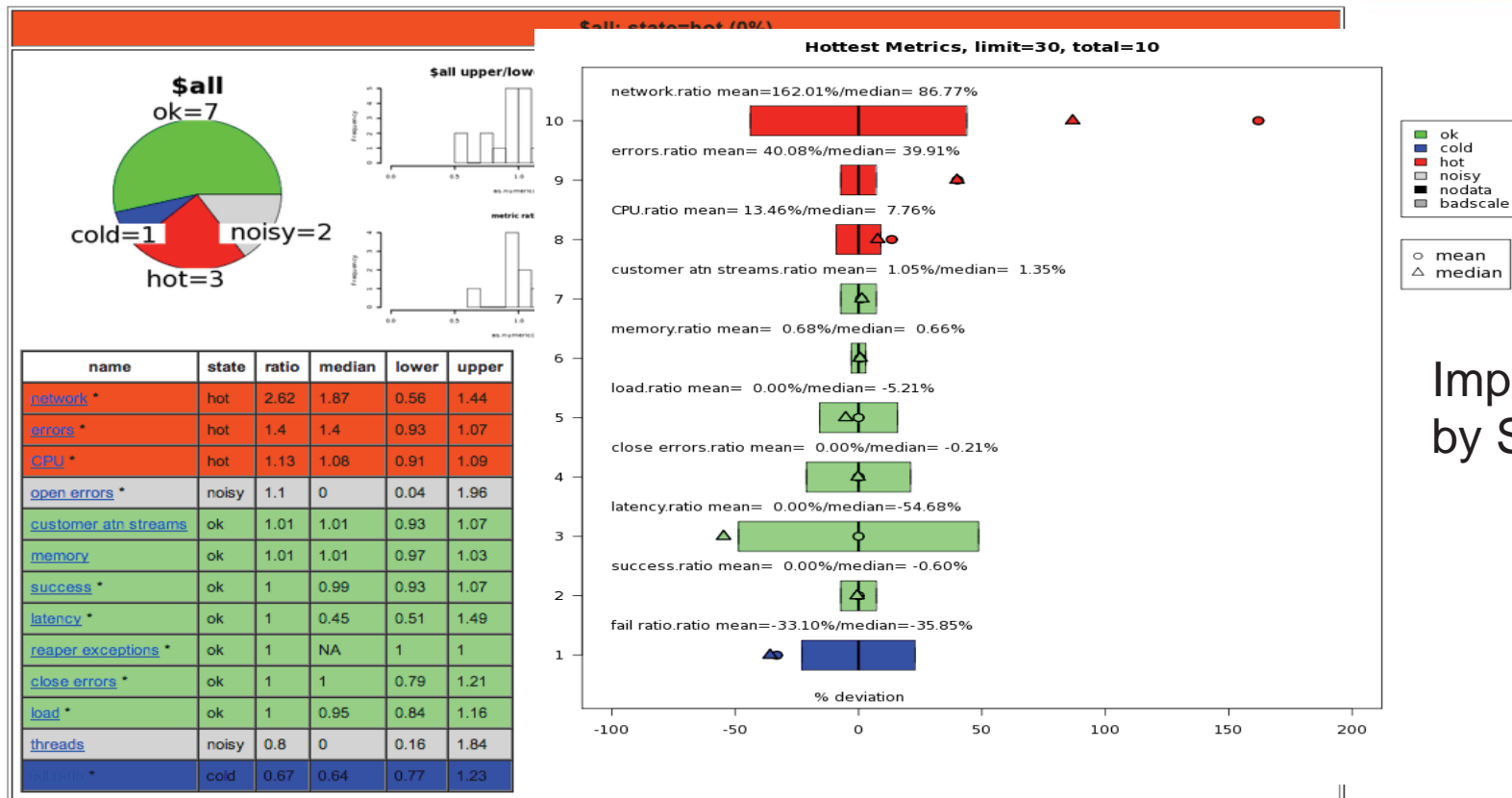
Netflix Automatic Code Deployment Canary

Bad Signature



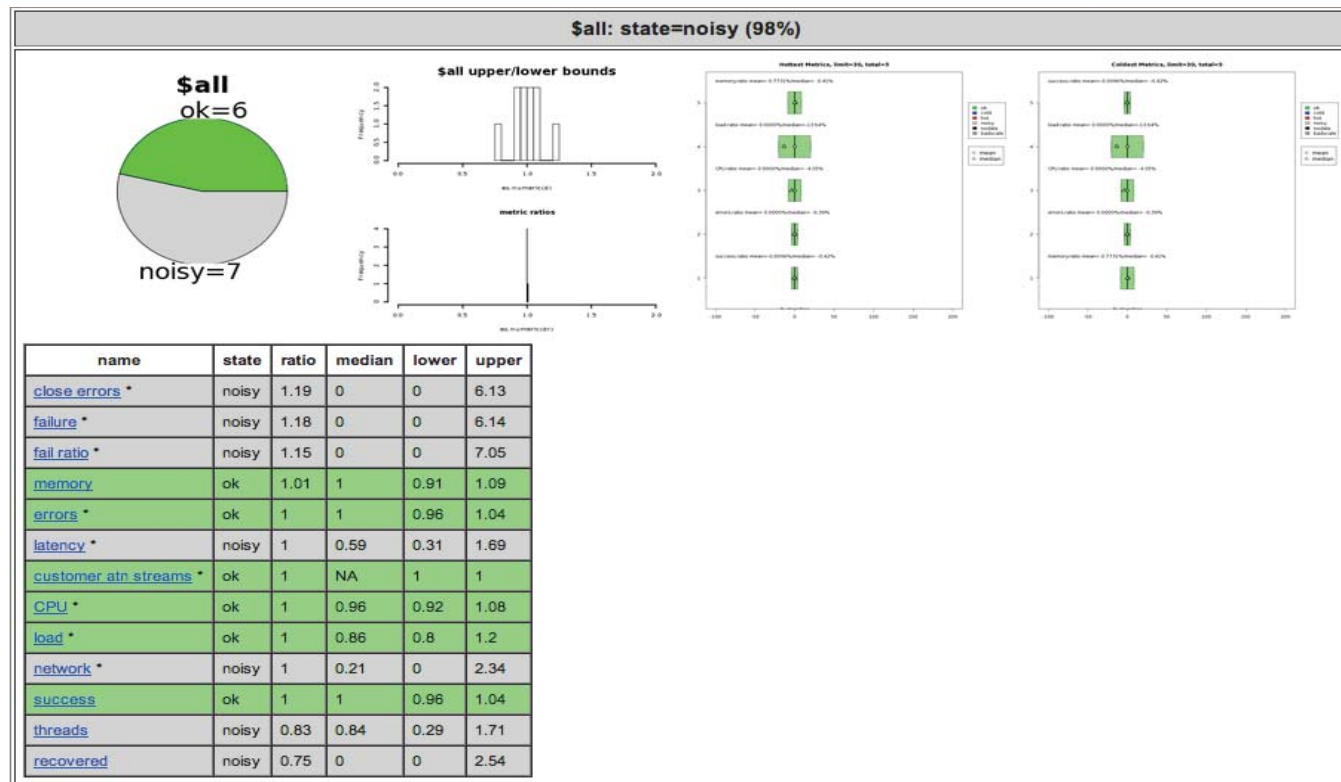
Implemented
by Simon Tuffs

Netflix Automatic Code Deployment Canary Bad Signature



Implemented
by Simon Tuffs

Happy Canary Signature



Speeding Up The Platform



Datacenter Snowflakes

- Deploy in months
- Live for years

Speeding Up The Platform



Datacenter Snowflakes

- Deploy in months
- Live for years



Virtualized and Cloud

- Deploy in minutes
- Live for weeks

Speeding Up The Platform



Datacenter Snowflakes

- Deploy in months
- Live for years



Virtualized and Cloud

- Deploy in minutes
- Live for weeks



Docker Containers

- Deploy in seconds
- Live for minutes/hours

Speeding Up The Platform



Datacenter Snowflakes

- Deploy in months
- Live for years



Virtualized and Cloud

- Deploy in minutes
- Live for weeks



Docker Containers

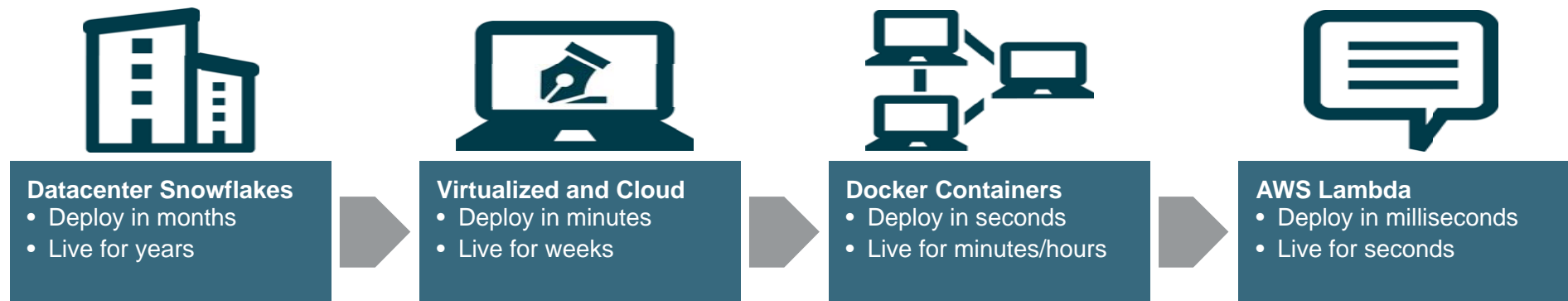
- Deploy in seconds
- Live for minutes/hours



AWS Lambda

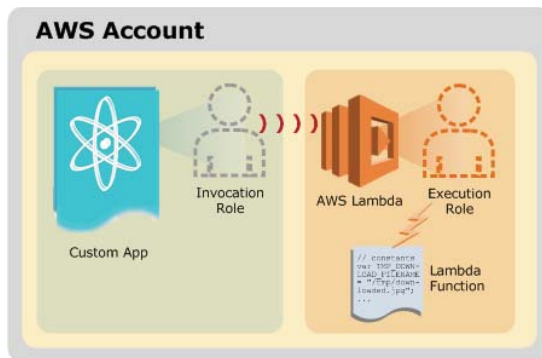
- Deploy in milliseconds
- Live for seconds

Speeding Up The Platform

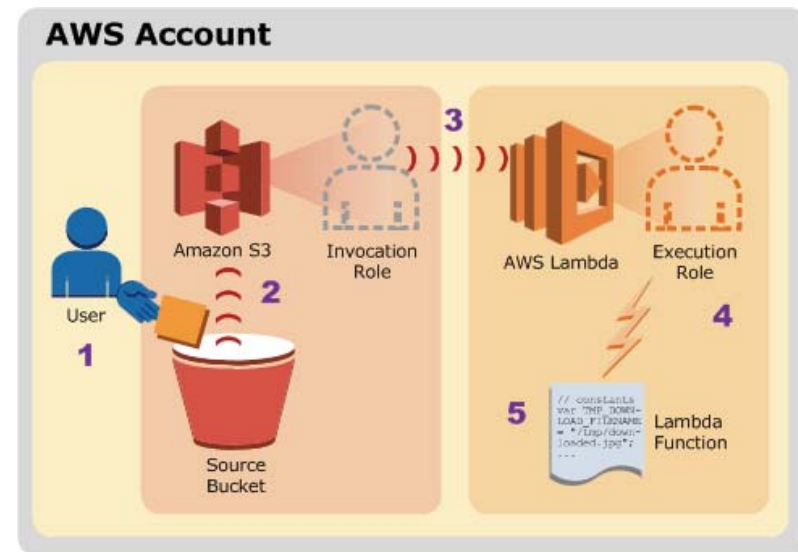
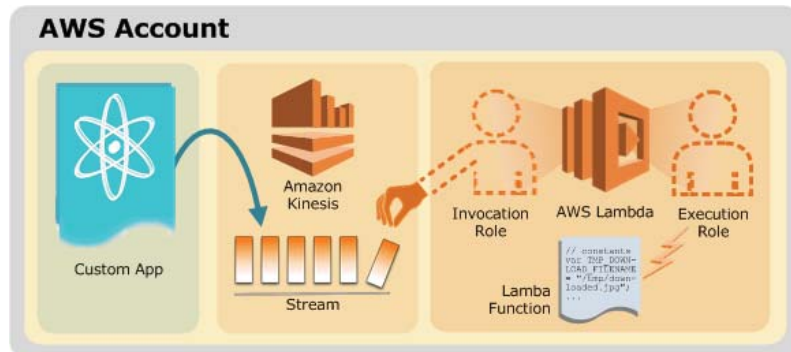


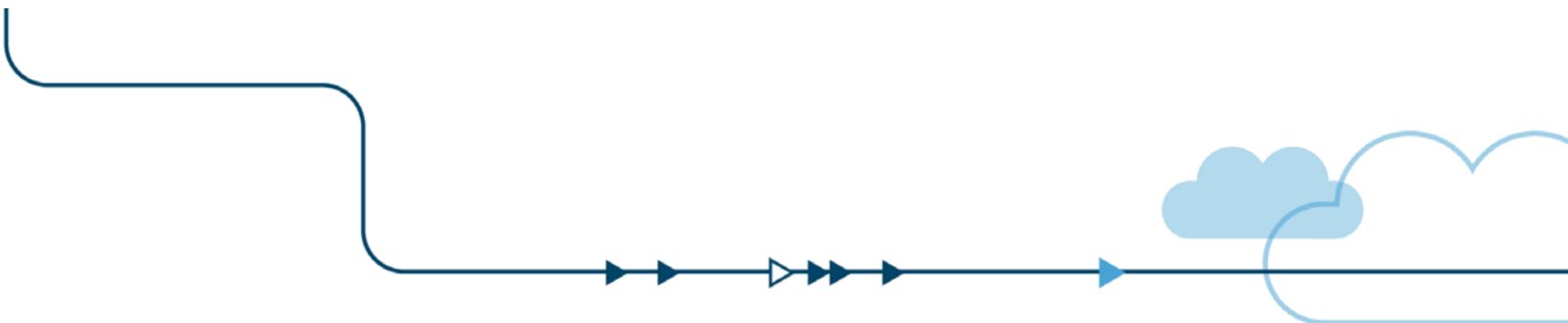
➤ *Speed enables and encourages new microservice architectures*

*With AWS Lambda
compute resources are charged
by the 100ms, not the hour*



*First 1,000,000 node.js executions/month are free
First 400,000 GB-seconds of RAM-CPU are free*





Monitoring Requirements

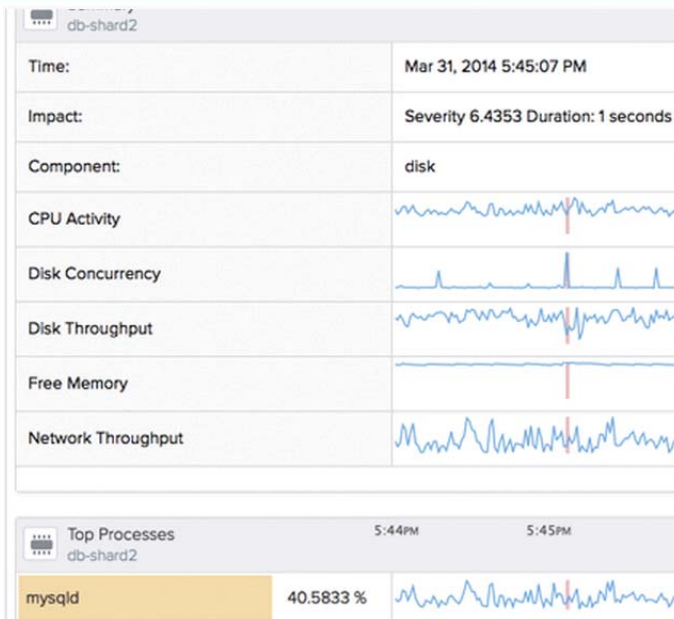
Metric resolution microseconds

Metric update rate 1 second

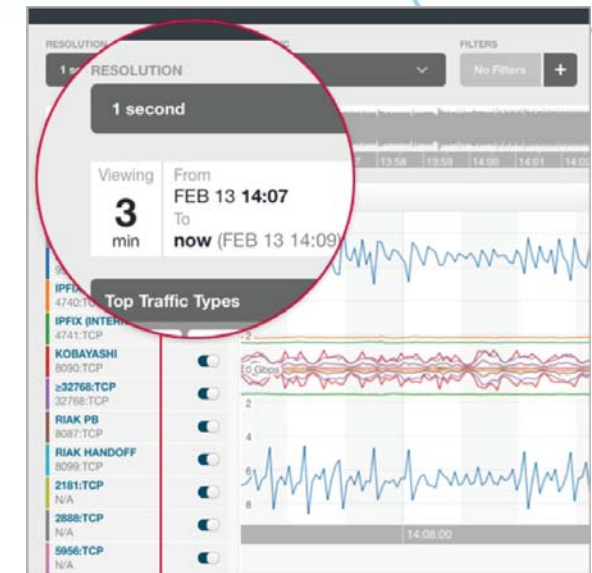
*Metric to display latency less than human
attention span (<10s)*



Low Latency SaaS Based Monitors



boundary

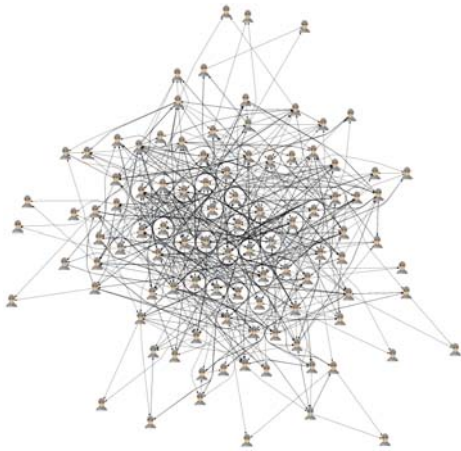


1-second data collection and real-time streaming processing on all components of the application stack

www.vividcortex.com and www.boundary.com

@adriano
BV
Battery Ventures

Adrian's Tinkering Projects

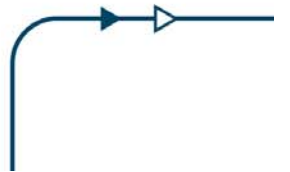


*Model and visualize microservices
Simulate interesting architectures*

See github.com/adrianco/spigo
Simulate Protocol Interactions in Go



See github.com/adrianco/d3grow
Dynamic visualization





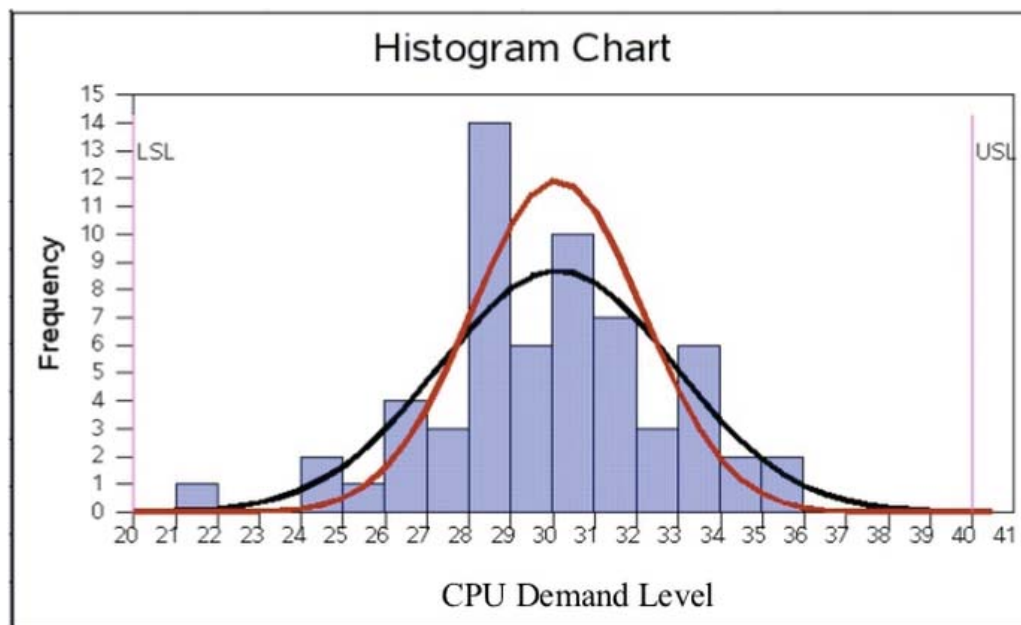
Cost Optimization



Lower Spec Limit

When demand probability is below USL by 3.0 sigma scale down resource to save money

Capacity Optimization for a Single System Bottleneck



Upper Spec Limit

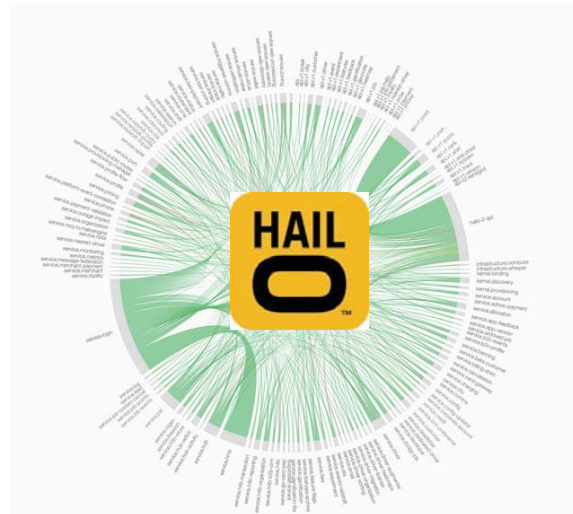
When demand probability exceeds USL by 4.0 sigma scale up resource to maintain low latency

[Documentation on Capability Plots](#)

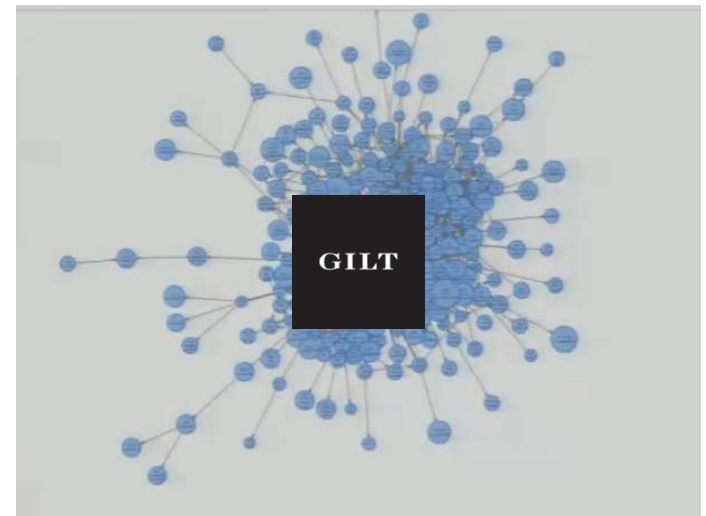
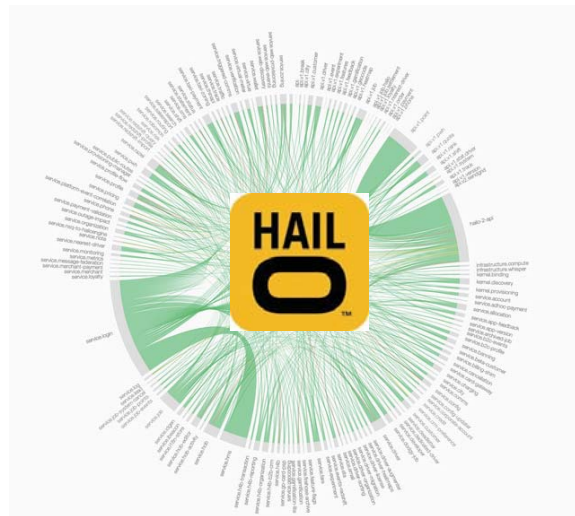
To get accurate high dynamic range histograms see <http://hdrhistogram.org/>

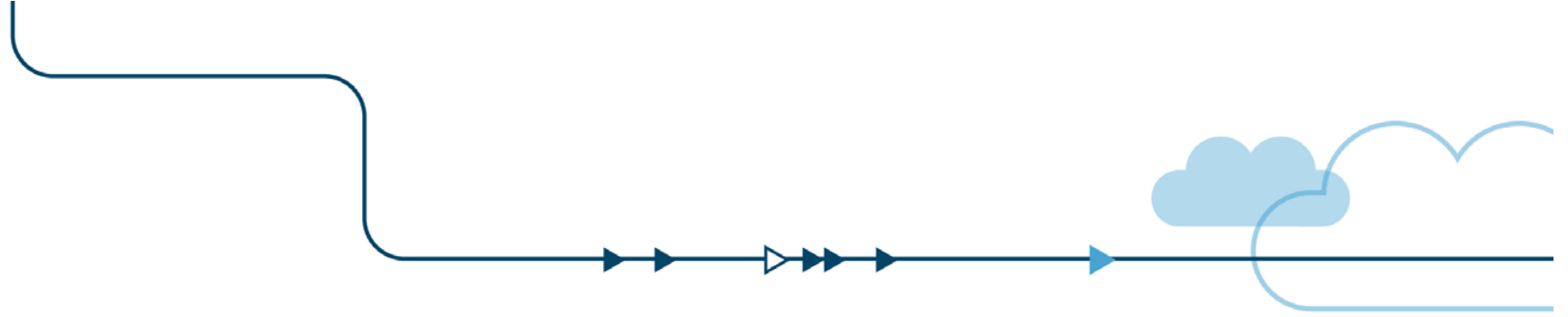
[Slideshare: 2003 Presentation on Capacity Planning Methods](#) [See US Patent: 7467291](#)

*But interesting systems
don't have a single
bottleneck nowadays...*



*But interesting systems
don't have a single
bottleneck nowadays...*





What about cloud costs?

Cloud Native Cost Optimization

\$ \$ \$

Optimize for speed first

Turn it off!

Capacity on demand

Consolidate and Reserve

Plan for price cuts

FOSS tooling



The Capacity Planning Problem



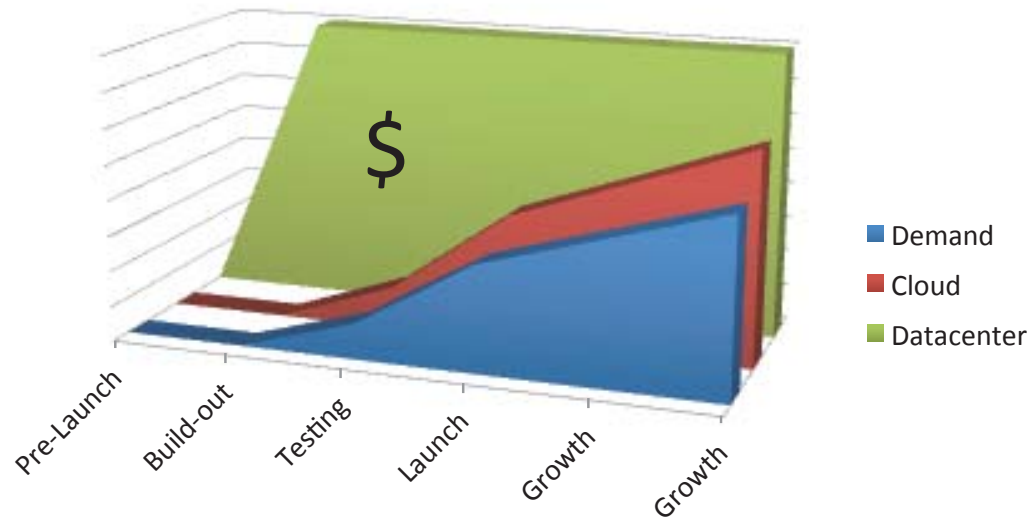
@adrianco

BV
Battery Ventures

Best Case Waste



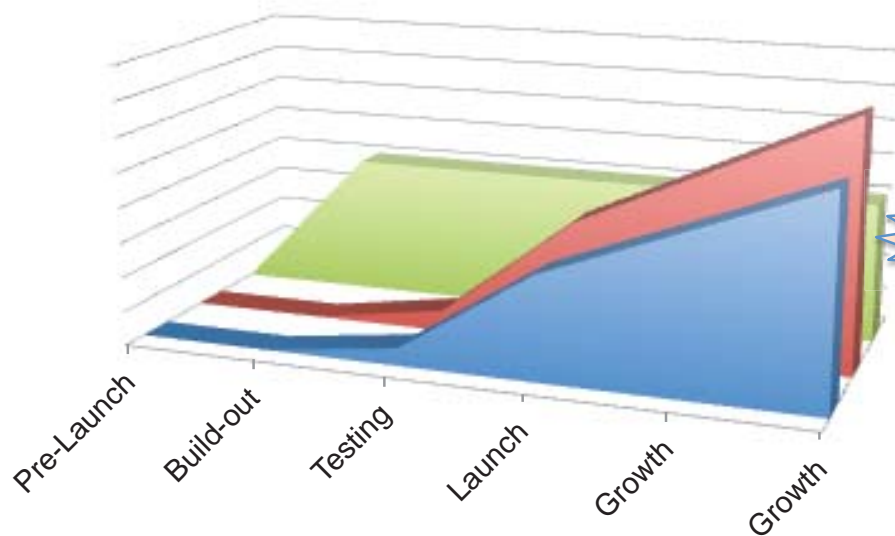
Product Launch Agility - Rightsized



Cloud capacity used is maybe half average DC capacity

Failure to Launch

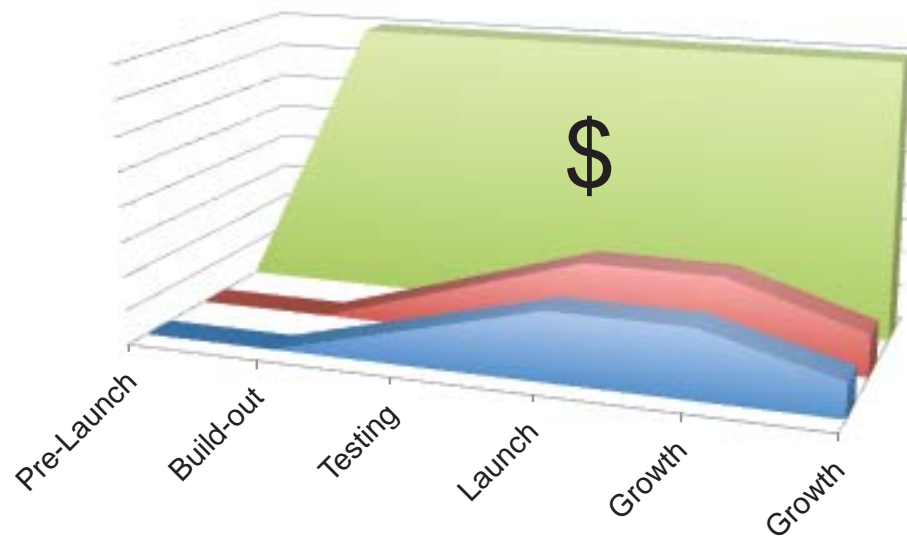
Product Launch - Under-estimated



*Mad scramble
to add more DC
capacity during
launch phase
outages*

Over the Top Losses

Product Launch Agility – Over-estimated



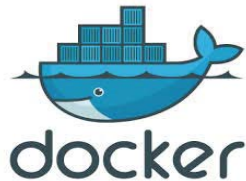
*Capacity wasted
on failed launch
magnifies the
losses*

Turning off Capacity



Off-peak production
Test environments
Dev out of hours
Dormant Data Science

Containerize Test Environments



Snapshot or freeze

Fast restart needed

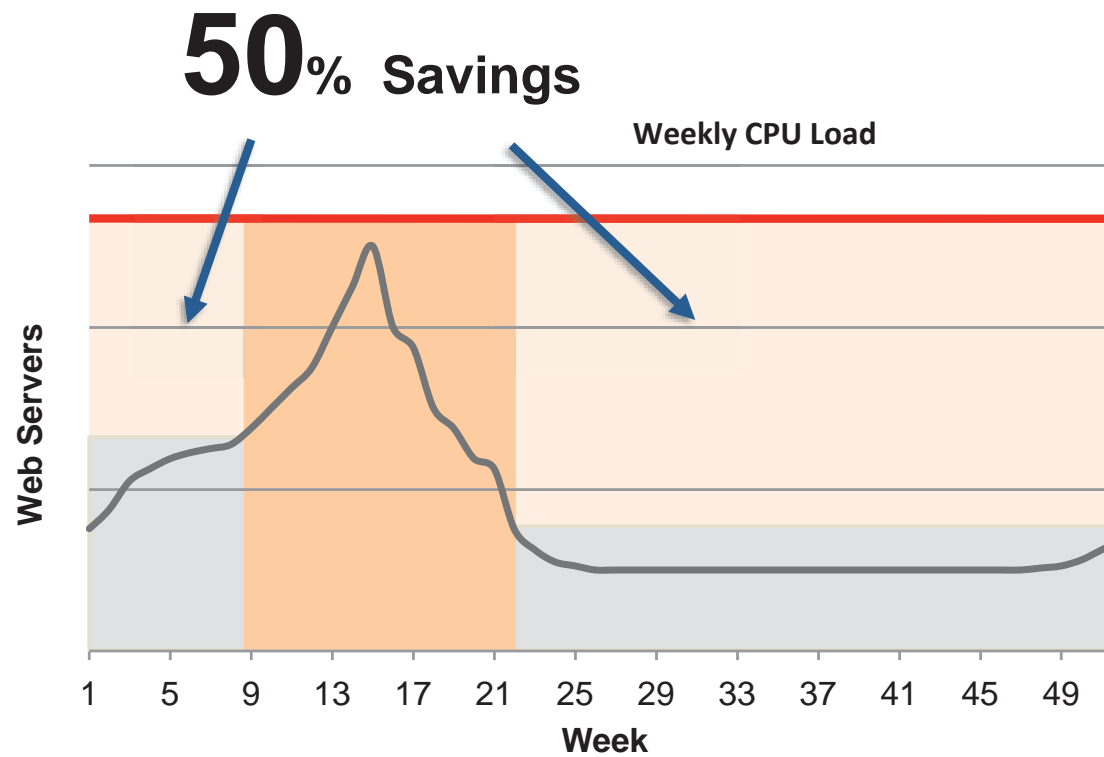
Persistent storage

40 of 168 hrs/wk

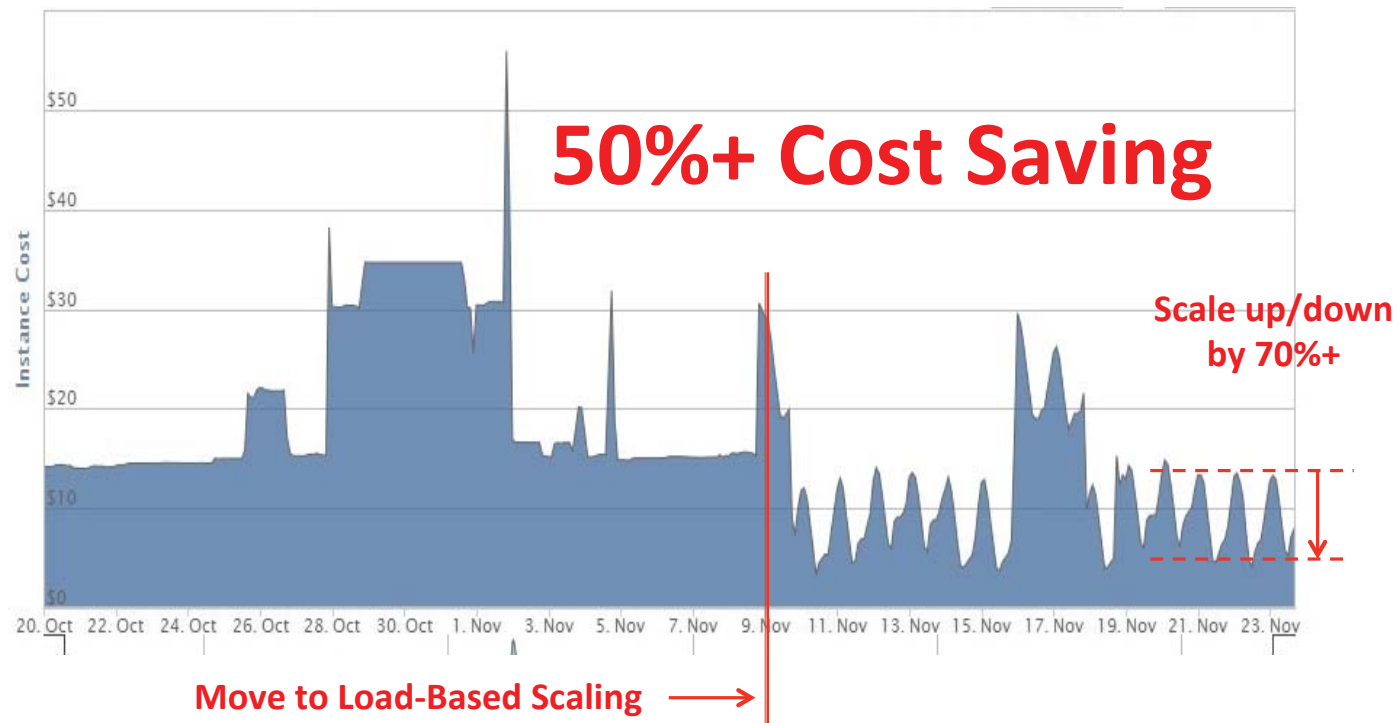
Bin-packed containers

shippable.com saved 70%

Seasonal Savings



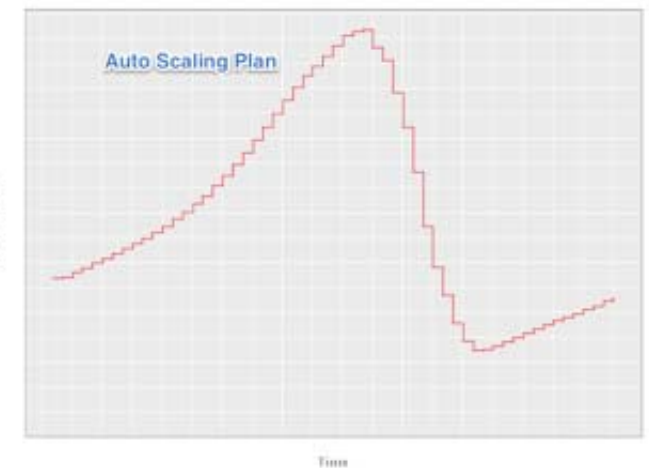
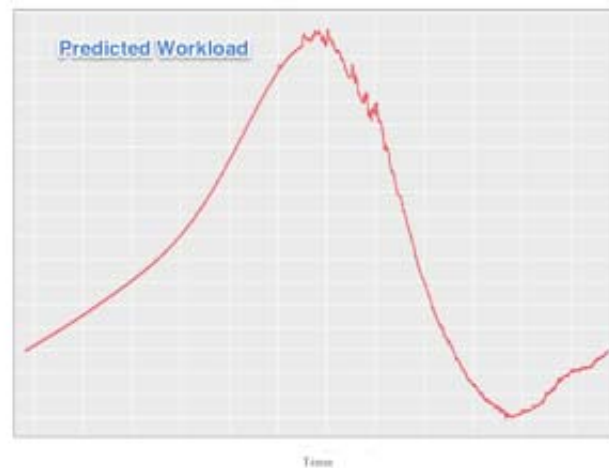
Autoscale the Costs Away



Daily Duty Cycle



*Reactive Autoscaling
saves around 50%*



*Predictive Autoscaling saves around 70%
See Stryer on Netflix Tech Blog*

Underutilized and Unused

AWS Support – Trusted Advisor – Your personal cloud assistant

Trusted Advisor Beta

[Expand All](#)[Download Excel](#)[Refresh All](#)[Contact Support](#)

The AWS Trusted Advisor program monitors AWS infrastructure services, identifies customer configurations, compares them to known best practices, and then notifies customers when opportunities may exist to save money, improve system performance, or close security gaps.



No issue detected



Investigation Recommended



Action Recommended

▼ Cost Optimizing Checks

✓ Unused Elastic IPs ?

Updated: 2012-06-14 00:00 PDT ↻

➤ Summary: 0 of 6 Elastic IPs are not in use

⚠ Underutilized EC2 Instances ?

Updated: 2012-06-13 22:27 PDT ↻

➤ Summary: 27 EC2 instances are potentially underutilized

Clean Up the Crud

Other simple optimization tips

- **Don't forget to...**
 - Disassociate unused EIPs
 - Delete unassociated Amazon EBS volumes
 - Delete older Amazon EBS snapshots
 - Leverage Amazon S3 Object Expiration



Janitor Monkey cleans up unused resources



Total Cost of Oranges

When Comparing TCO...

Make sure that
you are including
all the cost factors
into consideration

Place
Power
Pipes
People
Patterns



Total Cost of Oranges

When Comparing TCO...

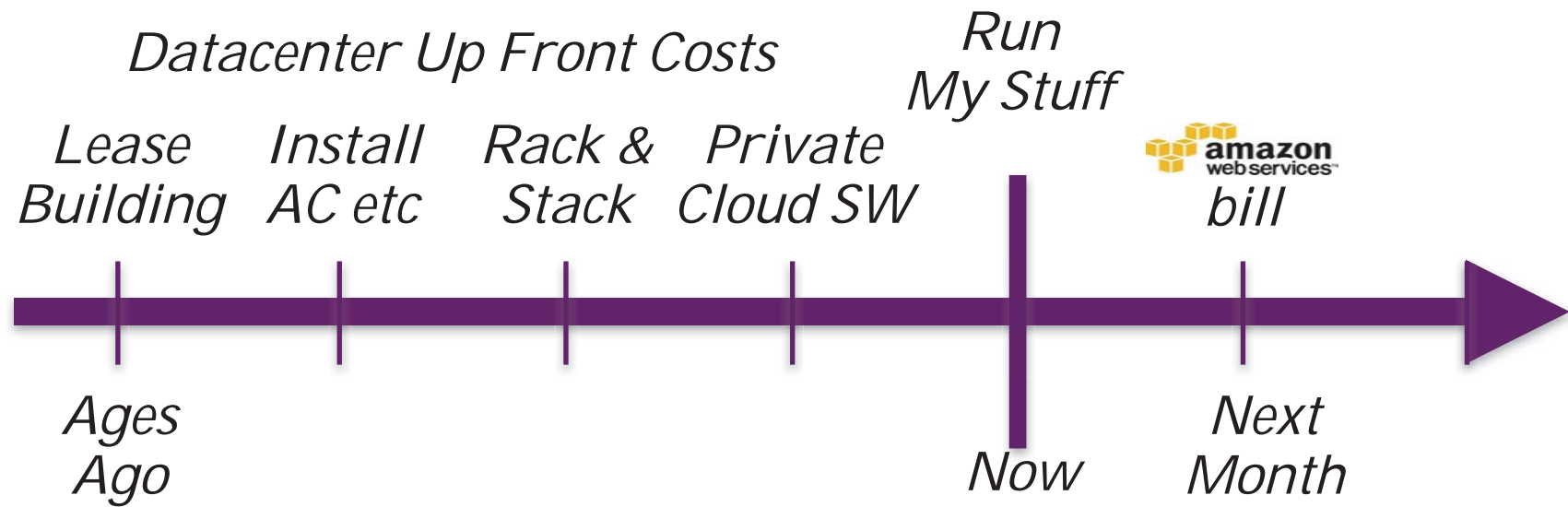
Make sure that
you are including
all the cost factors
into consideration

Place
Power
Pipes
People
Patterns

*How much does
datacenter automation
software and support
cost per instance?*



When Do You Pay?





Cost Model Comparisons




AWS has most complex model

- *Both highest and lowest cost options!*

CPU/Memory Ratios Vary

- *Can't get same config everywhere*

Features Vary

- *Local SSD included on some vendors, not others*
 - *Network and storage charges also vary*
- 

Digital Ocean Flat Pricing



<i>Hourly Price (\$0.06/hr)</i>	<i>Monthly Price (\$40/mo)</i>
<i>\$ No Upfront</i>	<i>\$ No Upfront</i>
<i>\$0.060/hr</i>	<i>\$0.056/hr</i>
<i>\$1555/36mo</i>	<i>\$1440/36mo</i>
<i>Savings</i>	<i>7%</i>

Prices on Dec 7th, for 2 Core, 4G RAM, SSD, purely to show typical savings

Google Sustained Usage



<i>Full Price Without Sustained Usage</i>	<i>Typical Sustained Usage Each Month</i>	<i>Full Sustained Usage Each Month</i>
<i>\$ No Upfront</i>	<i>\$ No Upfront</i>	<i>\$ No Upfront</i>
<i>\$0.063/hr</i>	<i>\$0.049/hr</i>	<i>\$0.045/hr</i>
<i>\$1633/36mo</i>	<i>\$1270/36mo</i>	<i>\$1166/36mo</i>
<i>Savings</i>	<i>22%</i>	<i>29%</i>

Prices on Dec 7th, for n1.standard-1 (1 vCPU, 3.75G RAM, no disk) purely to show typical savings

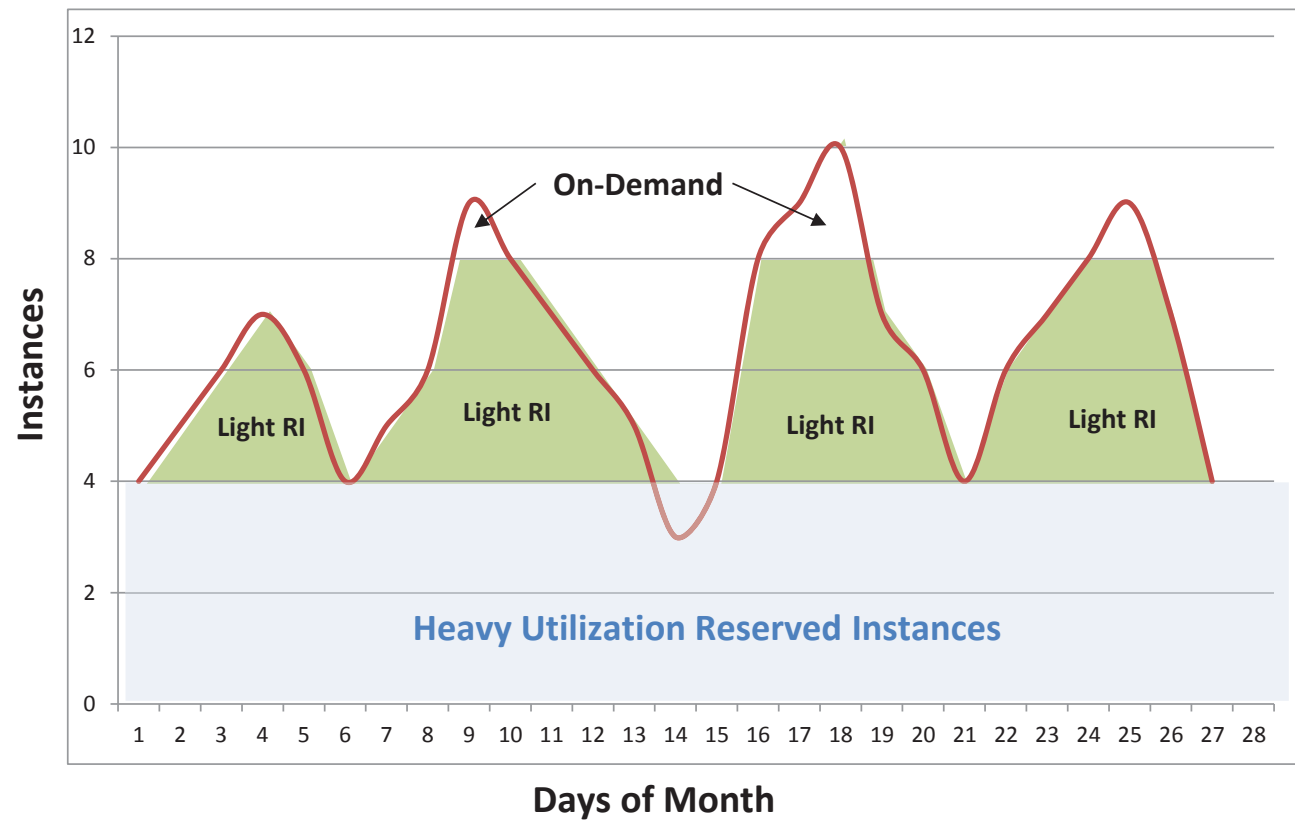
AWS Reservations



<i>On Demand</i>	<i>No Upfront 1 year</i>	<i>Partial Upfront 3 year</i>	<i>All Upfront 3 year</i>
<i>\$ No Upfront</i>	<i>\$No Upfront</i>	<i>\$337 Upfront</i>	<i>\$687 Upfront</i>
<i>\$0.070/hr</i>	<i>\$0.050/hr</i>	<i>\$0.0278/hr</i>	<i>\$0.00/hr</i>
<i>\$1840/36mo</i>	<i>\$1314/36mo</i>	<i>\$731/36mo</i>	<i>\$687/36mo</i>
<i>Savings</i>	<i>29%</i>	<i>60%</i>	<i>63%</i>

Prices on Dec 7th, for m3.medium (1 vCPU, 3.75G RAM, SSD) purely to show typical savings

Blended Benefits



On Demand

Partial Upfront

All Upfront

Consolidated Reservations

Burst capacity guarantee

Higher availability with lower cost

Other accounts soak up any extra

Monthly billing roll-up

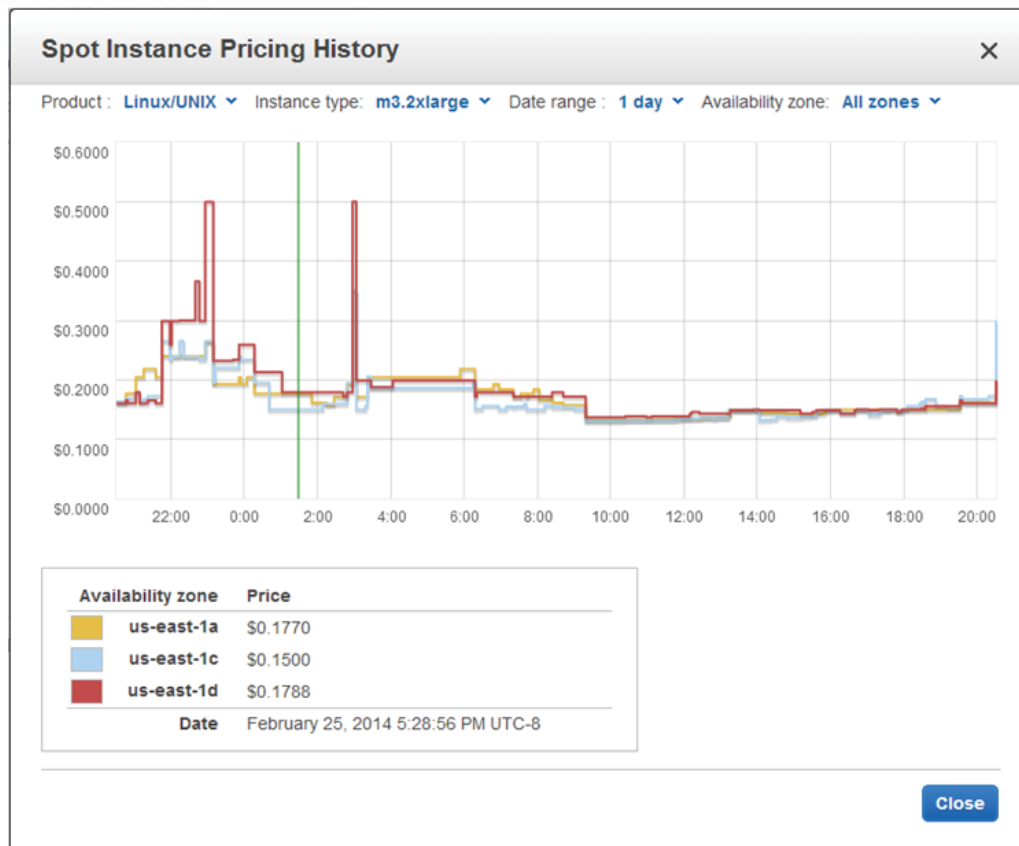
Capitalize upfront charges!

But: Fixed location and instance type

@adrianco

BV
Battery Ventures

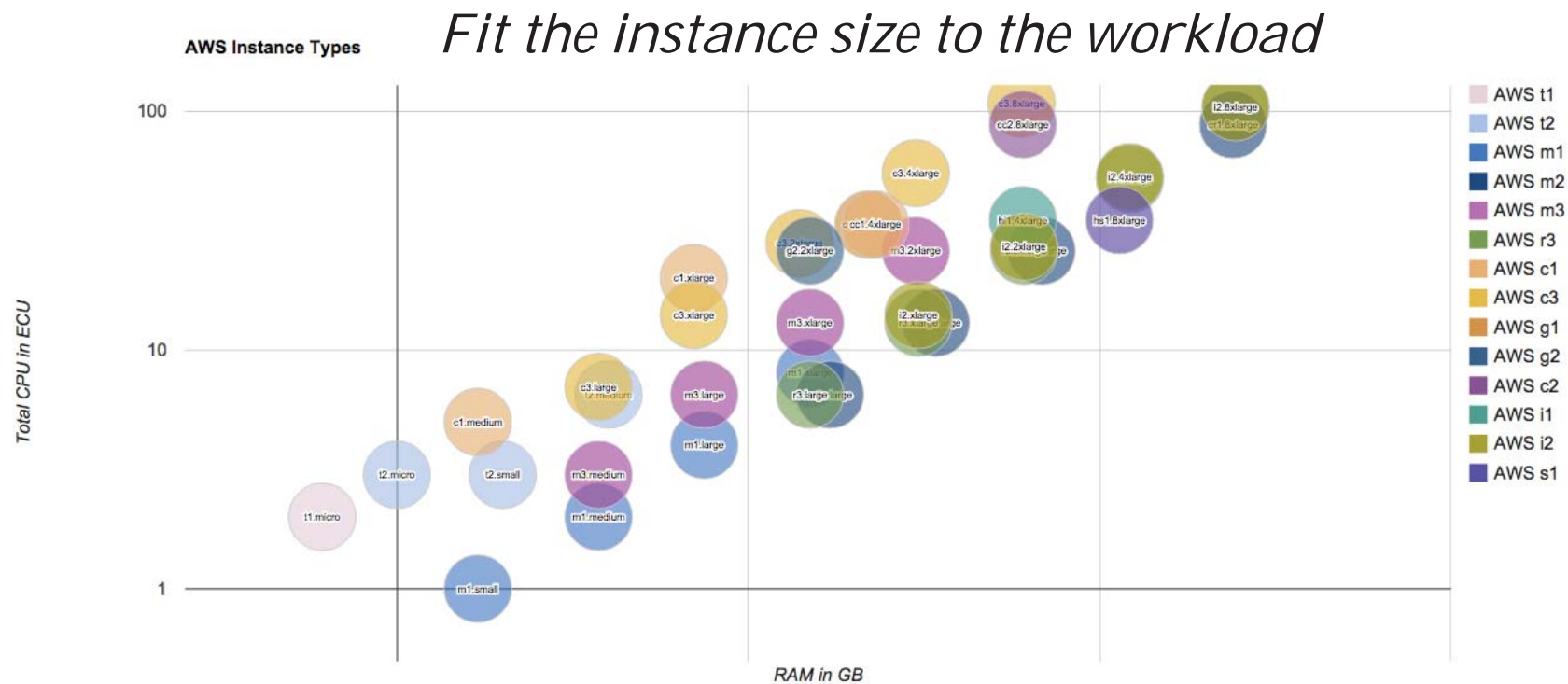
Use EC2 Spot Instances



*Cloud native
dynamic autoscaled
spot instances*

*Real world total
savings up to 50%*

Right Sizing Instances



Six Ways to Cut Costs



#1 Business Agility by Rapid Experimentation = Profit

#2 Business-driven Auto Scaling Architectures = Savings

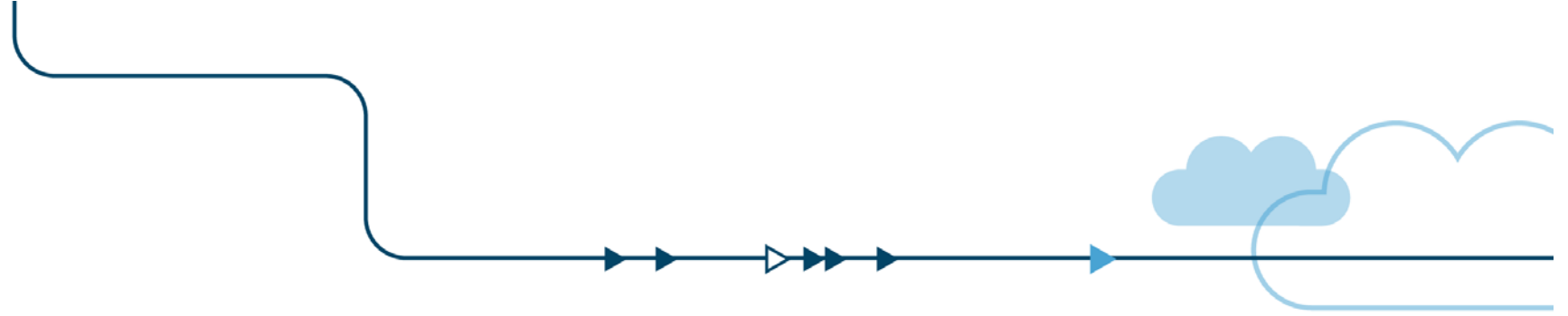
#3 Mix and Match Reserved Instances with On-Demand = Savings

#4 Consolidated Billing and Shared Reservations = Savings

#5 Always-on Instance Type Optimization = Recurring Savings

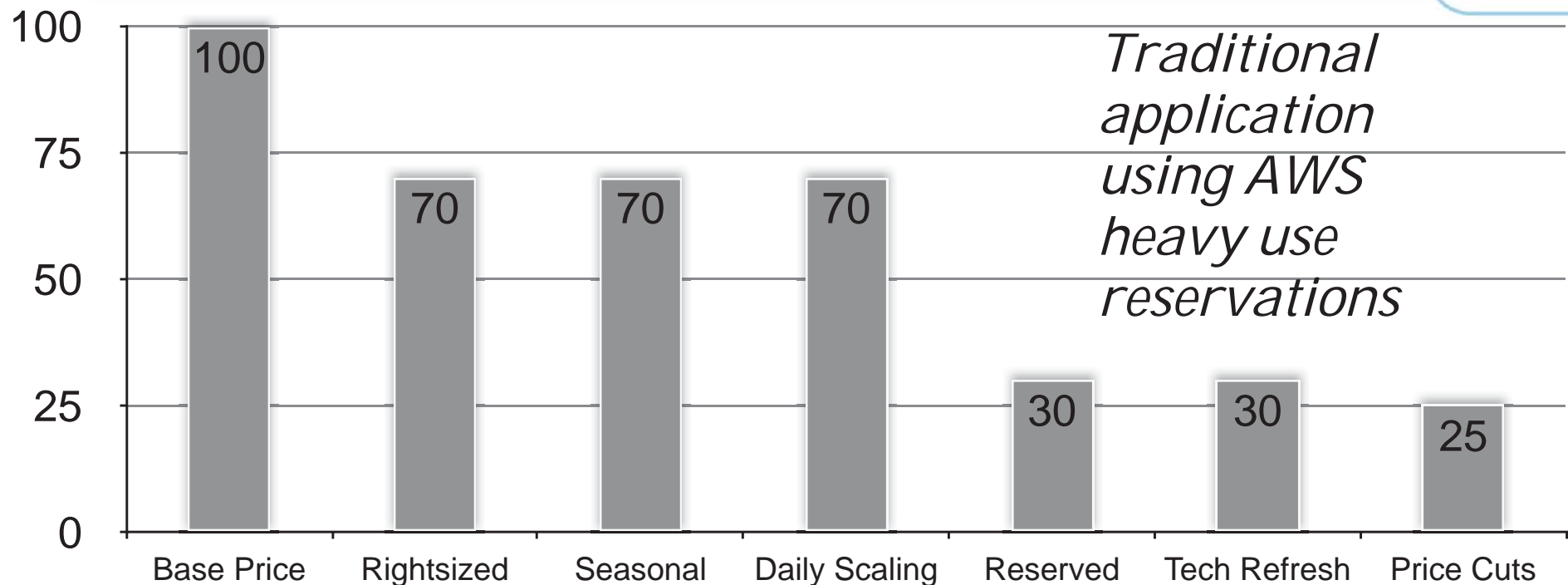
**#6 Follow the Customer (Run web servers) during the day
Follow the Money (Run Hadoop clusters) at night**

Credit to Jinesh Varia of AWS for this summary



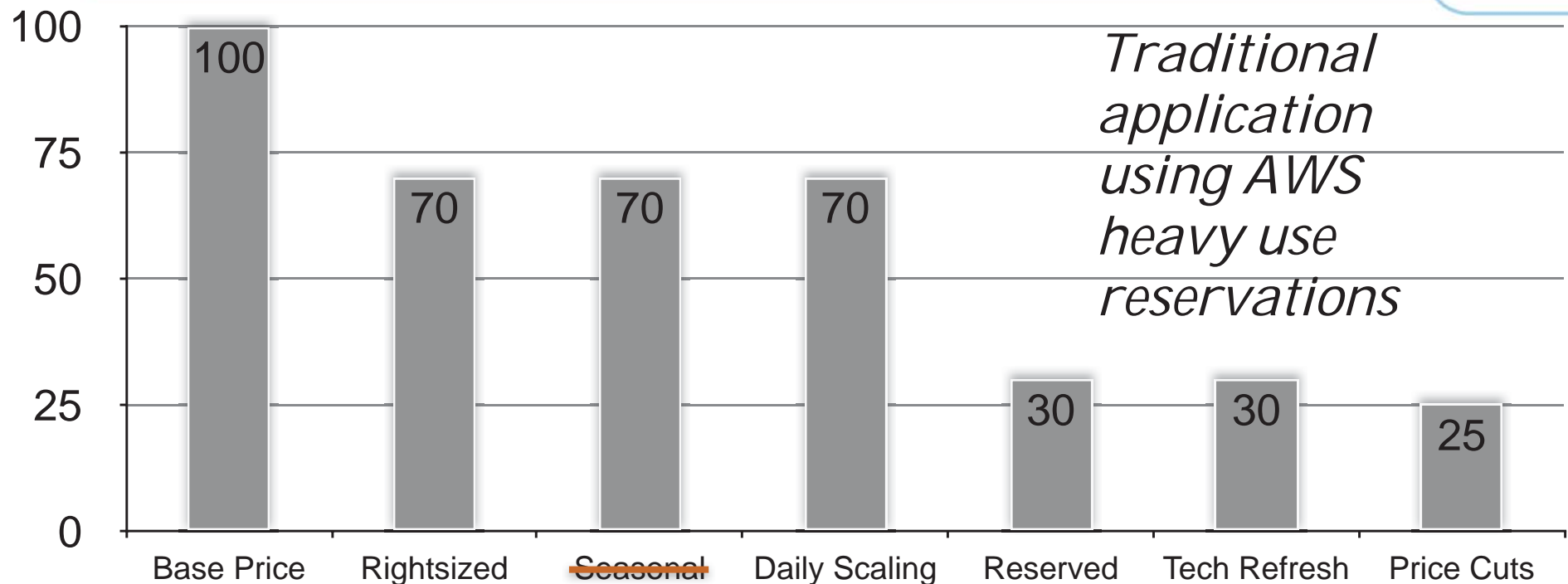
Compounded Savings

Lift and Shift Compounding



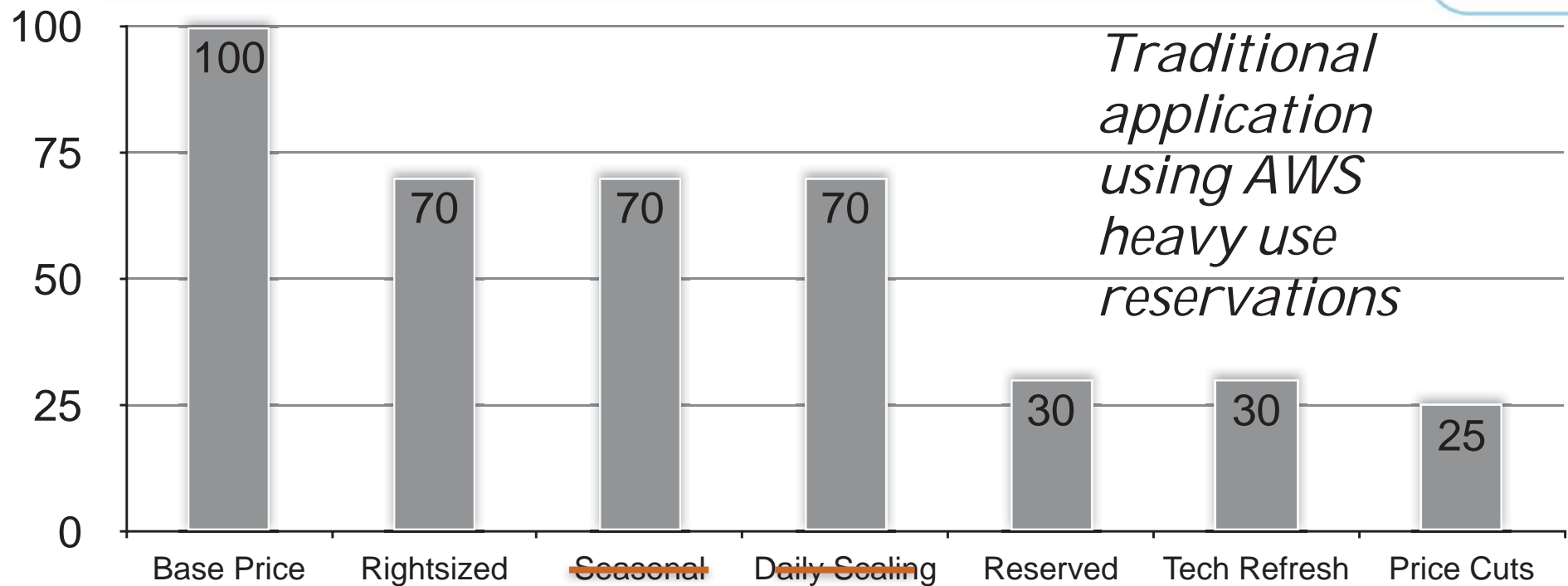
Base price is for capacity bought up-front

Lift and Shift Compounding



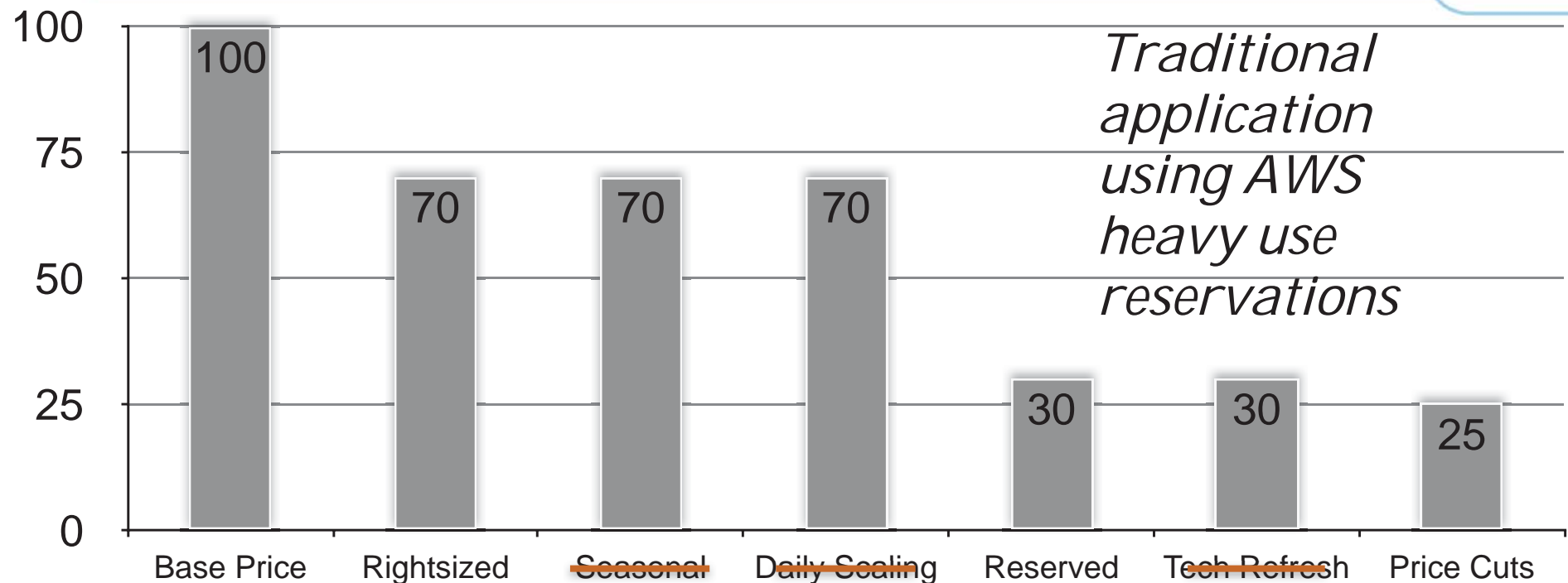
Base price is for capacity bought up-front

Lift and Shift Compounding



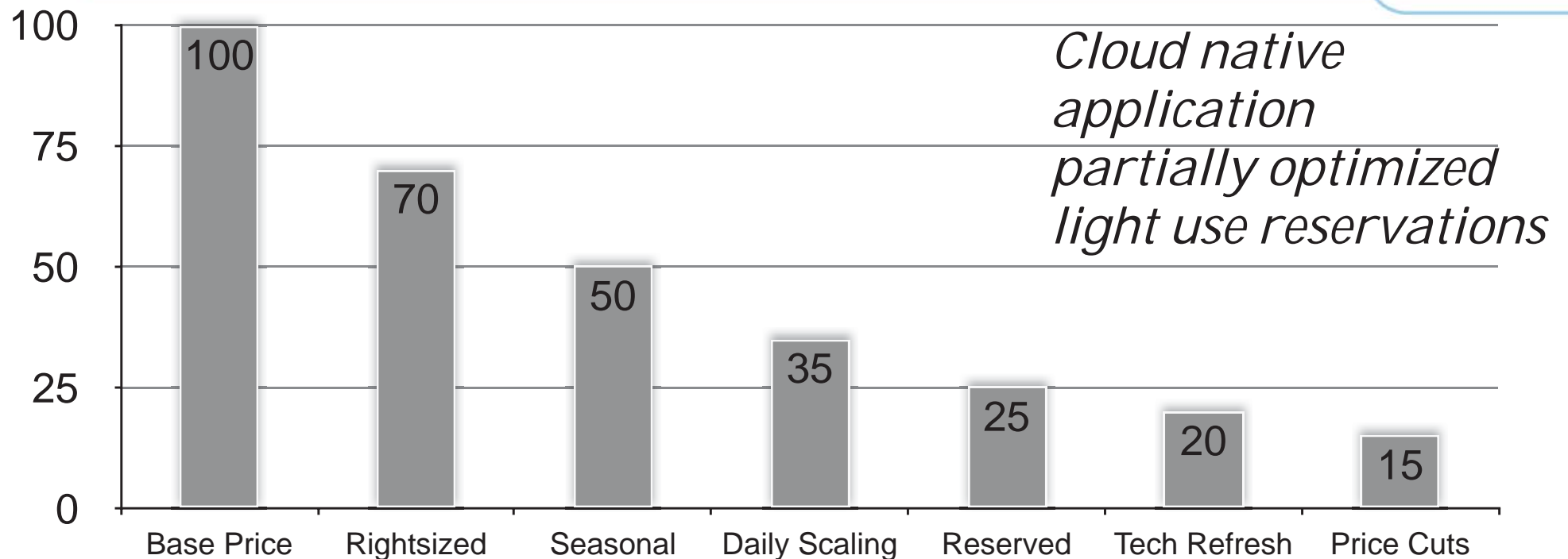
Base price is for capacity bought up-front

Lift and Shift Compounding

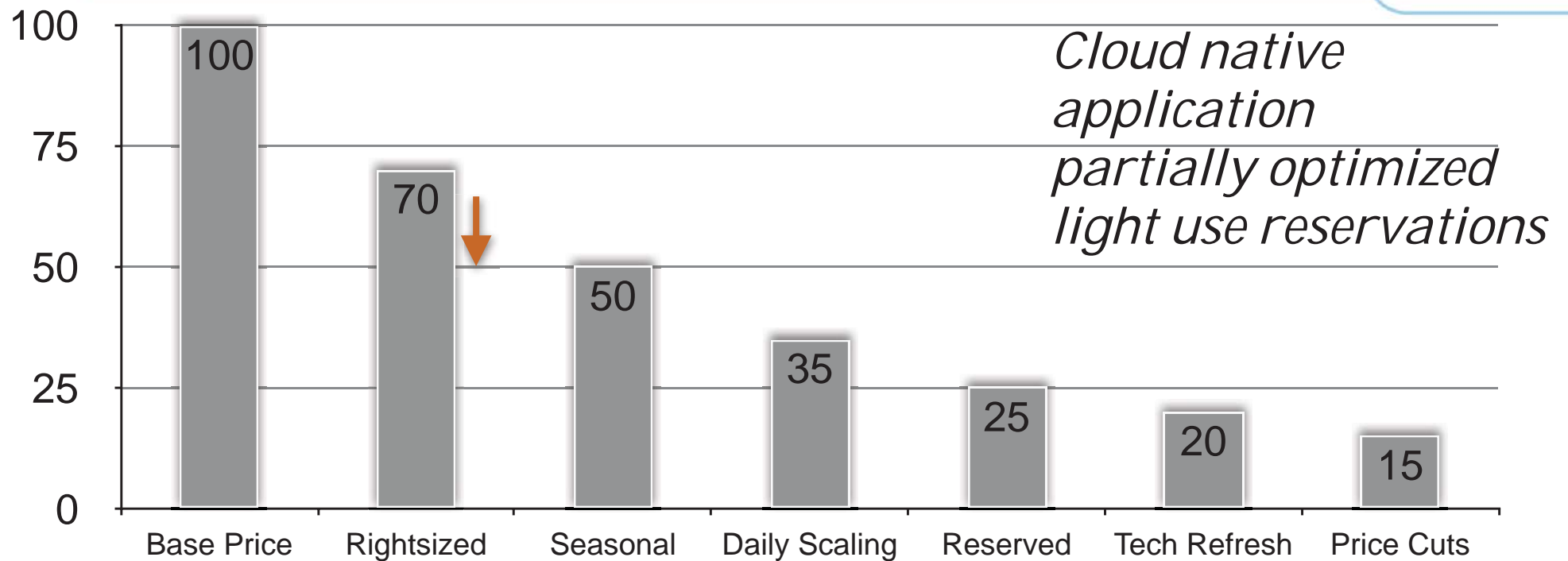


Base price is for capacity bought up-front

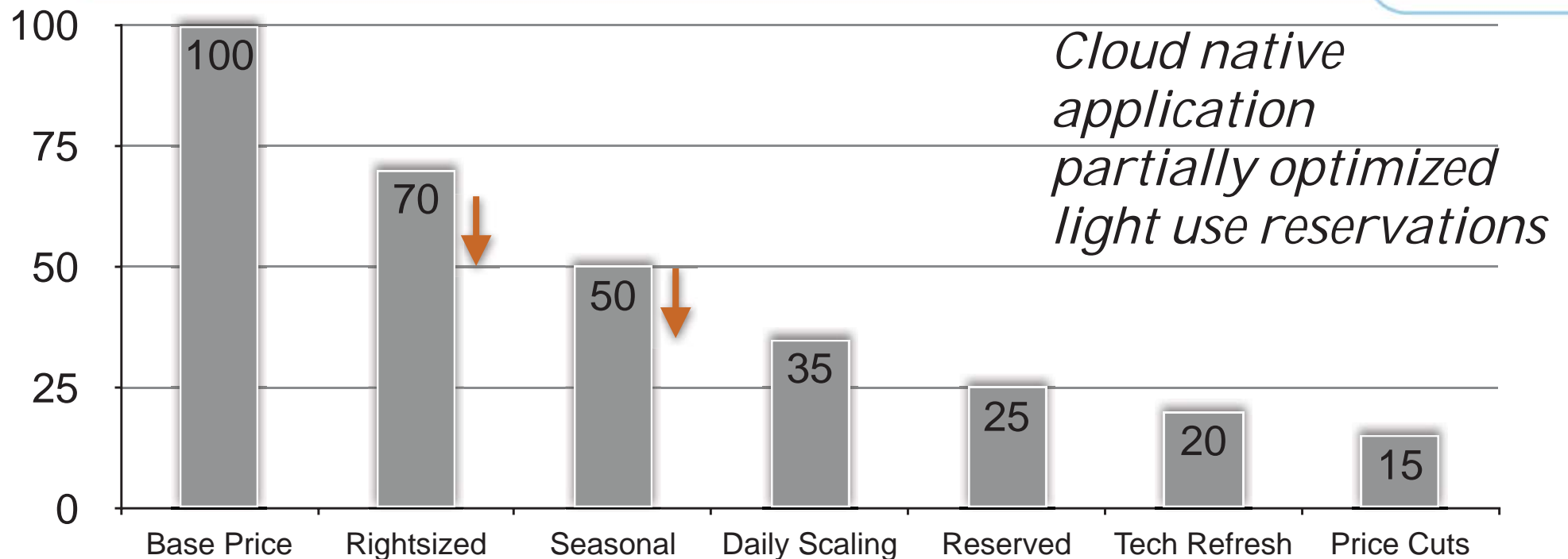
Conservative Compounding



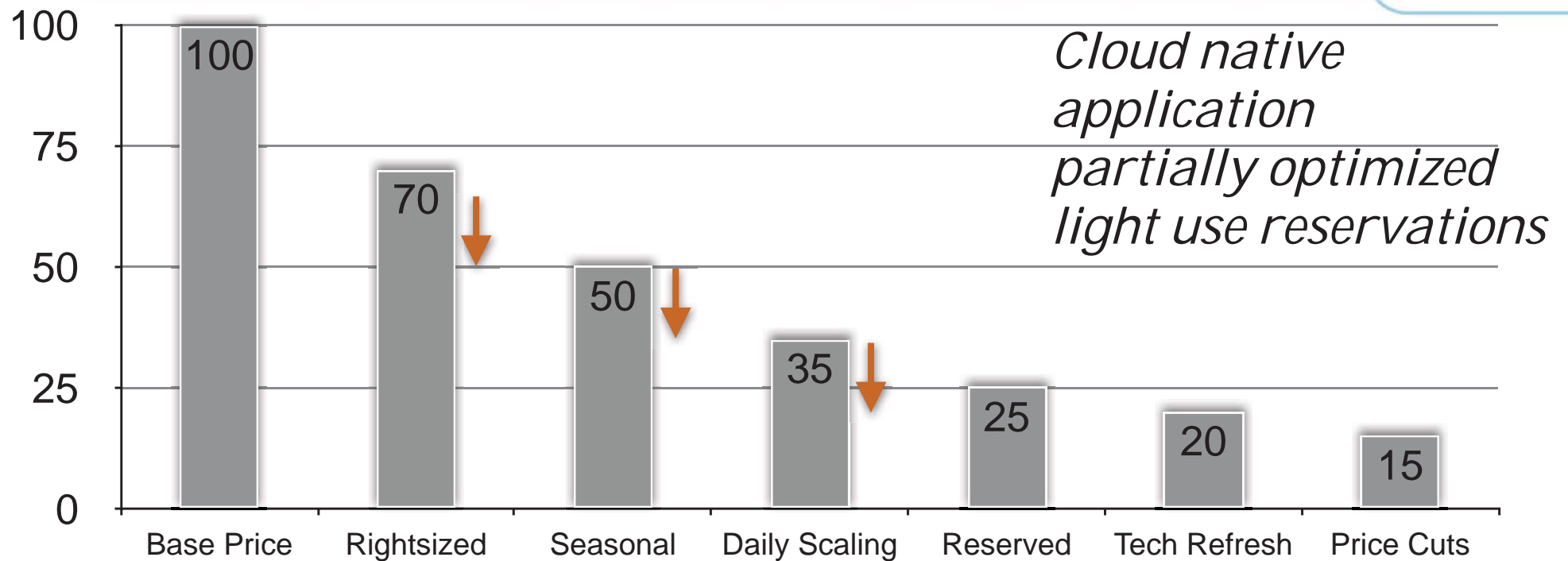
Conservative Compounding



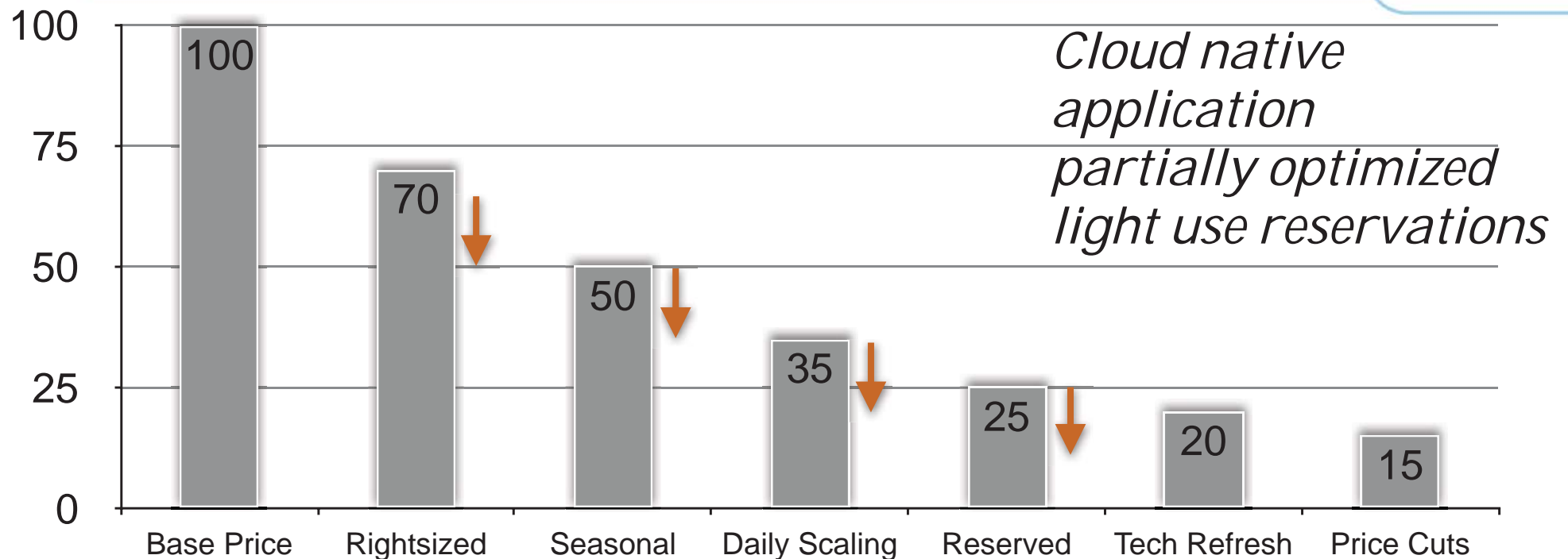
Conservative Compounding



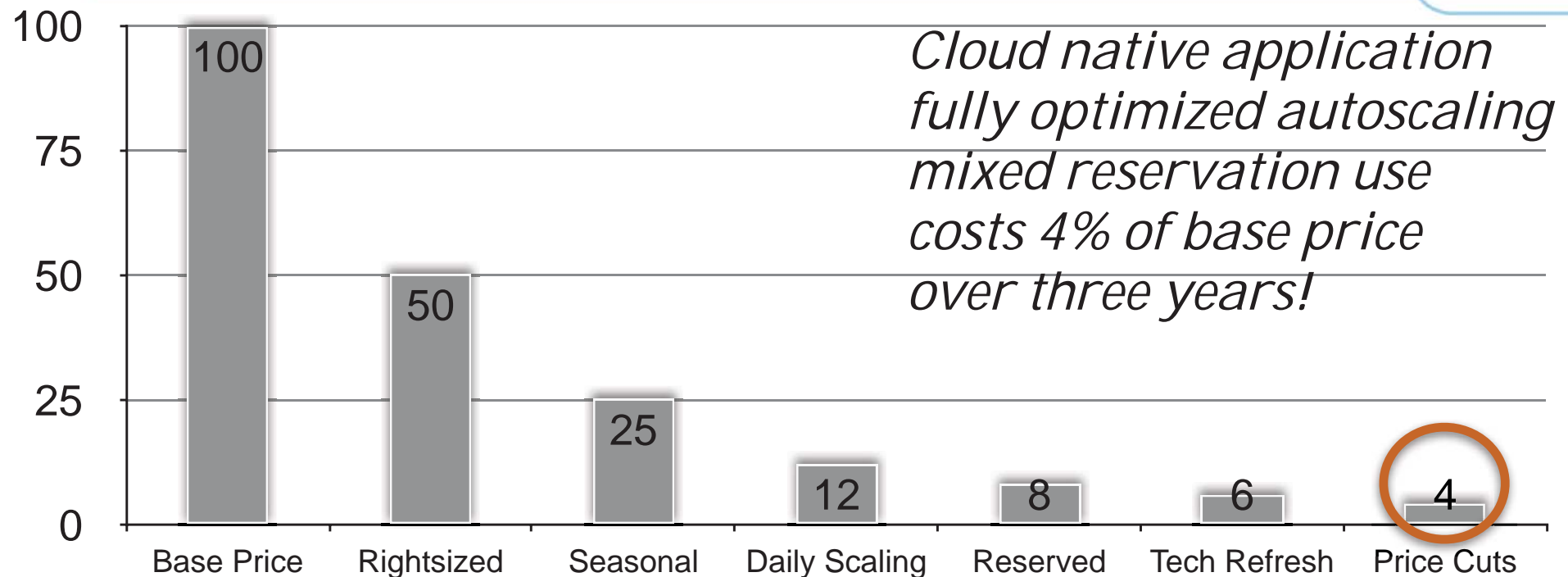
Conservative Compounding



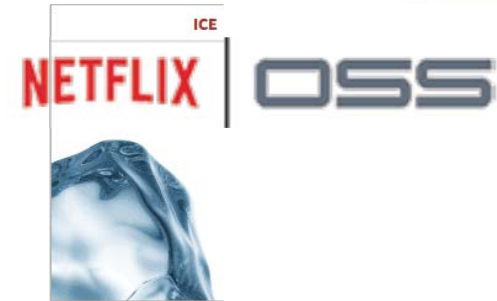
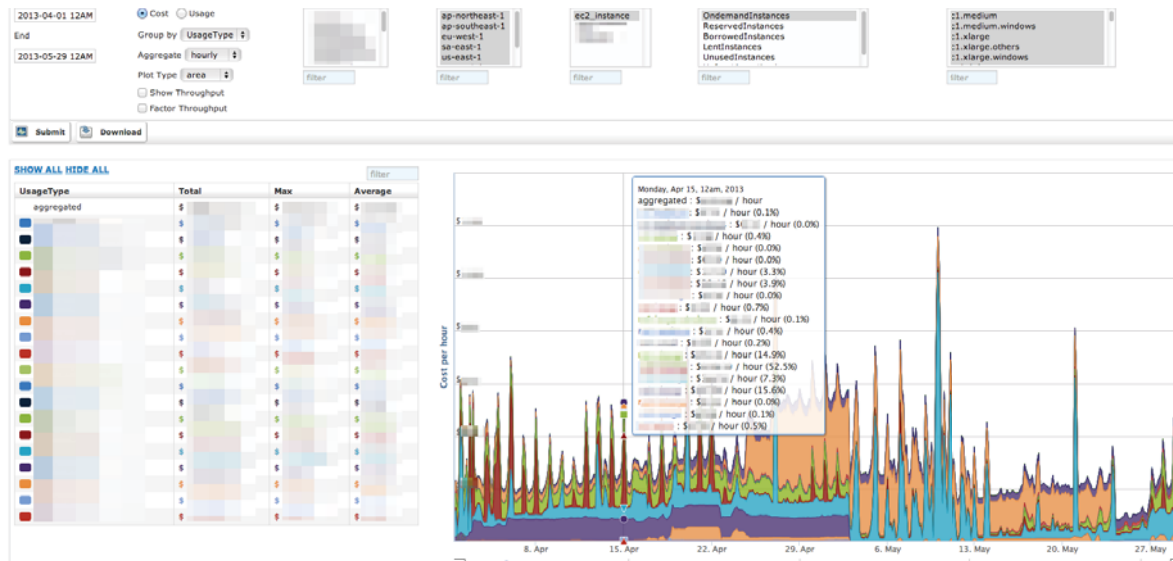
Conservative Compounding



Agressive Compounding



Cost Monitoring and Optimization



Final Thoughts



Turn off idle instances
Clean up unused stuff
Optimize for pricing model
Assume prices will go down
Go cloud native to be fast and save
Complex dynamic control issues!

Any Questions?



- Battery Ventures <http://www.battery.com>
- Adrian's Tweets @adrianco and Blog <http://perfcap.blogspot.com>
- Slideshare <http://slideshare.com/adriancockcroft>

- Monitorama Opening Keynote Portland OR - May 7th, 2014
- GOTO Chicago Opening Keynote May 20th, 2014
- Qcon New York – Speed and Scale - June 11th, 2014
- Structure - Cloud Trends - San Francisco - June 19th, 2014
- GOTO Copenhagen/Aarhus – Fast Delivery - Denmark – Sept 25th, 2014
- DevOps Enterprise Summit - San Francisco - Oct 21-23rd, 2014 #DOES14
- GOTO Berlin - Migrating to Microservices - Germany - Nov 6th, 2014
- AWS Re:Invent - Cloud Native Cost Optimization - Las Vegas - November 14th, 2014
- O'Reilly Software Architecture Conference - Fast Delivery - Boston March 16th 2015
- High Performance Transaction Systems Workshop - <http://hpts.ws> September 2015

Disclosure: some of the companies mentioned may be Battery Ventures Portfolio Companies
See www.battery.com for a list of portfolio investments

