# Cost Optimization at Scale:

## Building and Realizing the
## Economic Case for the AWS Cloud

Shahbaz Alam

AWS Professional Services

September 2016

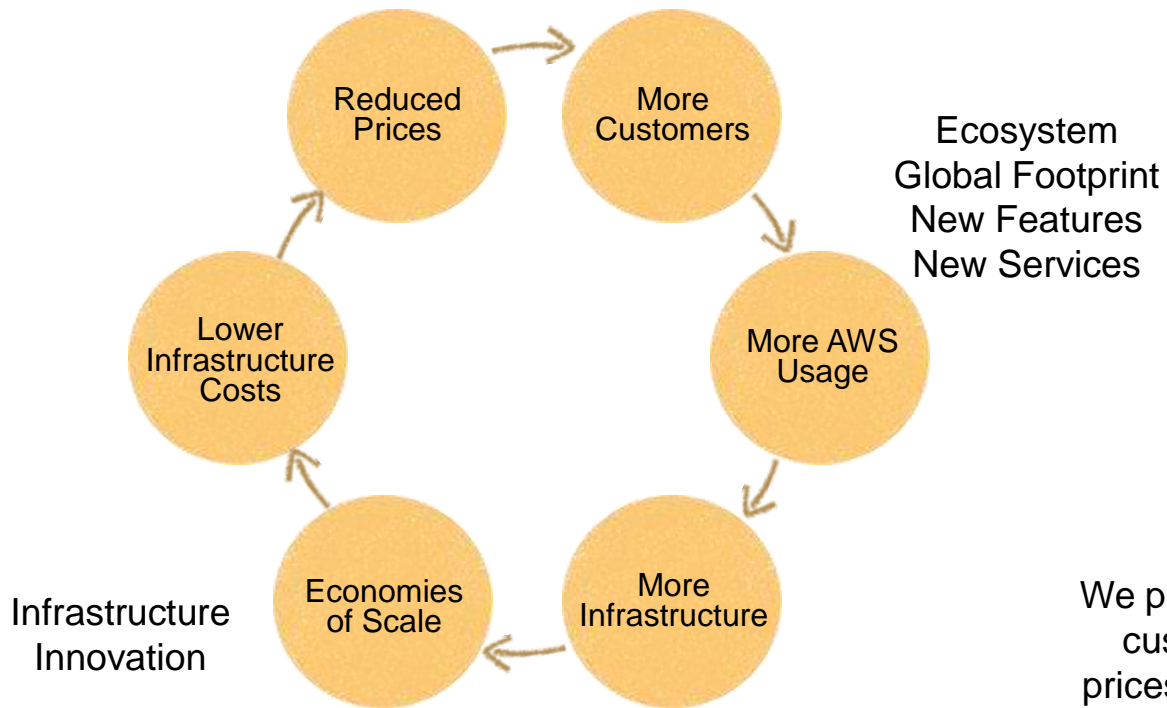# A Couple Assumptions…

1. You're using AWS…

2. You like it!!

# But maybe you are spending more than you planned…

Or you'd just like to spend less

# AWS Pricing Philosophy



Reduced Prices → More Customers → More AWS Usage → More Infrastructure → Economies of Scale → Lower Infrastructure Costs → Reduced Prices

Infrastructure Innovation

Ecosystem
Global Footprint
New Features
New Services
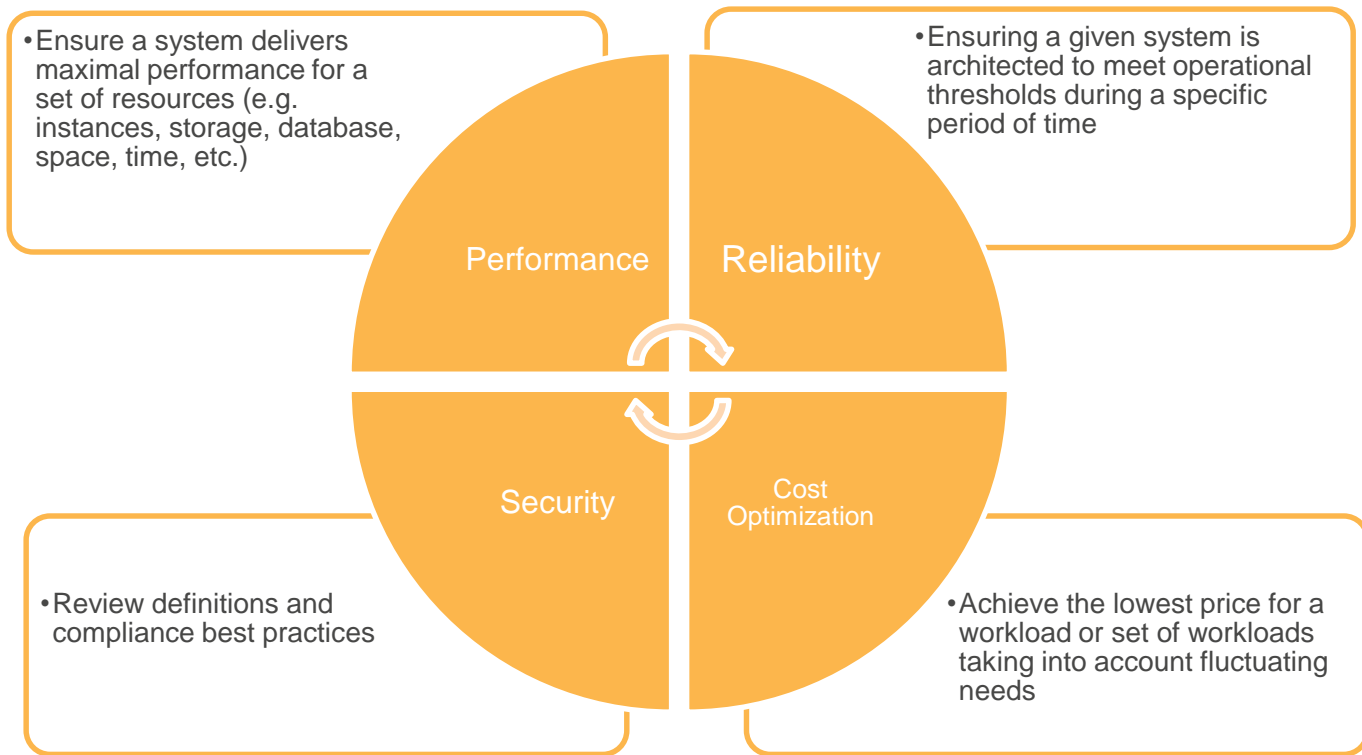
## 50+
### PRICE REDUCTIONS

We pass the savings along to our customers in the form of low prices and continuous reductions

# 4 Components of AWS Architecture Best Practices

- Ensure a system delivers maximal performance for a set of resources (e.g. instances, storage, database, space, time, etc.)

- Ensuring a given system is architected to meet operational thresholds during a specific period of time

Performance

Reliability

Security

Cost Optimization

- Review definitions and compliance best practices

- Achieve the lowest price for a workload or set of workloads taking into account fluctuating needs

In the beginning . . .

…there was **TCO**

# What is TCO?

**Definition:** *Comparative* **total cost of ownership analysis** (acquisition and operating costs) for running an infrastructure environment end-to-end on-premises vs. on AWS.

**Used for:**

1) Comparing the costs of running an **entire infrastructure environment or specific workload** on-premises or in a co-location facility vs. on AWS

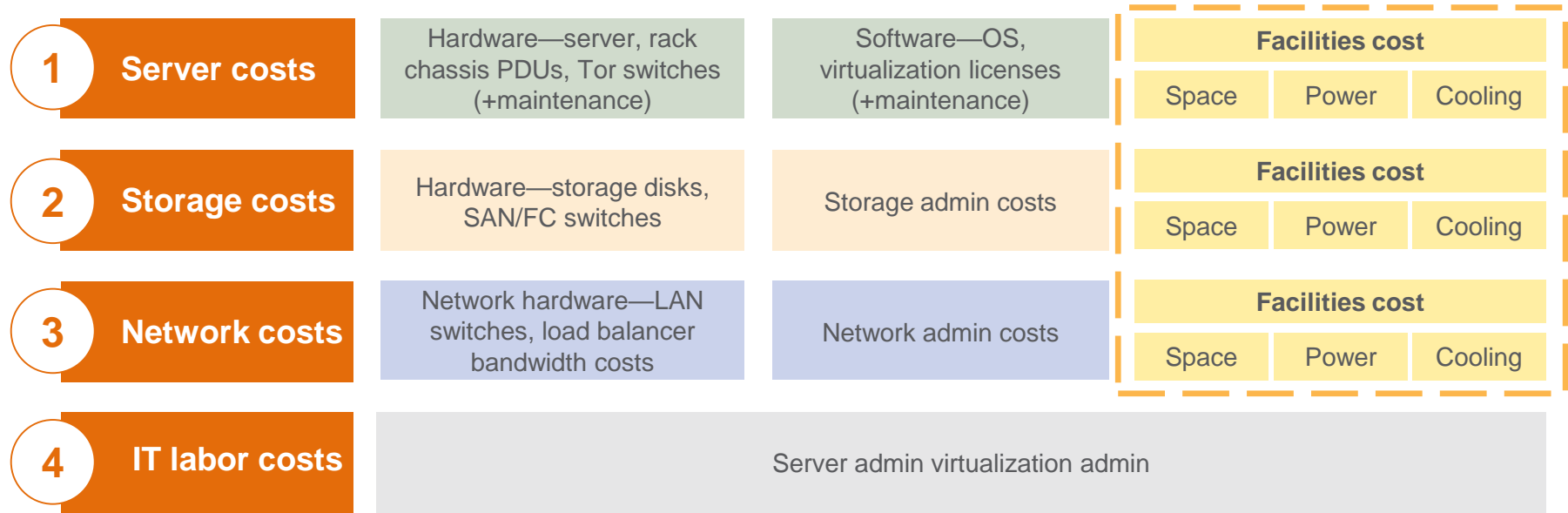2) Budgeting and **building the business case** for moving to AWS

# So how do we do it?

# TCO = Acquisition costs + Operations costs

| 1 | **Server costs** | Hardware—server, rack chassis PDUs, Tor switches (+maintenance) | Software—OS, virtualization licenses (+maintenance) | **Facilities cost** | | |
|---|---|---|---|---|---|---|
| | | | | Space | Power | Cooling |
| 2 | **Storage costs** | Hardware—storage disks, SAN/FC switches | Storage admin costs | **Facilities cost** | | |
| | | | | Space | Power | Cooling |
| 3 | **Network costs** | Network hardware—LAN switches, load balancer bandwidth costs | Network admin costs | **Facilities cost** | | |
| | | | | Space | Power | Cooling |
| 4 | **IT labor costs** | Server admin virtualization admin | | | | |

The diagram doesn't include every cost item. For example, software costs can include database, management, and middle-tier software costs. Facilities cost can include costs associated with upgrades, maintenance, building security, taxes, and so on. IT labor costs can include security admin and application admin costs.

# Questions to explore your existing footprint…

**1** **Capacity Planning**
- How do you plan for capacity?
- How many servers have you added in the past year? Anticipating next year?
- Can you switch your hardware on and off and only pay for what is used?

**2** **Utilization**
- What is your average server utilization?
- How much do you <u>overprovision</u> for peak load?

**3** **Operations**
- Will you run out of data center space some time in the future?
- What was your last year power utility bill for the Data Center(s)?
- Have you budgeted for both <u>average</u> and <u>peak power</u> requirements?

**4** **Optimization**
- Are you on AWS today?
- Is your architecture cost-optimized (Auto Scaling, RIs, Spot, Instances turn on/off)?

# And, make sure to…

**Consider**

- Power/Cooling (compute, storage, shared network)
- Data Center Administration (procurement, design, build, operate, network, security personnel)
- Rent/Real Estate (building deprecation, taxes)
- Software (OS, virtualization licensing & maintenance)
- RAW vs. USABLE storage capacity
- Storage Redundancy (RAID penalty, OS penalty)
- Storage Backup costs (tape, backup software)
- Bandwidth, Network Gear & Redundancy (routers, VPN, WAN, etc.)

**Understand**

- Procurement Time, Resource sitting on self
- Cost of Lost Customers
- RTO, RPO

# Resources to get you started



AWS TCO Calculator

https://awstcocalculator.com

Case studies and research

http://aws.amazon.com/economics/

# Lowering TCO through cost optimization



Economic Case Improves through Optimization

On-Premises | Lift & Shift | Instance Right-Sizing | Improved Elasticity | Measure, Monitor, Improve | Optimized EC2 | Storage Optimization | Serverless Architecture | Managed Services | True AWS Optimized

Traditional TCO Comparisons

So you're feeling pretty good.

**Until your CFO shows up with the bill.**

# Cost optimization is…



going from…

pay for what you *use*

to…

pay for what you *need*

# Key inputs to cost optimization on AWS

Where do you start?

# The four pillars of cost optimization

Right-sizing

Reserved Instances

Increase elasticity

Measure, monitor, and improve

# Right-sizing



## Right-sizing

- Selecting the cheapest instance available while meeting performance requirements

- Looking at CPU, RAM, storage, and network utilization to identify potential instances that can be downsized

- Leveraging Amazon CloudWatch metrics and setting up custom RAM metrics

Rule of thumb: Right size, then reserve.

*(But if you're in a pinch, reserve first.)*

# Reserved Instances

### Step 1: RI Coverage

- Cover always-on resources.

### Step 2: RI Utilization

- Leverage RI flexibility to increase utilization.
- Merge and split RIs as needed.

Rule of thumb: Target 70–80% always-on coverage and 95% RI utilization rate.

# EC2 Reserved Pricing

Steady State

Reserved Capacity

Upfront payments to reduce costs

# Reserved Instances

## Up to 75%+ savings* (and capacity reservation)

**Commitment level**
1 year
3 year

**AWS services offering RIs**
Amazon EC2
Amazon RDS
Amazon DynamoDB
Amazon Redshift
Amazon ElastiCache

* Dependent on specific AWS service, size/type, and region

# Increase elasticity



## Turn off nonproduction instances

- Look for dev/test, nonproduction instances that are running always-on and turn them off.

## Autoscale production

- Use Auto Scaling to scale up and down based on demand and usage (for example, spikes).

Rule of thumb: Shoot for 20–30% of Amazon EC2 instances running on demand to be able to handle elasticity needs.

# Using right-sizing and elasticity to lower cost

More smaller instances vs. fewer larger instances



29 m4.large @ $0.12 /hr
**$2,505.60 / mo***

59 t2.medium @ $0.052/hr
**$2,208.96 / mo***

*Assumes Amazon Linux instances in the US-East (N. Virginia) Region at 720 hours per month

# EC2 Spot Pricing

Time or instance flexible

Experiment and/or build cost sensitive businesses

Users with urgent computing needs or large amounts of additional capacity

# Consider Spot for Elastic Workloads

**90% Savings!***

## Options

- Spot Fleet to maintain instance availability
- Spot Block durations (1-6 hours) for workloads that must run continuously

## Commitment level

- None

* Compared to On Demand price based on specific EC2 instance type, region, and Availability Zone

# Spot Rules

Markets where the price of compute changes based on supply and demand

You'll never pay more than your bid.

## Spot Instance Pricing History ✕

Product : Linux/UNIX (Amazon VPC) ▾   Instance type: r3.4xlarge ▾   Date range : 1 week ▾   Availability zone: us-east-1d ▾

*75% of OD*

*50% of OD*

*25% of OD*

| Availability zone | Price |
|---|---|
| █ us-east-1d | $0.1788 |
| Date | September 17, 2015 at 10:33:37 PM UTC-7 |

You pay the market price **87% discount!**

# Strike a Balance

**Reserved Instances**

**On Demand**

**Spot**

Finding balance between pricing options

# Consumption model by industry

### Web Scale (e.g. Adtech) Company



Normalized Usage by Pricing Model ±

SP

OD

Mar 1, 15   Jun 1, 15   Sep 1, 15   Dec 1, 15   Mar 1, 16

### Enterprise SaaS Company



Normalized Usage by Pricing Model ±

SP

OD

RI

Mar 1, 15   Jun 1, 15   Sep 1, 15   Dec 1, 15   Mar 1, 16

# Consumption model by industry (cont…)

**Onboarding Enterprise**



Normalized Usage by Pricing Model ±

**Gaming Company**



Normalized Usage by Pricing Model ±

# Consumption model workload…

### Dev Test



### Enterprise Applications



### Data Science



### New app development

# EC2 cost optimization options

| Cost Savings | EC2 | Benefit |
|---|---|---|
| **Base Price** | On Demand | • No Commitment<br>• Pay only what you use<br>• No capacity reservation<br>• No interruption |
| **< 10%** | Scheduled Reserved Instances | • Commitment of 1,200 hours for one year<br>• Specified schedule<br>• Capacity reservation; no interruption |
| **30% – 75%** | Standard Reserved Instances | • Commitment of one year or three years<br>• Capacity reservation; no interruption |
| **40% – 60%** | Spot Blocks | • Bid for 2-6 hours blocks of time<br>• No long term commitment<br>• No interruption |
| **Up to 85%** | Spot | • Bid for instances<br>• Interrupted if market price higher than your bid price<br>• 2 minute advanced notice |

# Putting it all together: case study
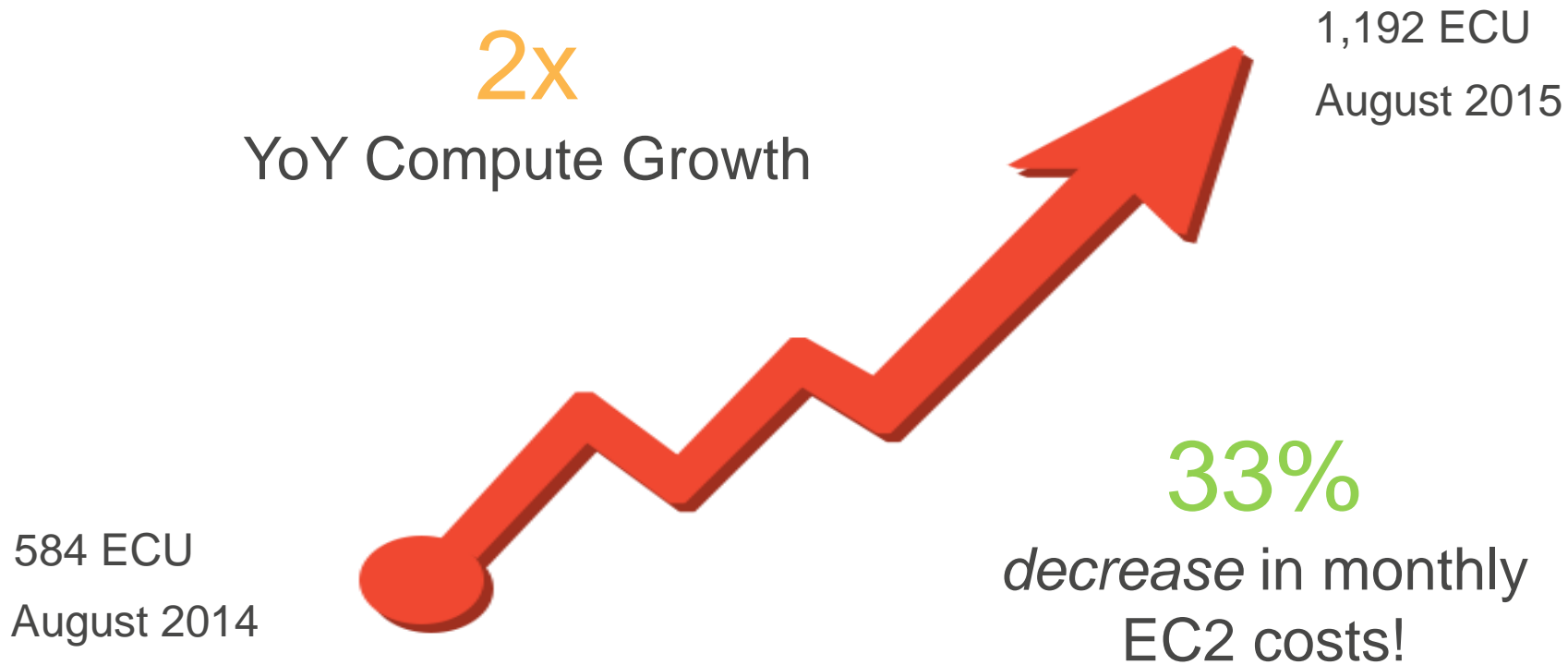
# Challenge:

**Minimizing *unit costs* during a period of massive growth.**

**Elastic compute unit (ECU)**
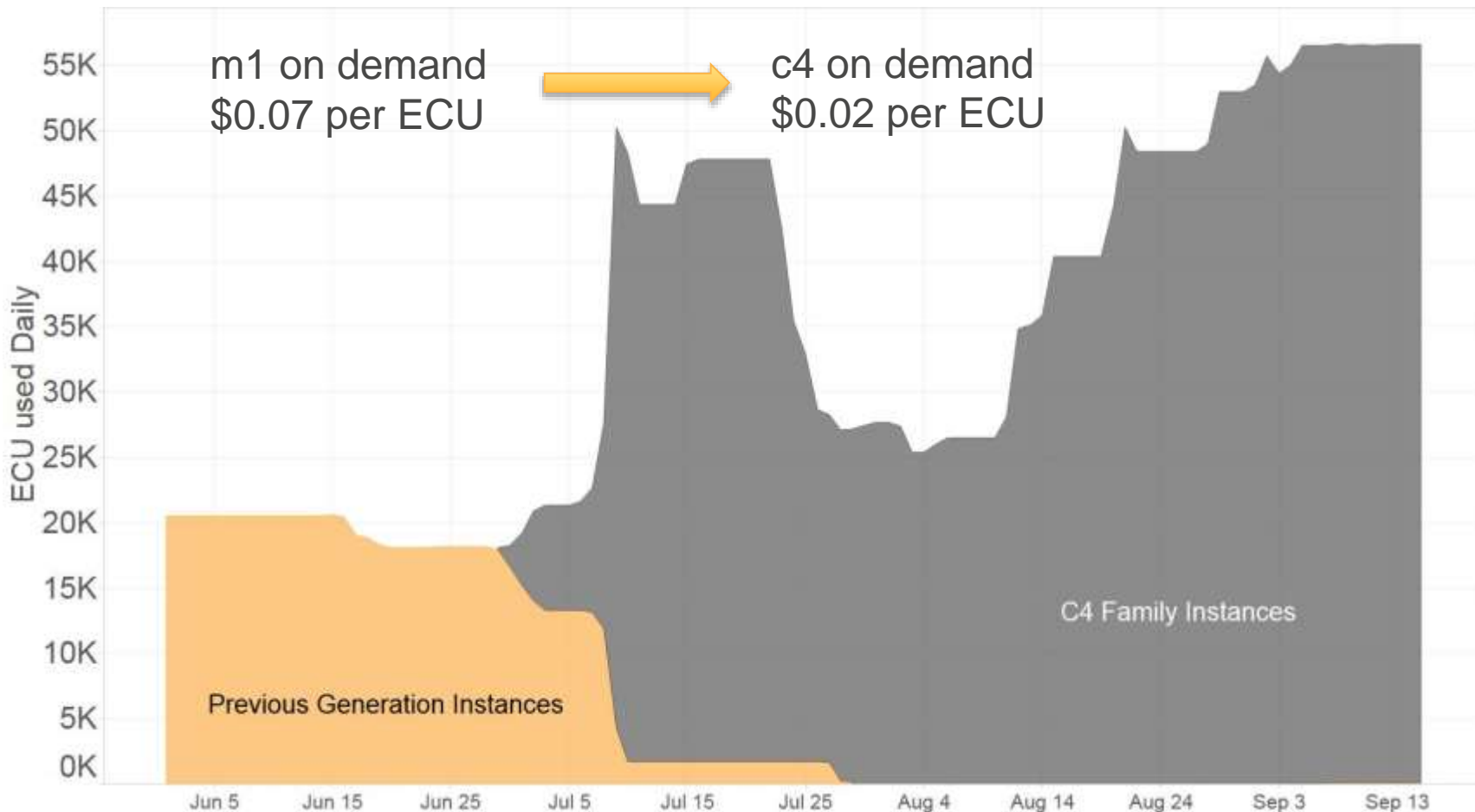


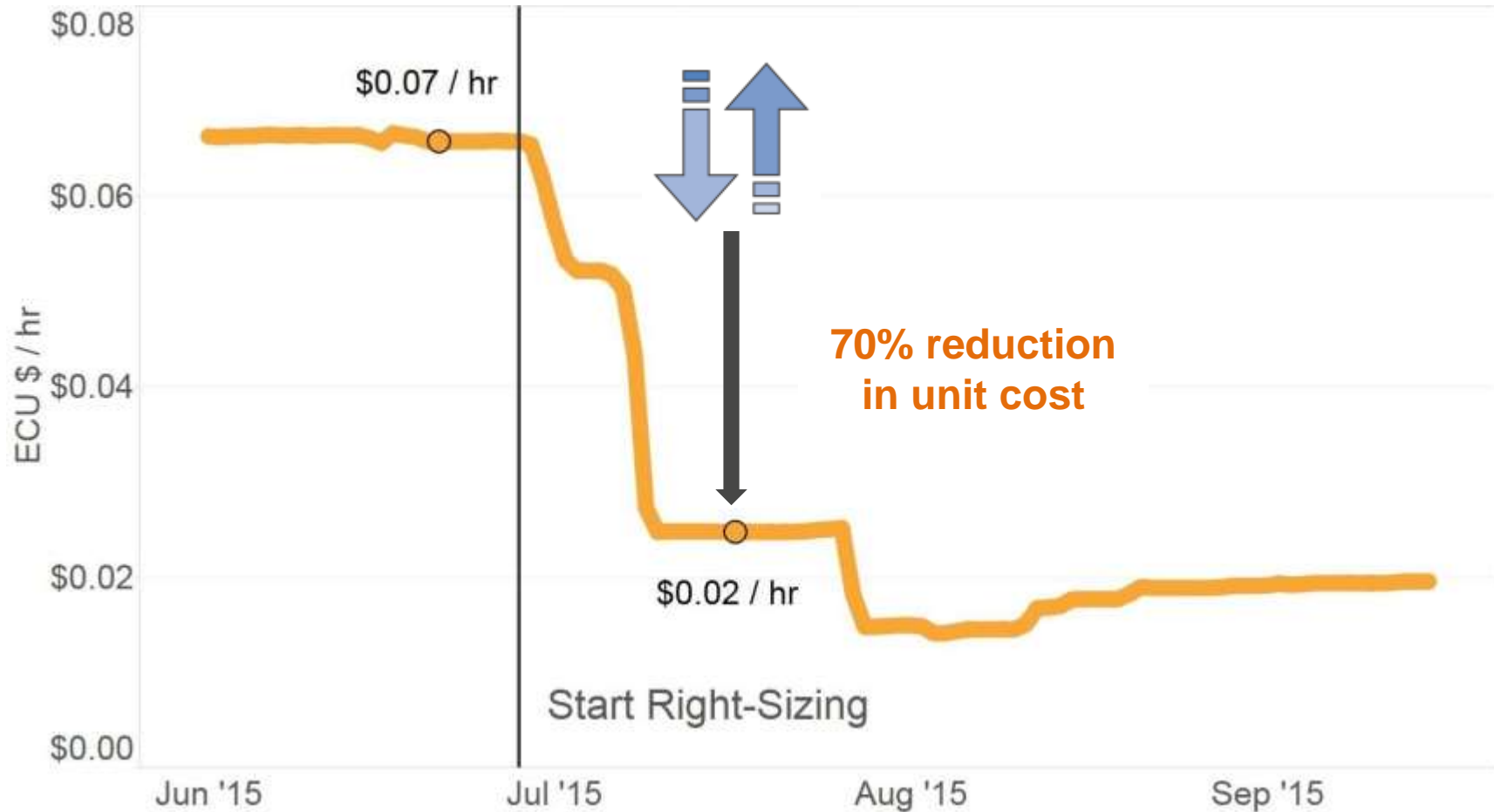A consistent measure of CPU processing power

# The growth challenge

**2x**

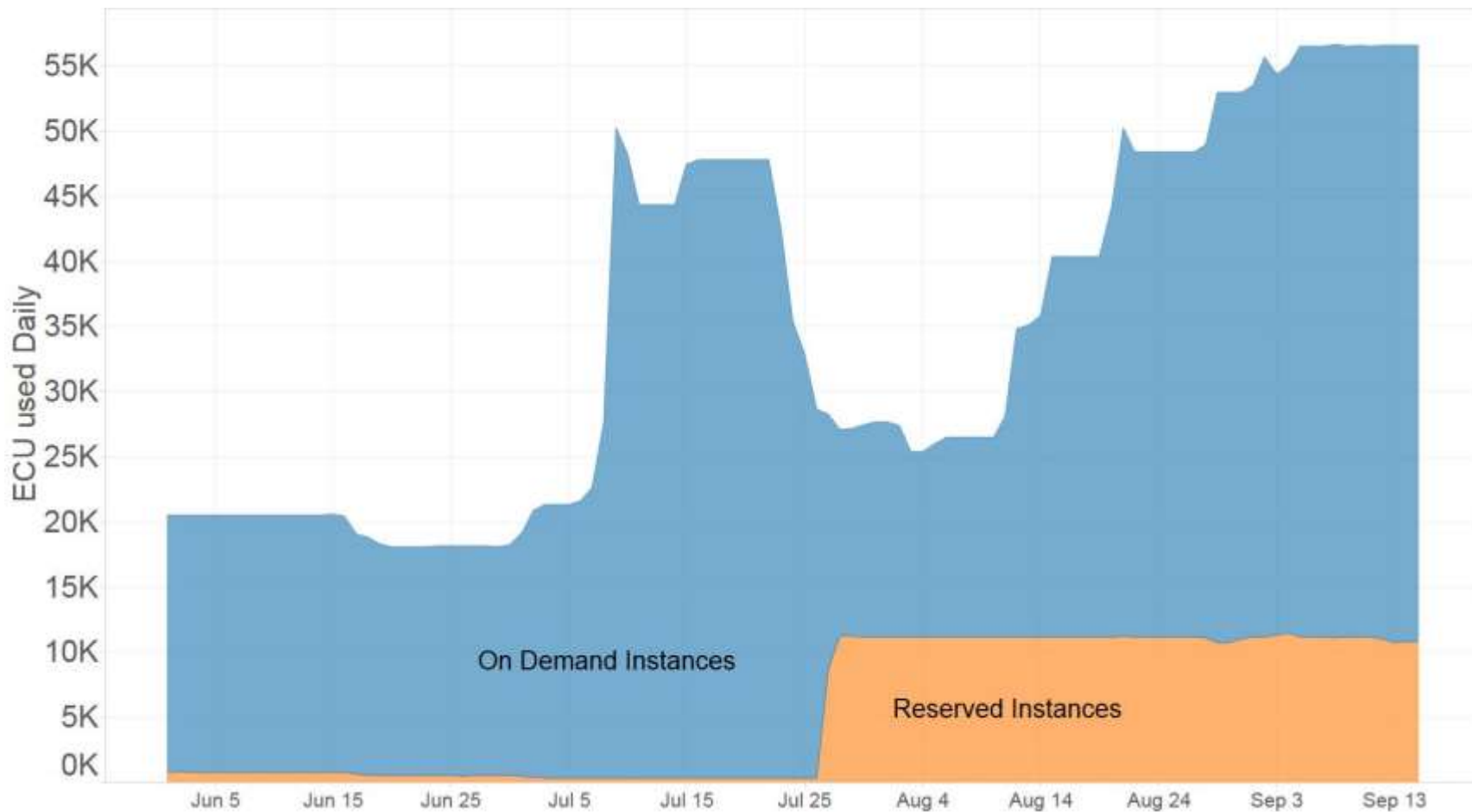YoY Compute Growth

1,192 ECU

August 2015

584 ECU

August 2014

**33%**

*decrease* in monthly
EC2 costs!

# Solving the growth challenge

# Step 1: Right-size and update instances



m1 on demand
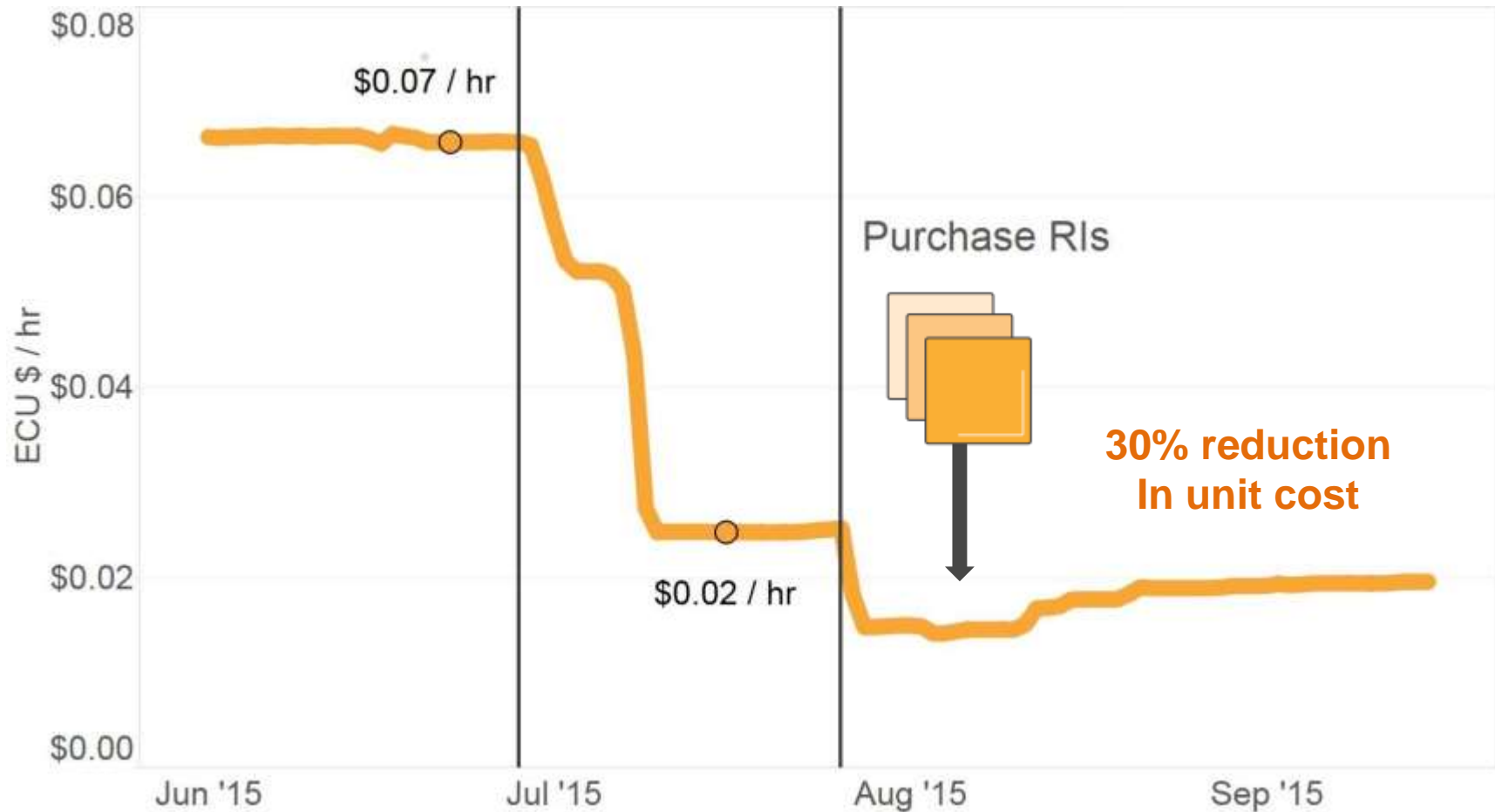$0.07 per ECU

→

c4 on demand
$0.02 per ECU

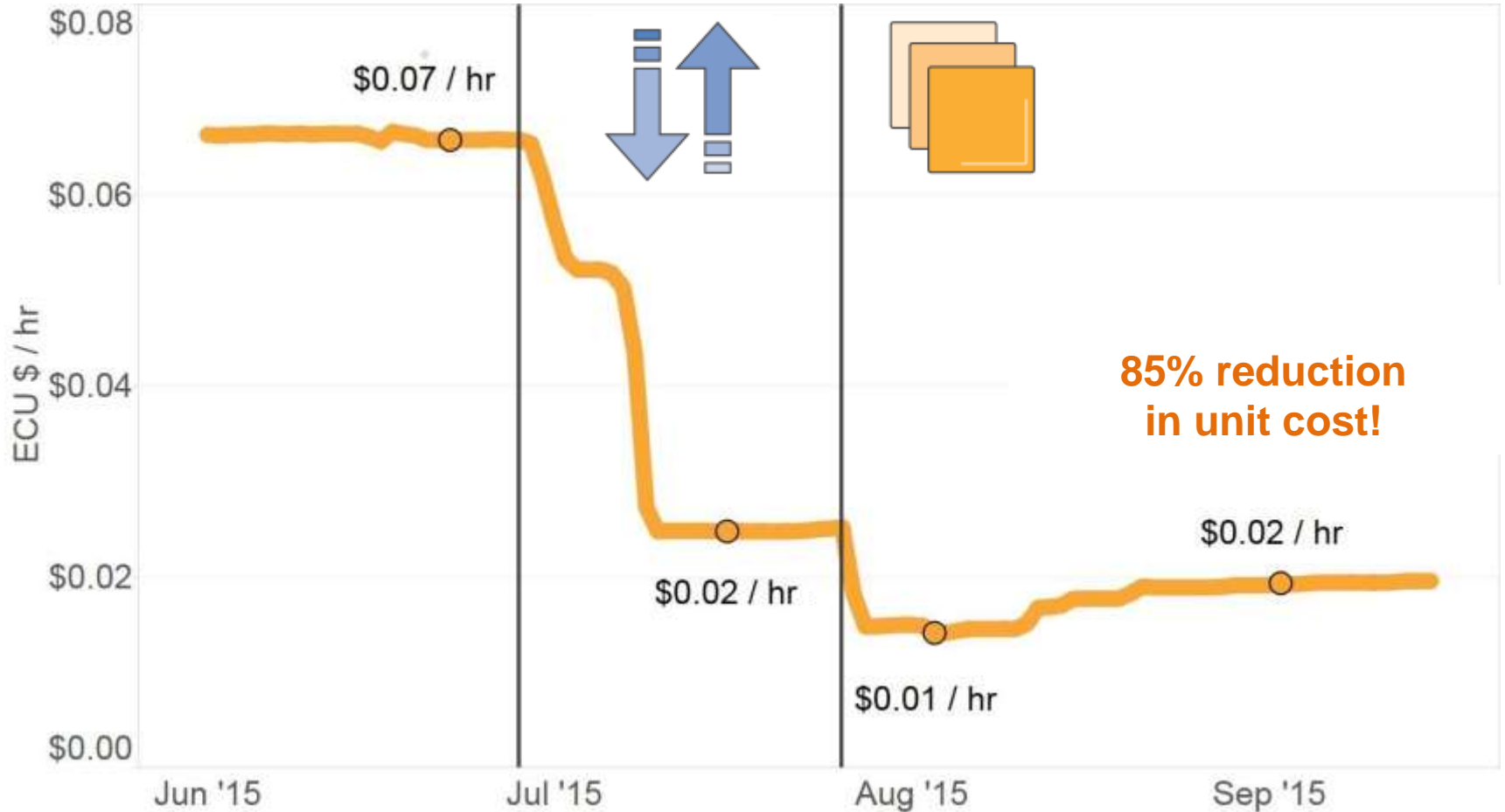# The impact of right-sizing

# Step 2: Reserve

# The impact of reservations

# Putting it together

# Sounds pretty easy, right?

**Not really.**

**In reality, it is very complex.**

- Scale
- Behavioral change
- Visibility
- Ownership

**Cost optimization governance (Remember the fourth pillar?)**
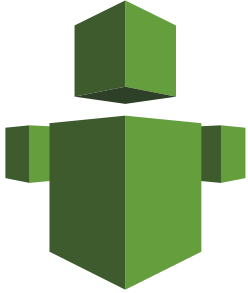
# Uncovering the cost optimization opportunities

1. Auto-tag resources.
2. Identify always-on nonprod.
3. Identify instances to downsize.
4. Recommend RIs to purchase.
5. Dashboard our status.
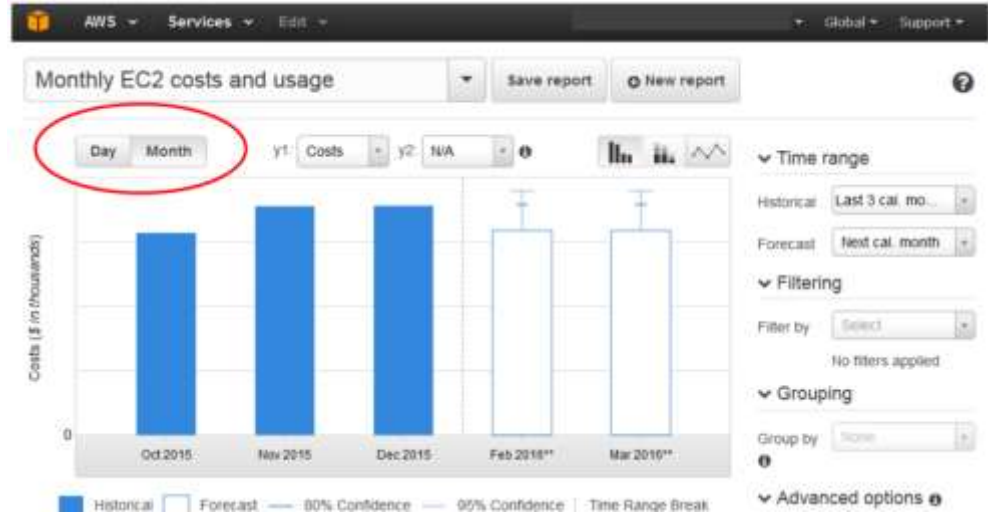6. Report on savings.

# AWS options



AWS Trusted Advisor

Cost Explorer

# Reserved Instances and right-sizing options

# Example: reasonable optimization dashboard

# Creating a culture of cost transparency

Targets and metrics

Cloud Competency Center

AWS Enterprise Support

# Cost Metrics

A company's overall AWS cost should be evaluated as a unit cost ratio with respect to another defined metric:

$$Unit\ Cost\ = \frac{Total\ Cost}{Individual\ or\ Business\ Metric}$$

Examples

- Unit cost per revenue generated
- Unit cost per product or business unit
- Unit cost per internal user
- Unit cost per customer or subscriber

# Putting it all together

# Where to start

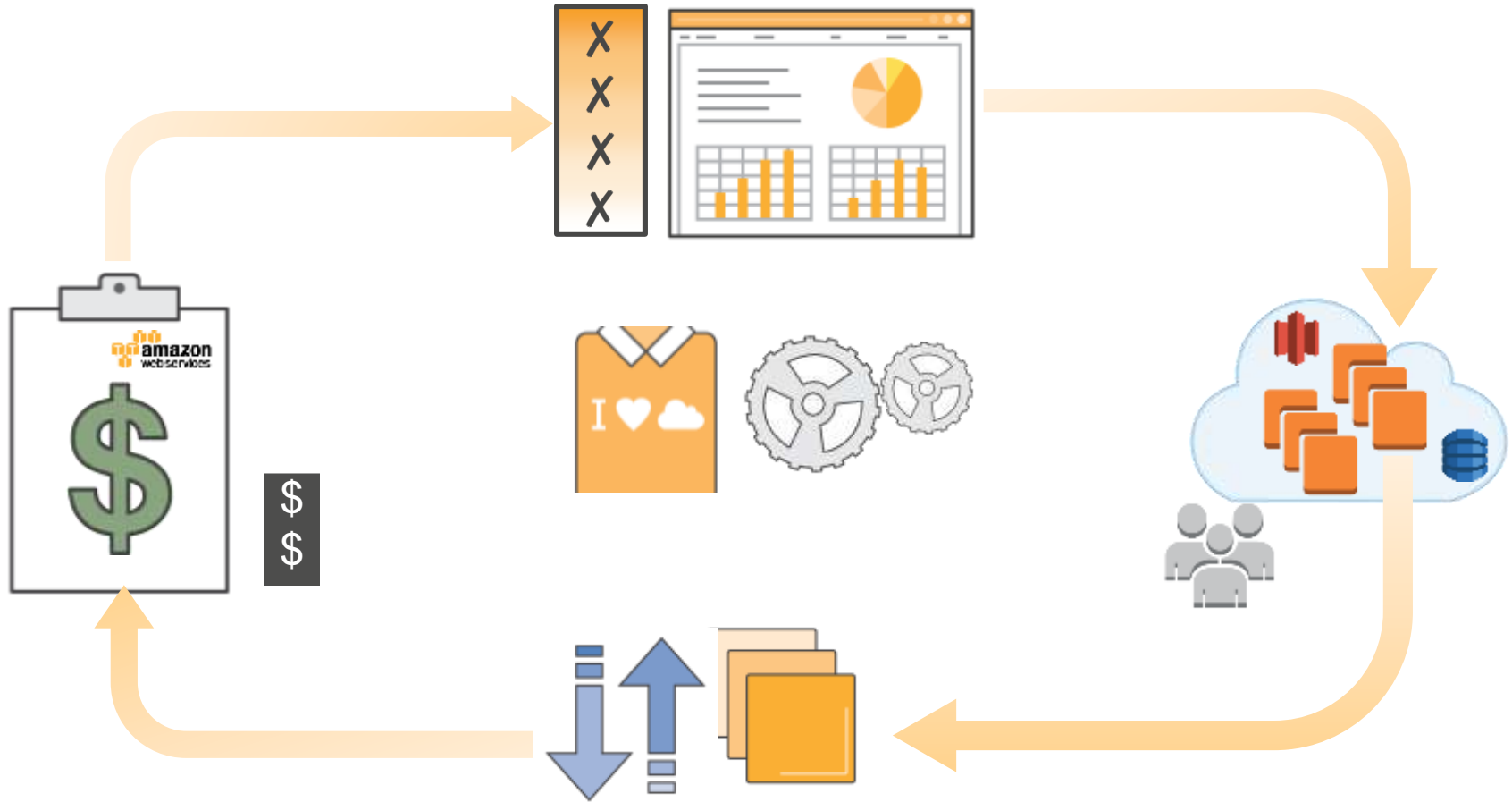Set up a Cloud
Competency Center

Bring in the right
tools

Use metrics to
reinforce behavior

Use partners to
accelerate!

# Cycle of cost optimization

**Remember to complete your evaluations!**