
Software Scalability

CS480 Software Engineering

Yu Sun, Ph.D.

<http://yusun.io>

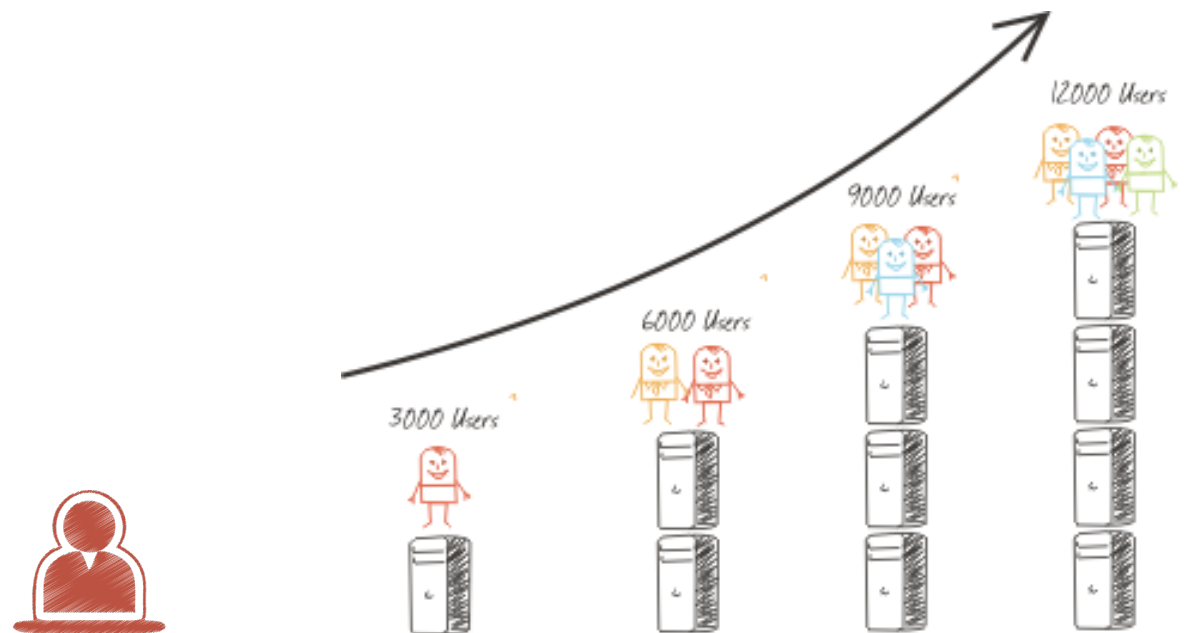
yusun@cpp.edu



CAL POLY POMONA

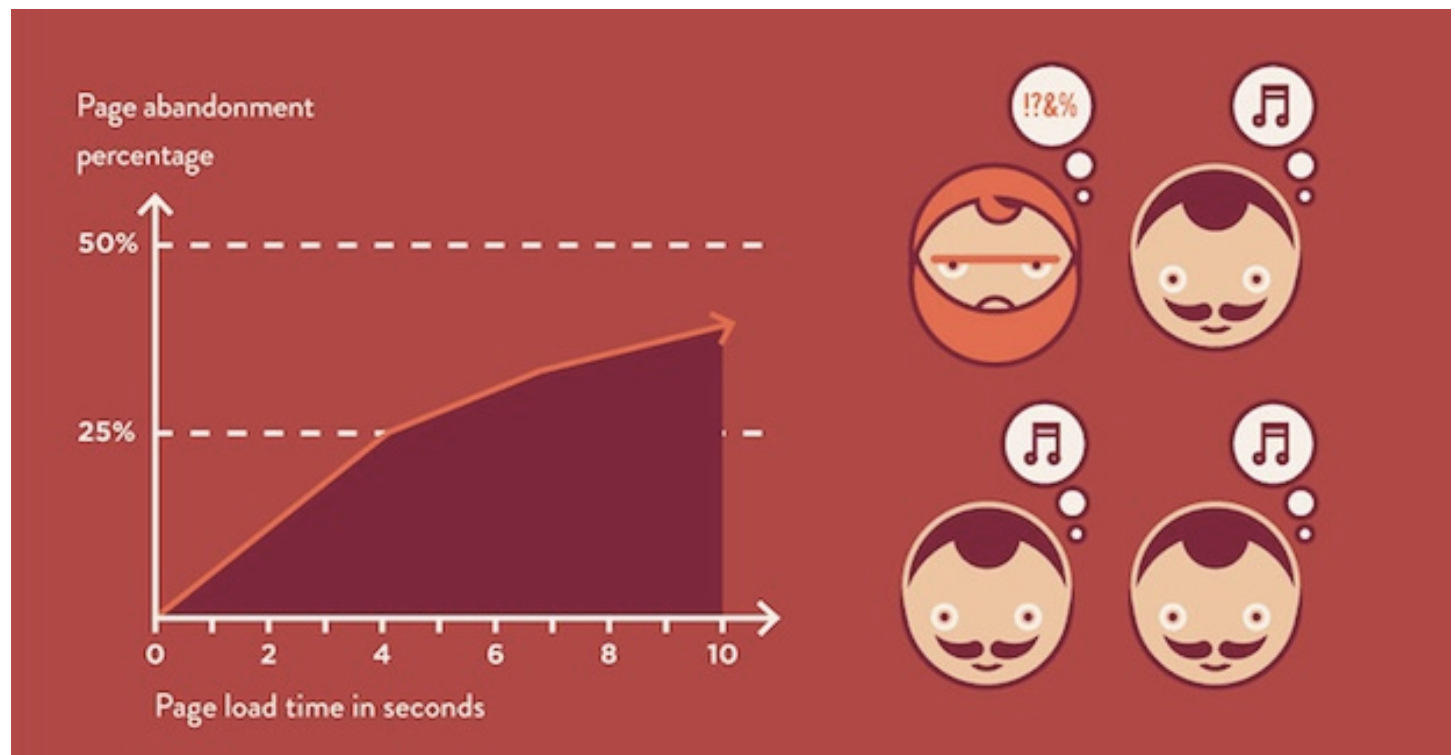
Software Scalability

- ◆ **Scalability** is the ability of a system to handle a growing amount of work in a capable manner or its ability to be enlarged to accommodate that growth



Amazon.com

- ◆ 426 items were sold per second during Christmas
- ◆ A page load slowdown of just one second could cost it \$1.6 billion in sales each year



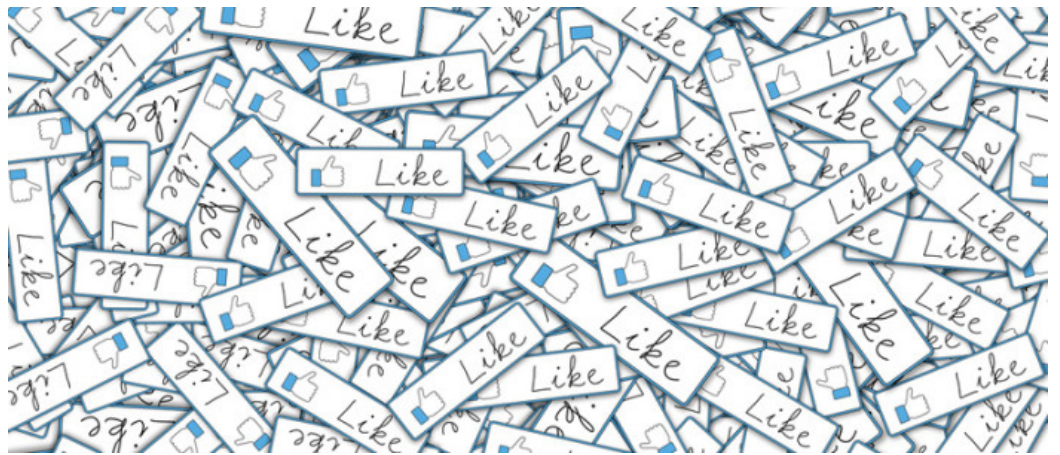
Google

- ◆ 3.5 billion searches / day
- ◆ 1.2 trillion searches / year
- ◆ “by slowing its search results by just 4/10 of a second they could lose 8 million searches per day”



Facebook

- ◆ People spend over 700 billion minutes per month on Facebook
- ◆ In 20 minutes 10.2 million comments are posted
- ◆ 750 million photos were uploaded to Facebook over New Year's weekend



Amazon S3

- ◆ 1.3 trillion objects stored
- ◆ 1.1 million requests / second



Uber Growth



BroncoDirect



Cal Poly Pomona

February 10 at 11:47am · 🌐

Broncos, if you're trying to register for classes right now, you've probably noticed that BroncoDirect is having some serious issues right now.

We apologize for this. Really, we do. We know that you don't need this extra aggravation while you're trying to get your classes for next quarter.

We are working as quickly as we can to fix the situation. We will keep you updated as much as possible.

Again, our sincere apologies.



Like · Comment · Share · 👍 106 💬 29 ➦ 1

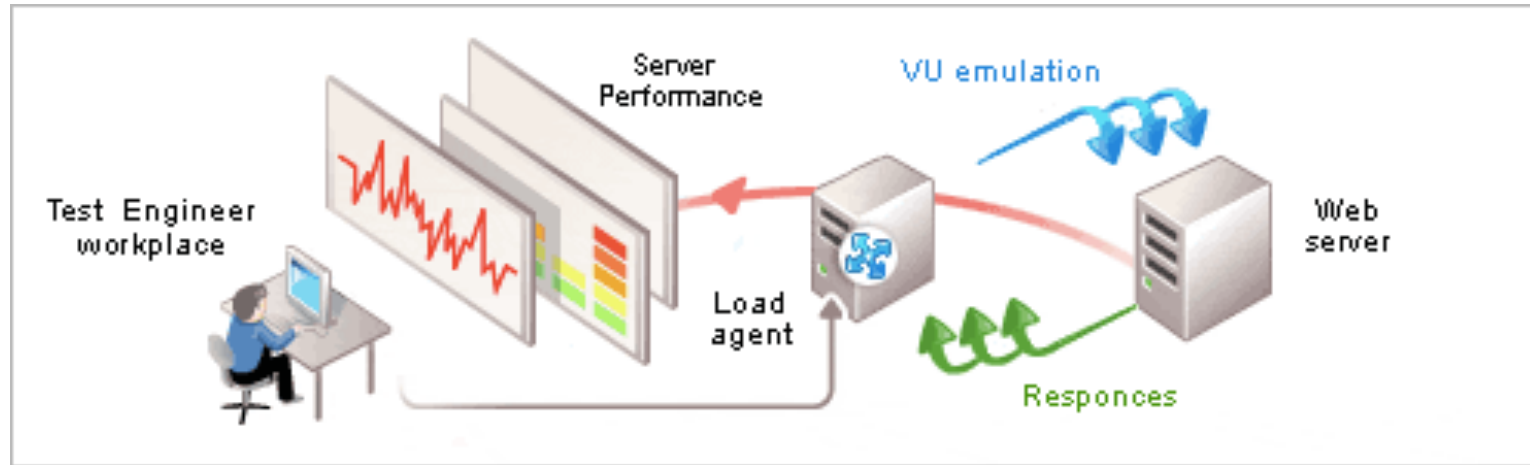
Jan I pray Lord for this matter to be resolved...Amen.
· February 11 at 11:56am

LOL Broncodirect... the worst thing about Cal Poly Pomona.
February 10 at 1:17pm

Gabriel Horowitz Lol one of the most stressful parts about going to school at CPP

Michelle Cassidy Every college has these problems during registration, I don't blame cpp, hope it's resolved soon though!

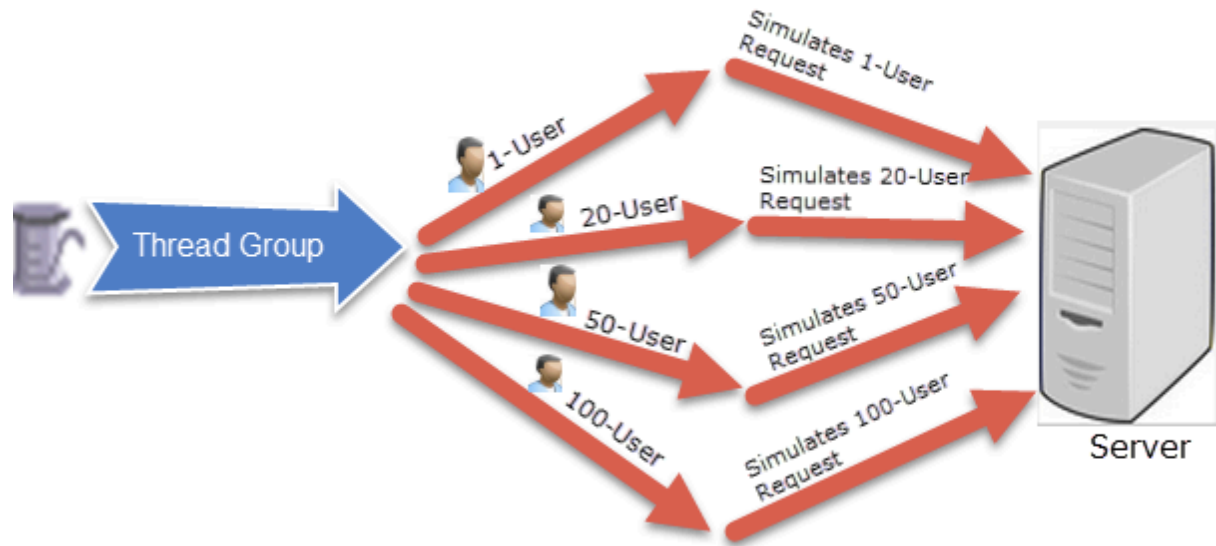
Scalability Verification – Load Test



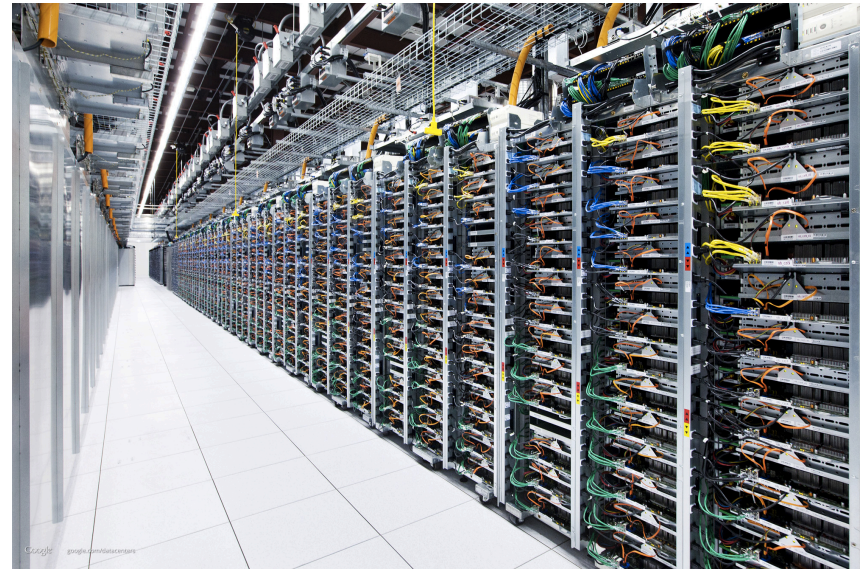
Scalability Verification – Load Test



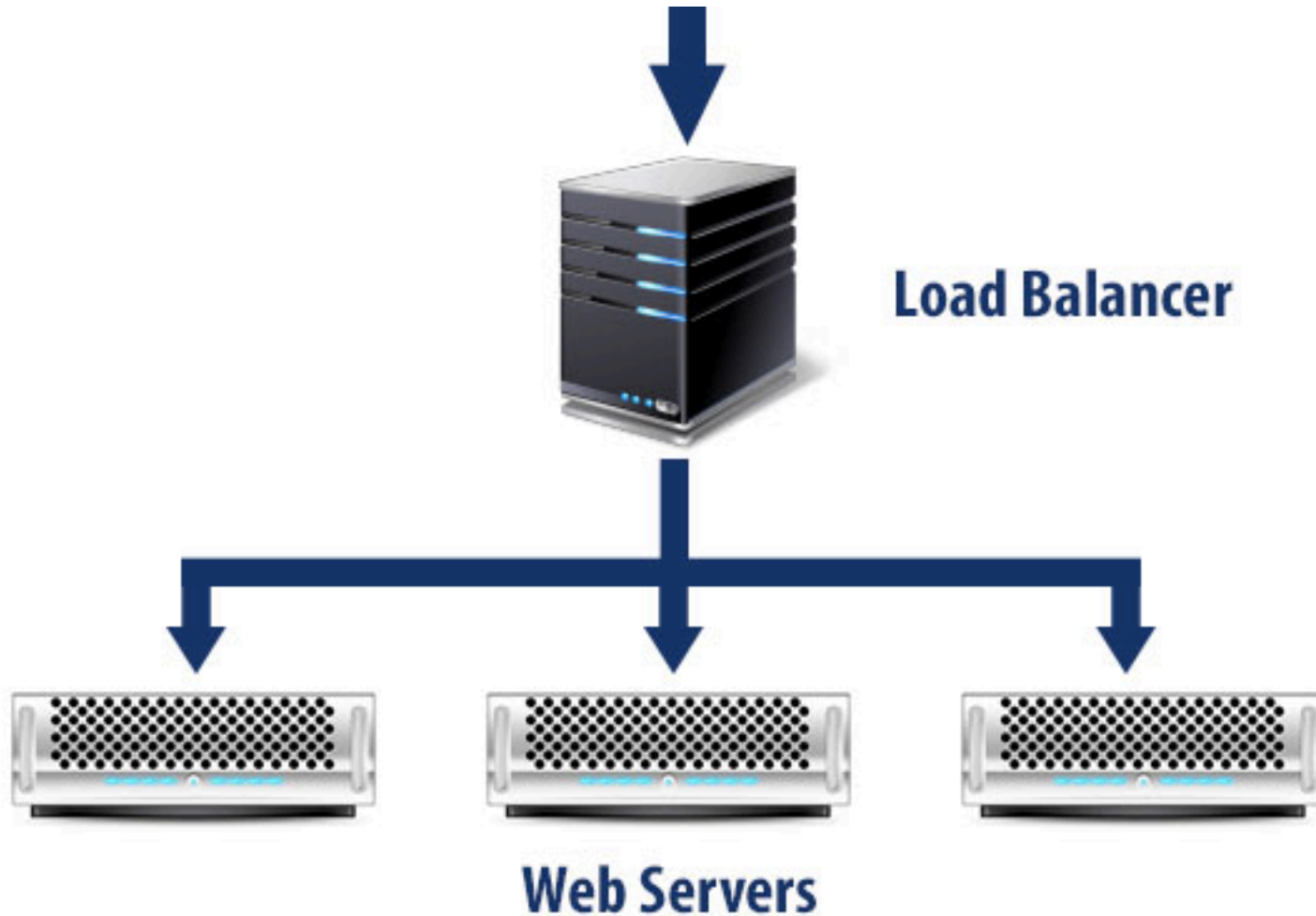
JMeter Demo



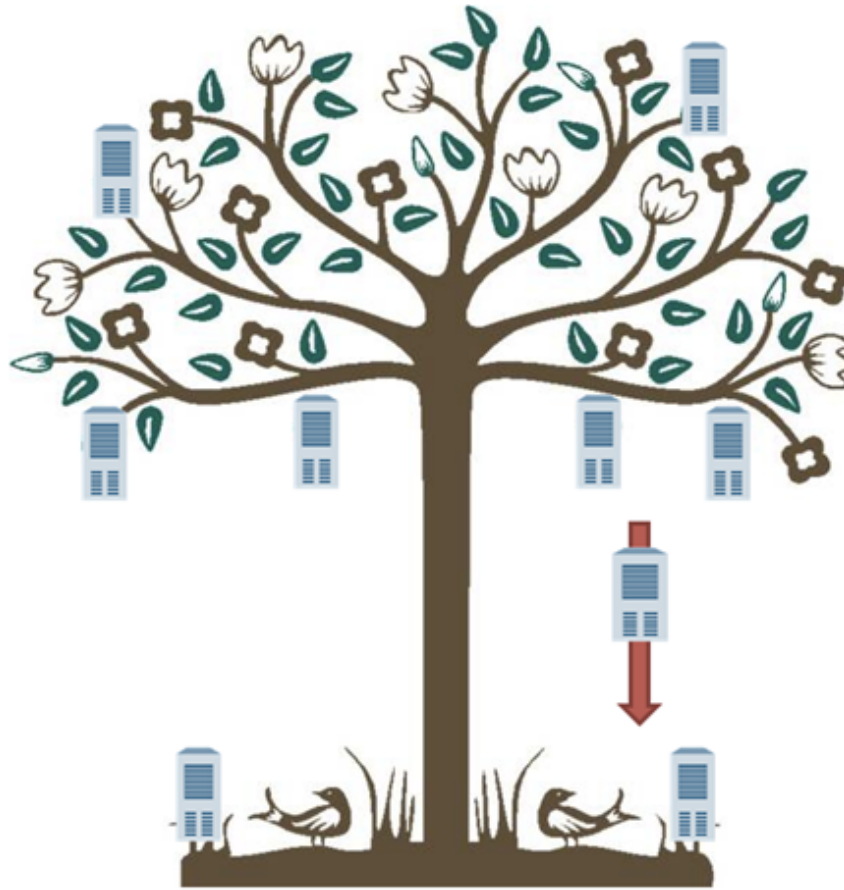
How to Improve Scalability?



How to Improve Scalability?



How to Make Scalability Easy?



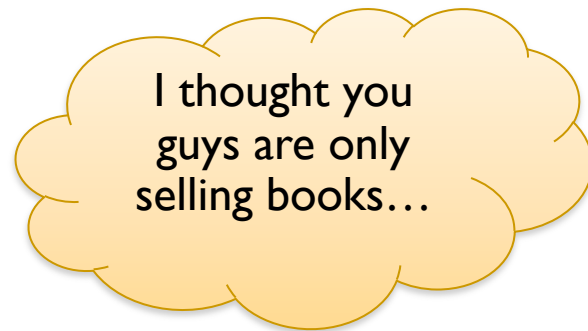
Cloud Computing

- ◆ Cloud computing shifts computing from local dedicated resources to distributed, virtual, elastic, multi-tenant resources
 - ◆ On-demand access to computing, storage, and software services
 - ◆ Based on a utility cost model



Cloud Computing & Amazon

- ◆ The popularization of the term can be traced to 2006 when Amazon.com introduced the Elastic Compute Cloud (EC2)



Case Study: Amazon Retail Website

The screenshot shows the Amazon website interface. At the top, there's the Amazon logo, navigation links like 'Today's Deals', 'Gift Cards', and 'Help', and a search bar. Below the search bar, there are links for 'Kindle Store', 'Buy a Kindle', 'Cloud Reader', 'Kindle eBooks', 'Kindle Singles', 'Newsstand', 'Popular Games', 'Accessories', 'Discussions', 'Manage Your Kindle', and 'Kindle Support'. A banner for 'Best of 2011 Best Books of 2011 So Far' is prominent. The main content area features a 'Kindle Store' section with the text 'Kindle The #1 Bestseller on Amazon' and an image of a Kindle device. Below this, there are three price options: Kindle for \$139, Kindle 3G with Special Offers for \$139, and Kindle 3G for \$189. To the right, there are sections for 'Kindle Daily Post', 'Great Deals on Kindle Accessories', and 'New on Kindle: Hidden Expedition: Amazon'. At the bottom, there's a 'Bestsellers' section for 'Kindle Store: All Kindle Content' and a 'Top 100 Paid' list.

amazon Today's Deals Gift Cards Help Save Up to 25% Off College Essentials

Shop by Department Kindle Store Hello, Derek Your Amazon 0 Wish List

Kindle Store Buy a Kindle Cloud Reader Kindle eBooks Kindle Singles Newsstand Popular Games Accessories Discussions Manage Your Kindle Kindle Support

Best of 2011 Best Books of 2011 So Far See all our editors' picks

Your Country or Region (What's this?) United States

Browse

Buy a Kindle

- Kindle (Wi-Fi, 6")
- Kindle 3G (Free 3G + Wi-Fi, 6")
- Kindle DX (Free 3G, 9.7", Graphite)
- Kindle DX (Free 3G, 9.7", White - 2nd Generation)

Kindle Reading Apps

- Kindle Cloud Reader
- Kindle for iPad
- Kindle for iPhone
- Kindle for PC
- Kindle for Mac
- Kindle for Android
- Kindle for BlackBerry
- Kindle for Windows Phone 7

Kindle for Windows Phone 7 See all

Give Kindle Gifts

- Give Kindle Books
- Give a Kindle Gift Card
- Redeem a Kindle Gift Card

Need Help?

- Getting Started
- Manage Your Kindle
- Kindle Support

Around the Store

- Kindle Singles

Kindle Store

Kindle (Wi-Fi, 6"), Kindle 3G (Free 3G + Wi-Fi, 6"), Accessories, and more than 950,000 Books, Kindle Singles, Newspapers, Magazines, Blogs, Audiobooks, and Games & Active Content

Kindle Daily Post

Our editors' blog Read new posts

Great Deals on Kindle Accessories

Save on select latest generation Kindle accessories from Cole Haan and kate spade new york. Shop now

New on Kindle: Hidden Expedition: Amazon

Join the Hidden Expedition Team in search of a professor lost in the Amazon rainforest in this hidden object game. Learn more

Bestsellers

Kindle Store: All Kindle Content

Updated hourly

Top 100 Paid Top 100 Free

1. 5 days in the top 100 Hidden in Plain View (Darryl Billups Mysteries) Blair S. Walker (Author)

Kindle \$114 with Special Offers Order now

Kindle \$139 Order now

Kindle 3G with Special Offers \$139 Order now

Kindle 3G \$189 Order now

More Items to Consider

You viewed Customers who viewed this also viewed

THE INNOVATOR'S DNA

BASIC ECONOMICS

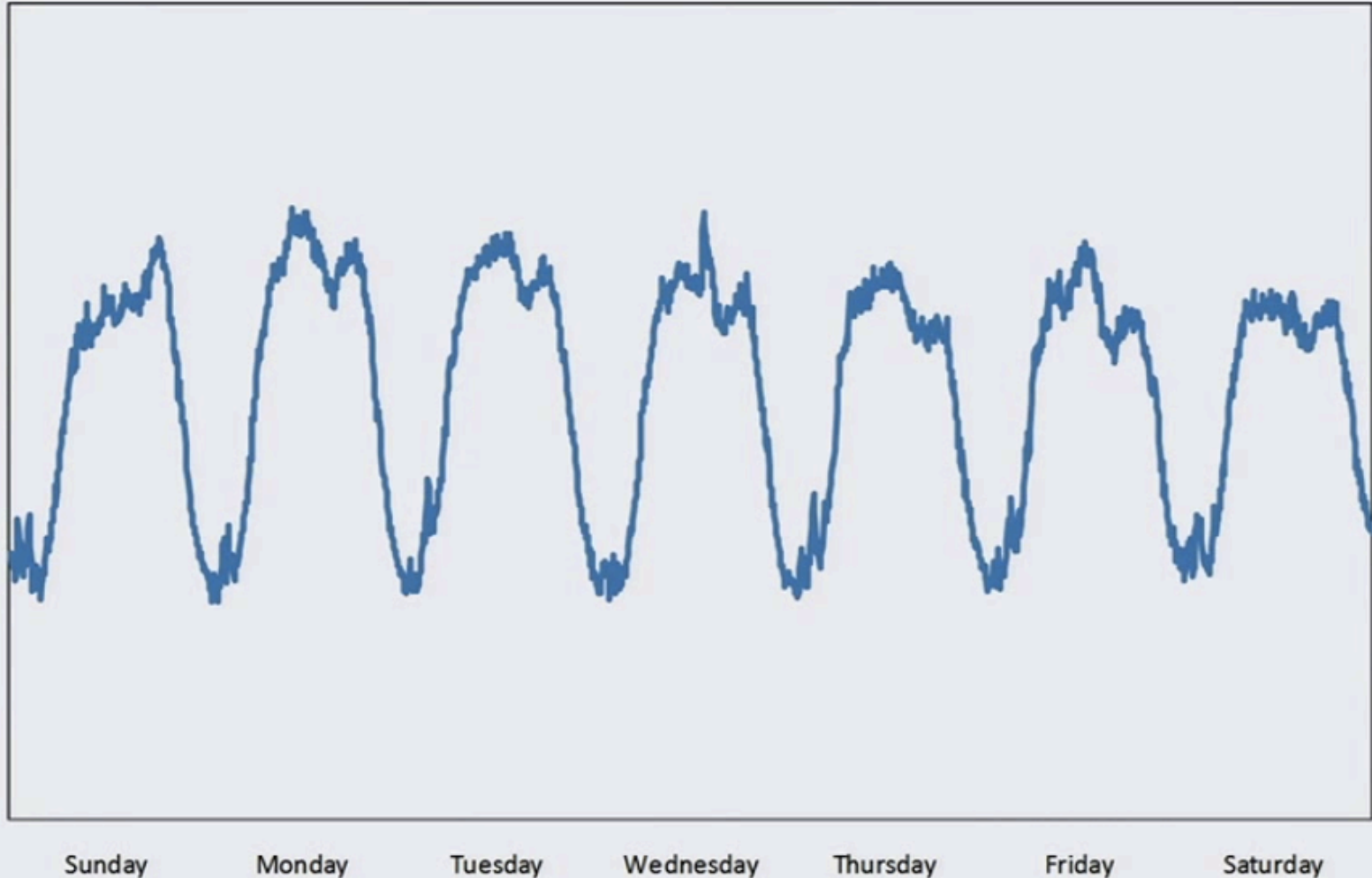
The Innovator's Dilemma

venture deals



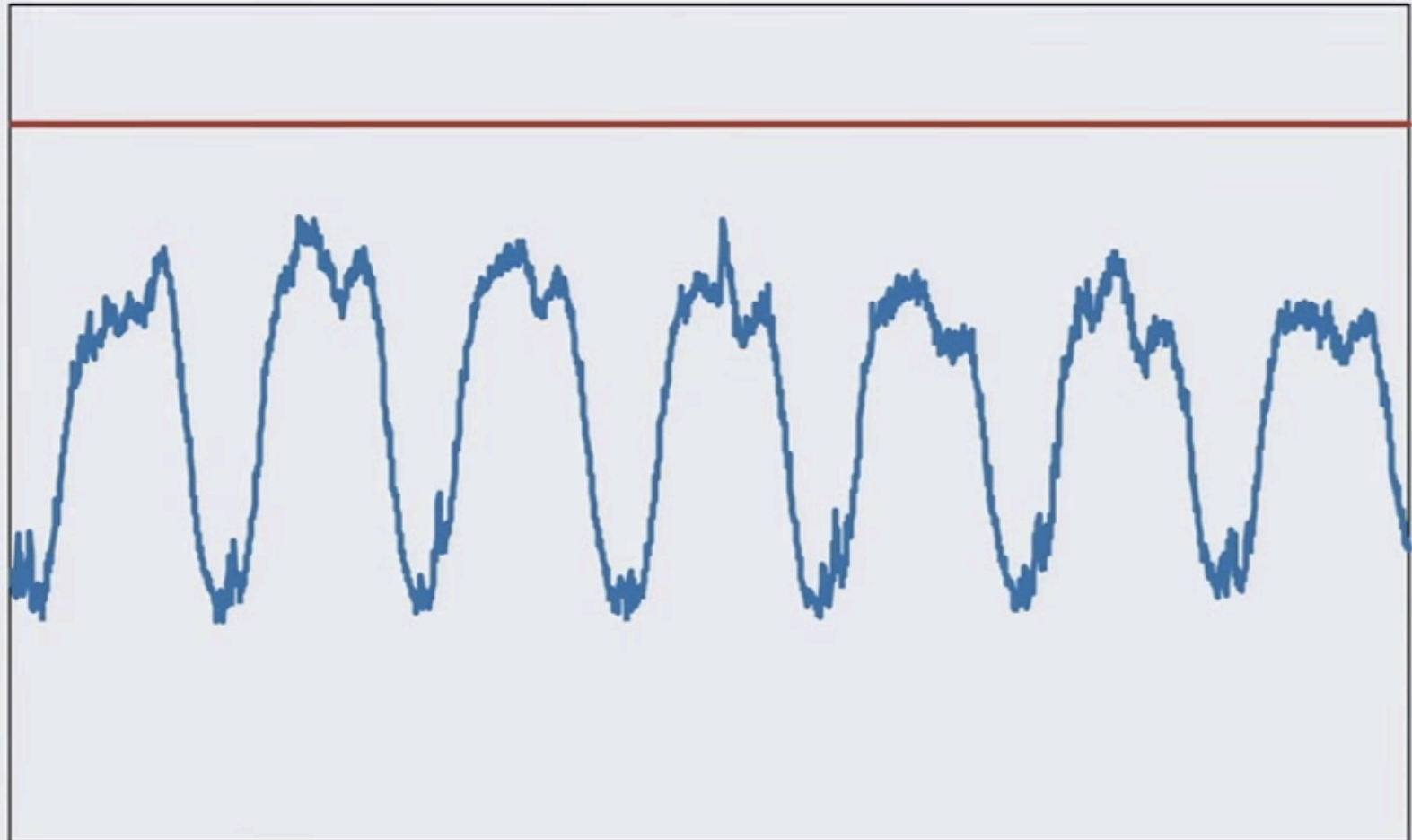
Case Study: Amazon Retail Website

Typical Weekly Traffic to amazon.com



Case Study: Amazon Retail Website

Typical Weekly Traffic to amazon.com



Sunday

Monday

Tuesday

Wednesday

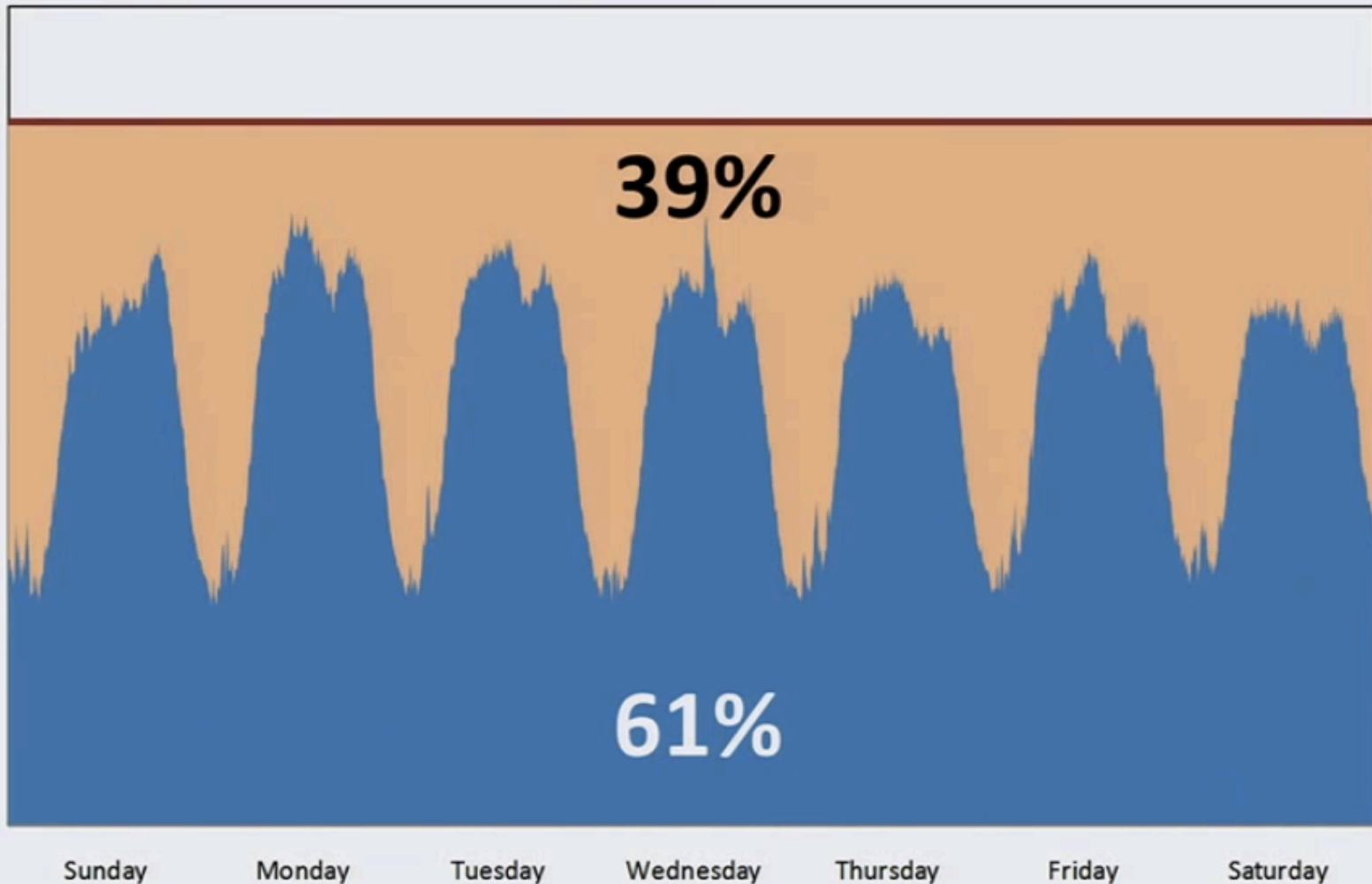
Thursday

Friday

Saturday

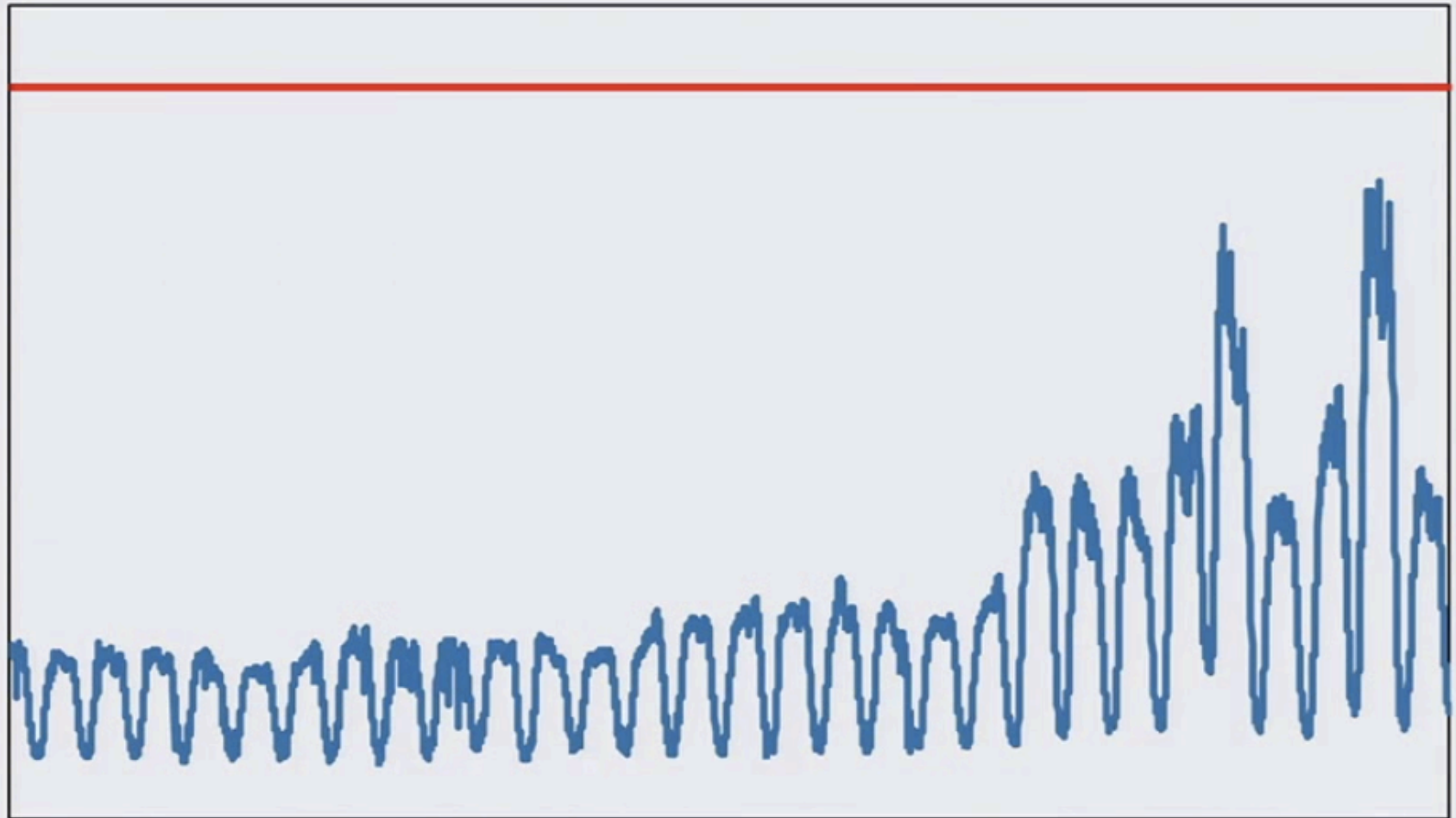
Case Study: Amazon Retail Website

Typical Weekly Traffic to amazon.com



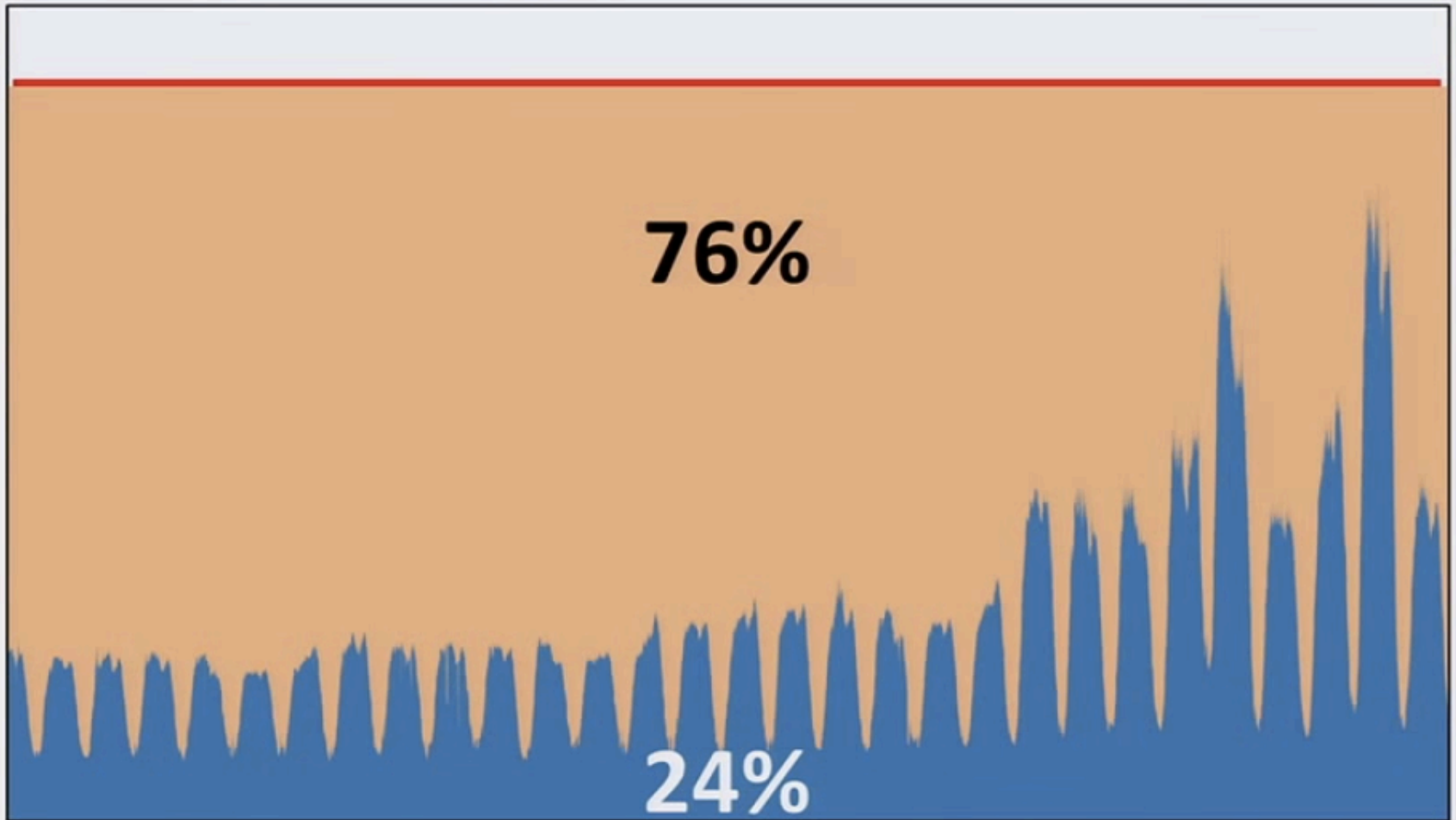
Case Study: Amazon Retail Website

November Traffic for amazon.com



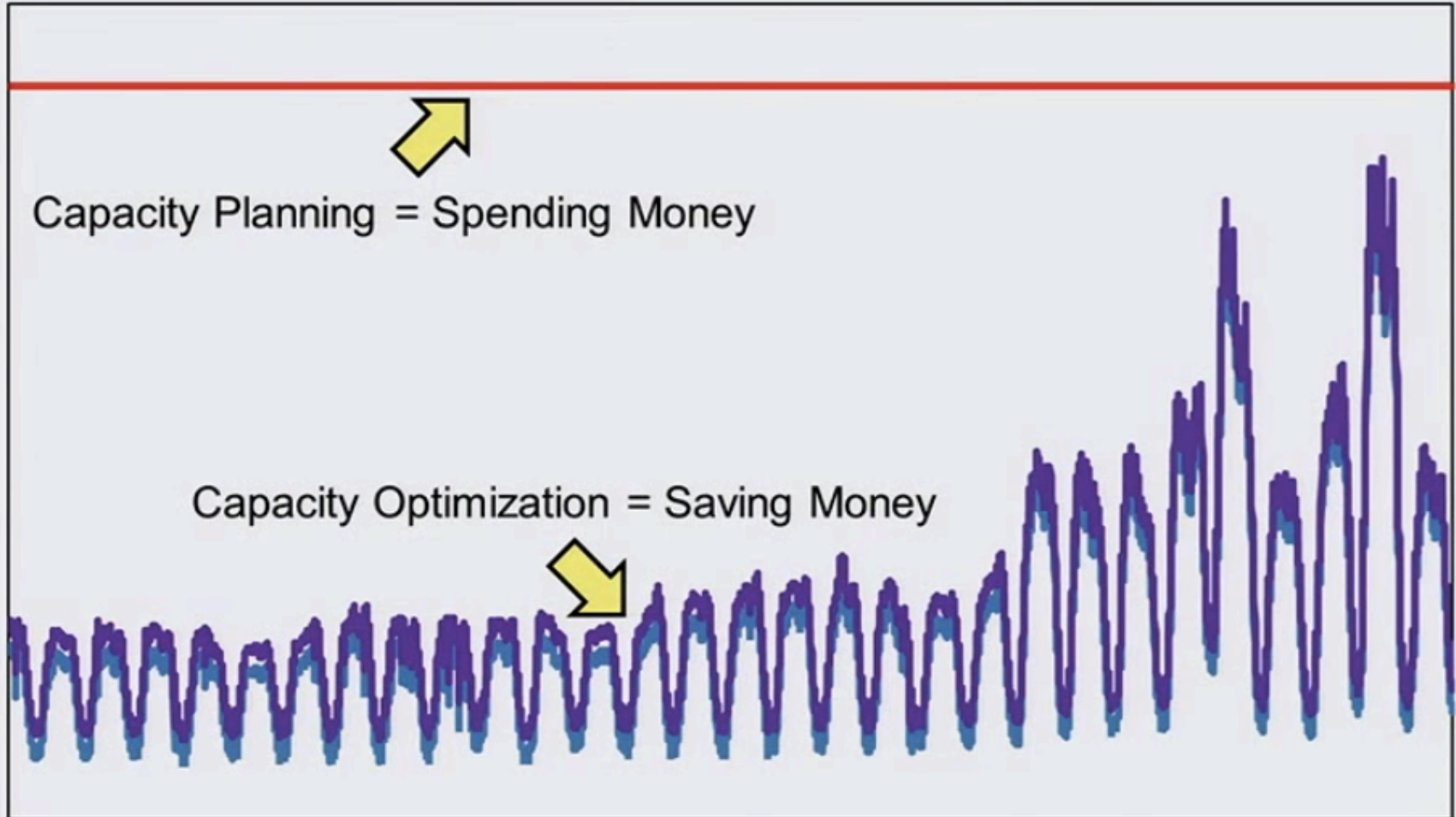
Case Study: Amazon Retail Website

November Traffic for amazon.com



Motivation of Cloud Computing (I)

November Traffic for amazon.com



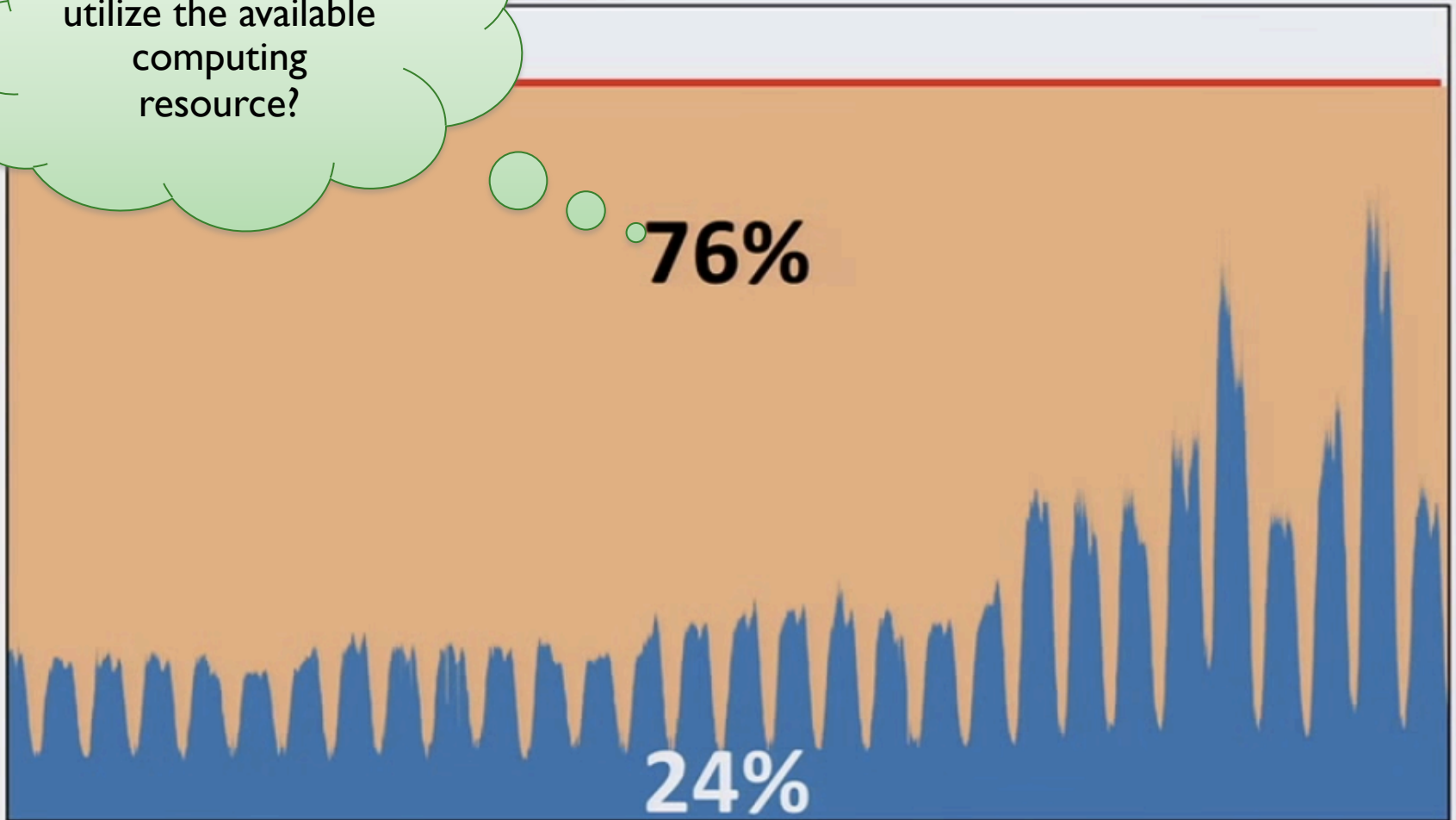
Motivation of Cloud Computing (2)

Number Traffic for amazon.com

How can we better utilize the available computing resource?

76%

24%



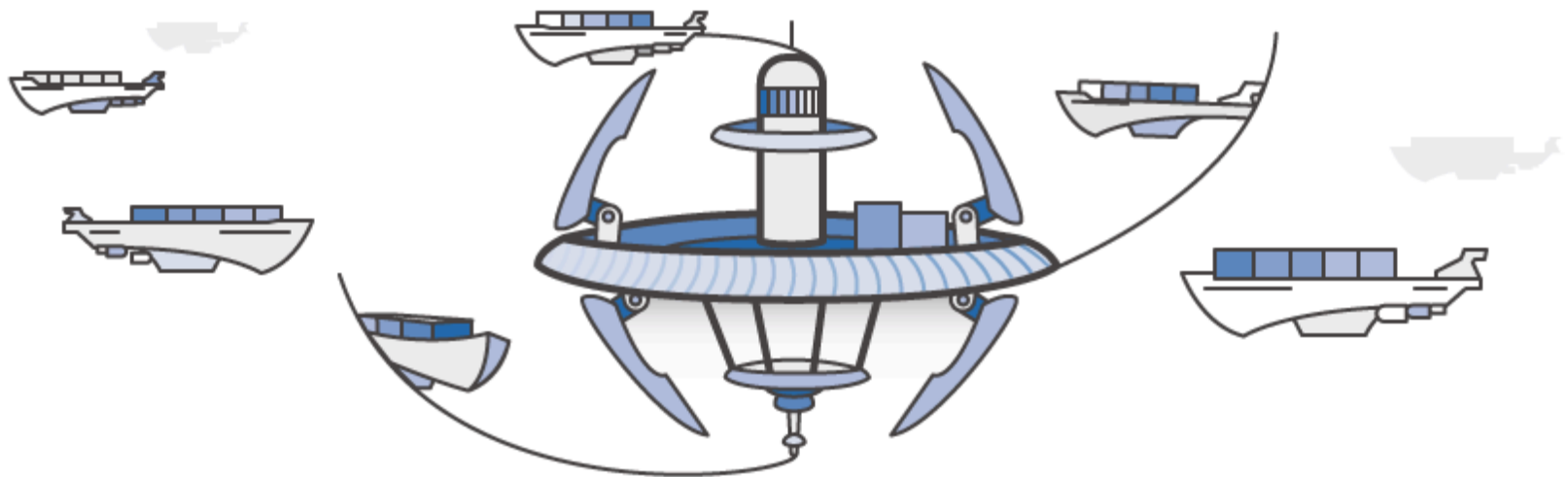
Cloud Computing

- ◆ Cloud computing shifts computing from local dedicated resources to distributed, virtual, elastic, multi-tenant resources
 - ◆ **On-demand access** to computing, storage, and software services
 - ◆ Based on a **utility cost model**

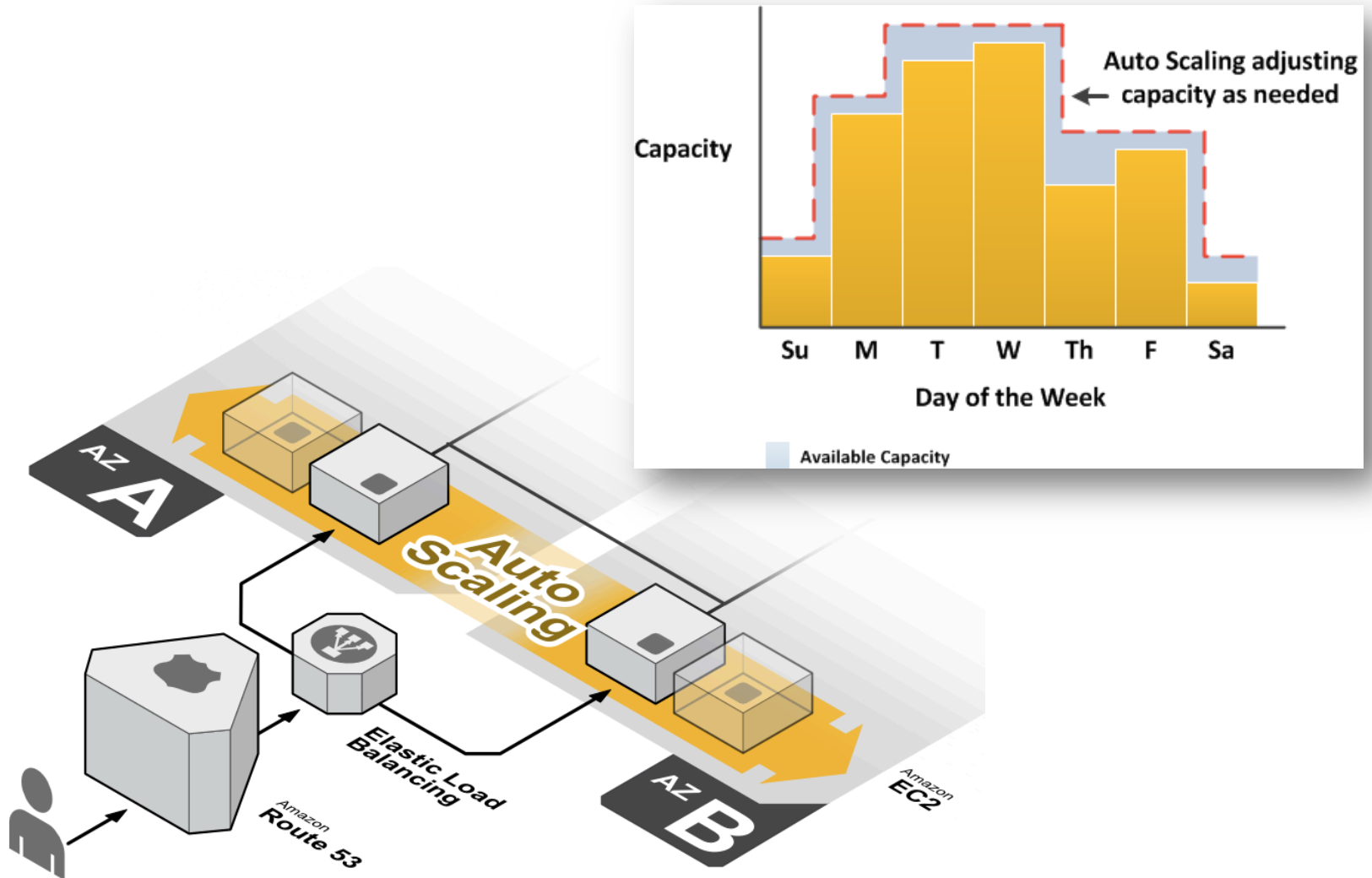


Amazon EC2 On-Demand Pricing

United States		Europe	
Standard On-Demand Instances		Linux/UNIX Usage	Windows Usage
Small (Default)		\$0.10 per hour	\$0.125 per hour
Large		\$0.40 per hour	\$0.50 per hour
Extra Large		\$0.80 per hour	\$1.00 per hour
High CPU On-Demand Instances		Linux/UNIX Usage	Windows Usage
Medium		\$0.20 per hour	\$0.30 per hour
Extra Large		\$0.80 per hour	\$1.20 per hour



Auto Scaling with Cloud Computing



Auto-Scaling Demo

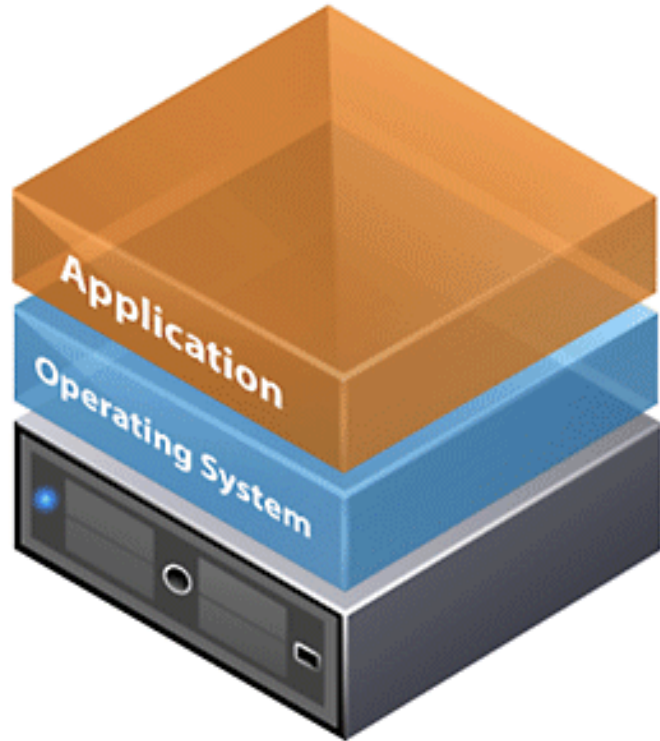


Amazon.com is Fully Served by EC2

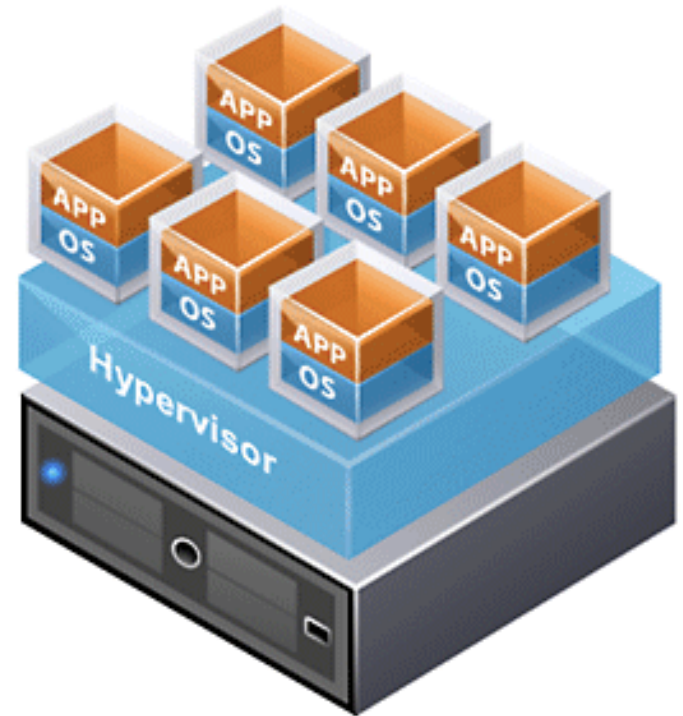
- ◆ Reduced spending on server capacity
- ◆ Fleet scales dynamically in increments as small as a single host
- ◆ Traffic spikes can be handled with ease
- ◆ Cultural change



Virtualization



Traditional Architecture



Virtual Architecture

Flexible Options

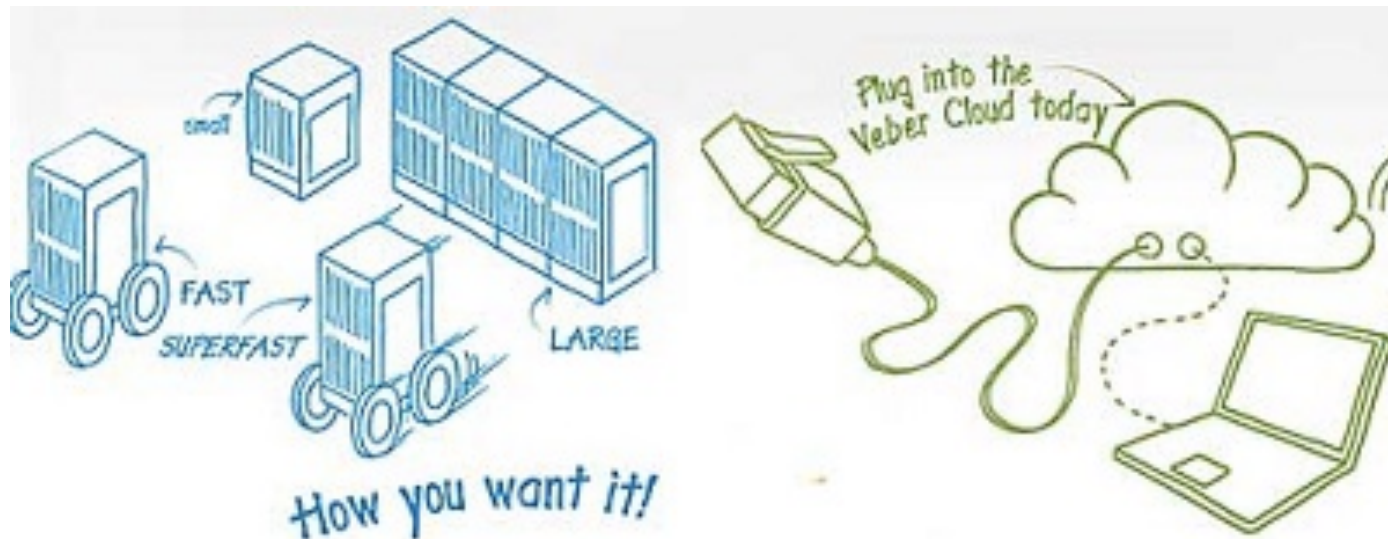
Google vs. AWS On-Demand Pricing

Google Instance Type	CPU Cores	RAM	AWS Instance Type	CPU Cores	RAM	Google New On-Demand (per hour)	AWS On-Demand (per hour)	New Google Price vs. AWS
n1-standard-1	1	3.75	m3.medium	1	3.75	\$ 0.070	\$ 0.113	-38.05%
n1-standard-2	2	7.5	m3.large	2	7.5	\$ 0.140	\$ 0.225	-37.78%
n1-standard-4	4	15	m3.xlarge	4	15	\$ 0.280	\$ 0.450	-37.78%
n1-standard-8	8	30	m3.2xlarge	8	30	\$ 0.560	\$ 0.900	-37.78%
n1-highmem-2	2	13	m2.xlarge	2	17.1	\$ 0.164	\$ 0.410	-60.00%
n1-highmem-4	4	26	m2.2xlarge	4	34.2	\$ 0.328	\$ 0.820	-60.00%
n1-highmem-8	8	52	m2.4xlarge	8	68.4	\$ 0.656	\$ 1.640	-60.00%
n1-highcpu-2	2	1.8	c3.large	2	3.75	\$ 0.088	\$ 0.150	-41.33%
n1-highcpu-4	4	3.6	c3.xlarge	4	7.5	\$ 0.176	\$ 0.300	-41.33%
n1-highcpu-8	8	7.2	c3.2xlarge	8	15	\$ 0.352	\$ 0.600	-41.33%
n1-highcpu-16	16	14.4	c3.4xlarge	16	30	\$ 0.704	\$ 1.200	-41.33%

March 25, 2014

Source: RightScale

Rapid Resource Allocation

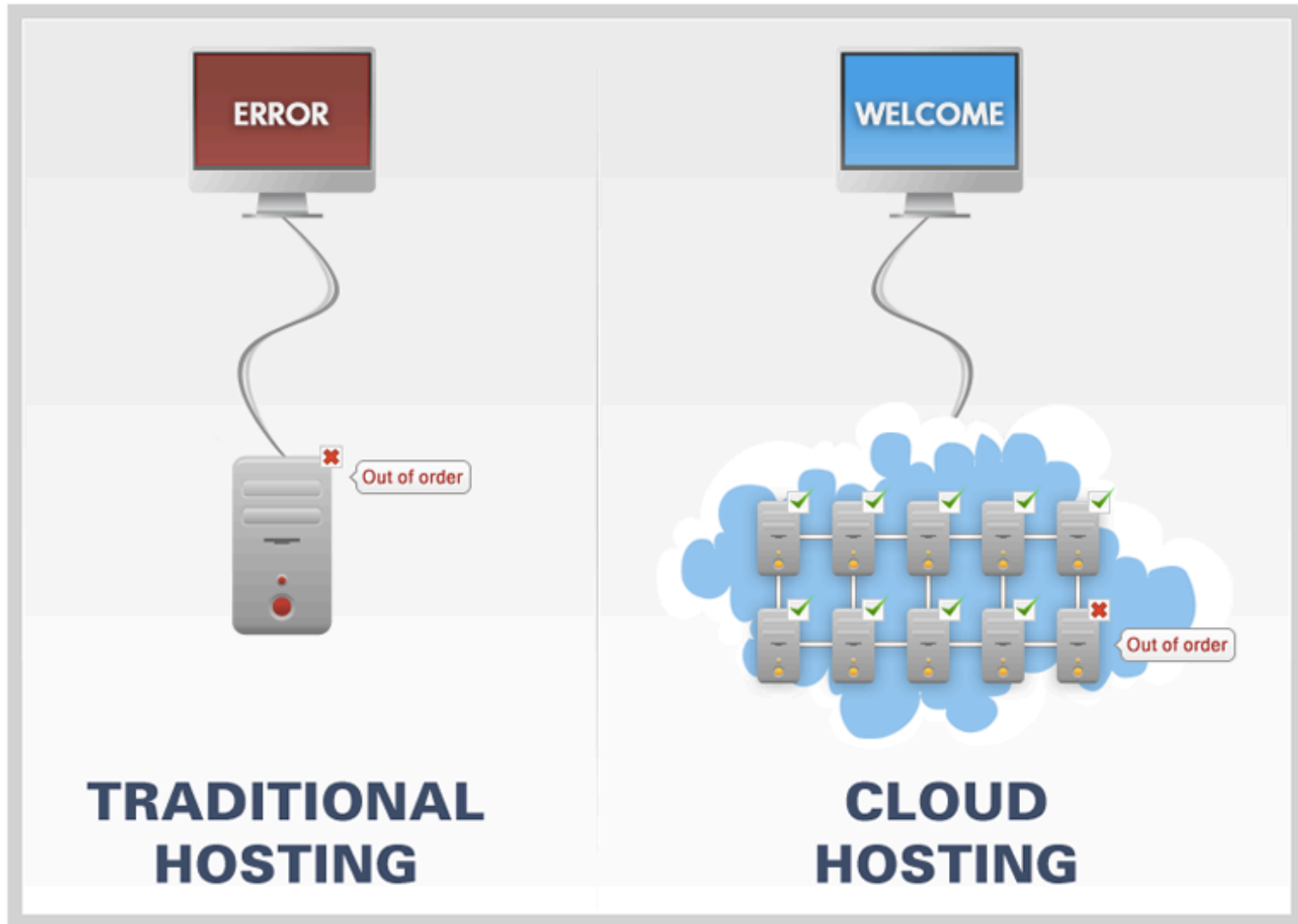


Dedicated

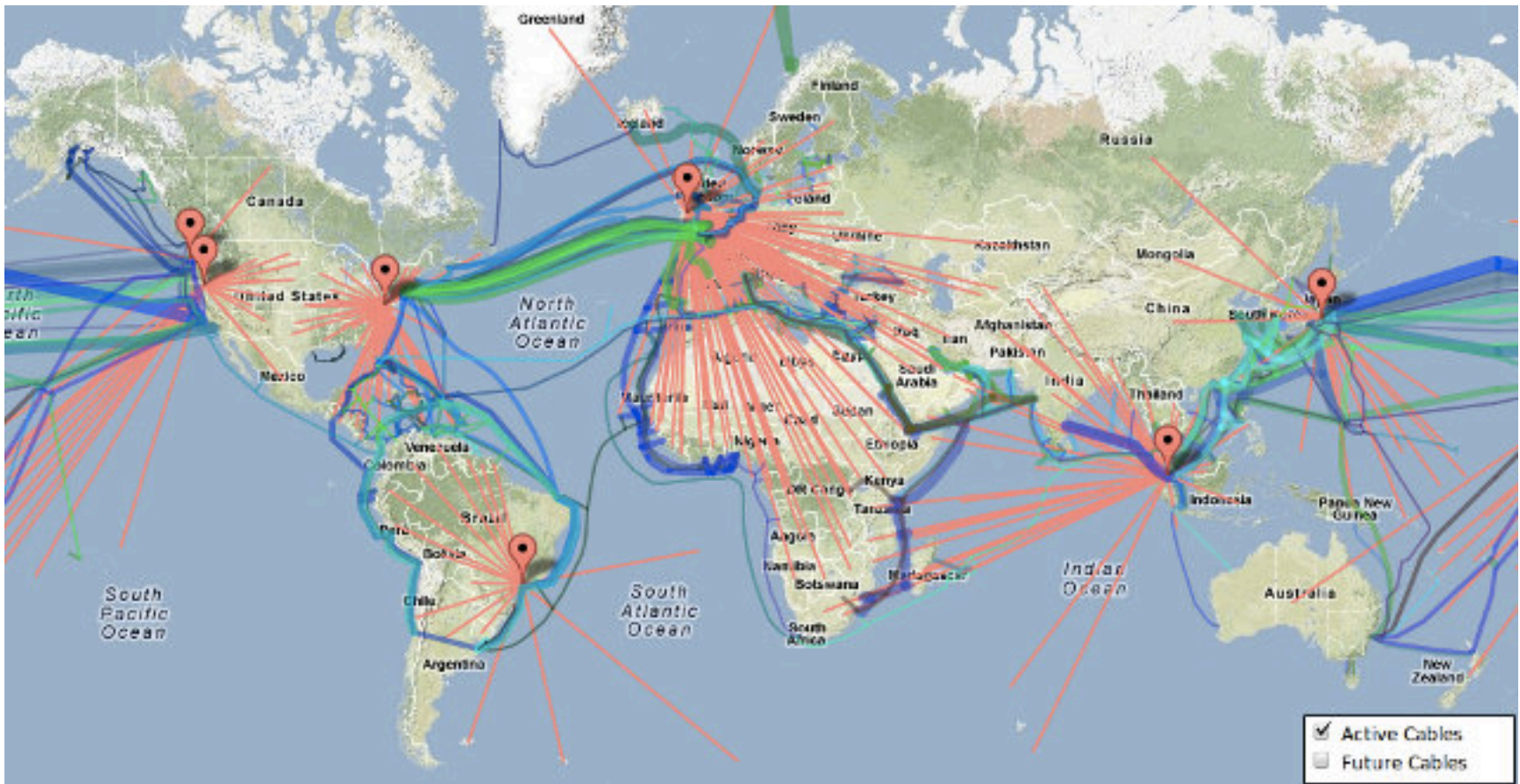
Vs

Cloud

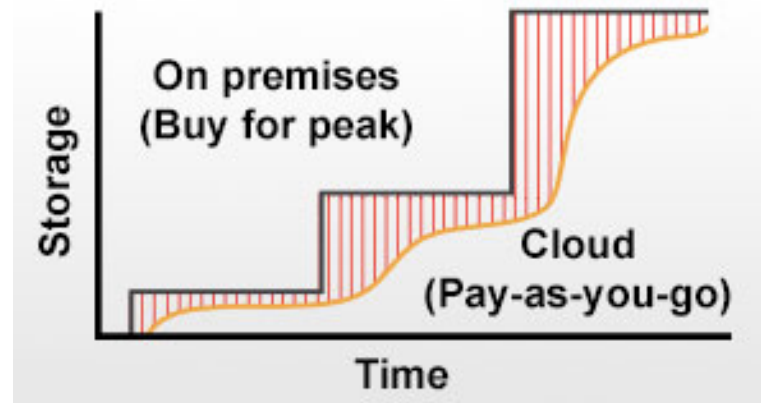
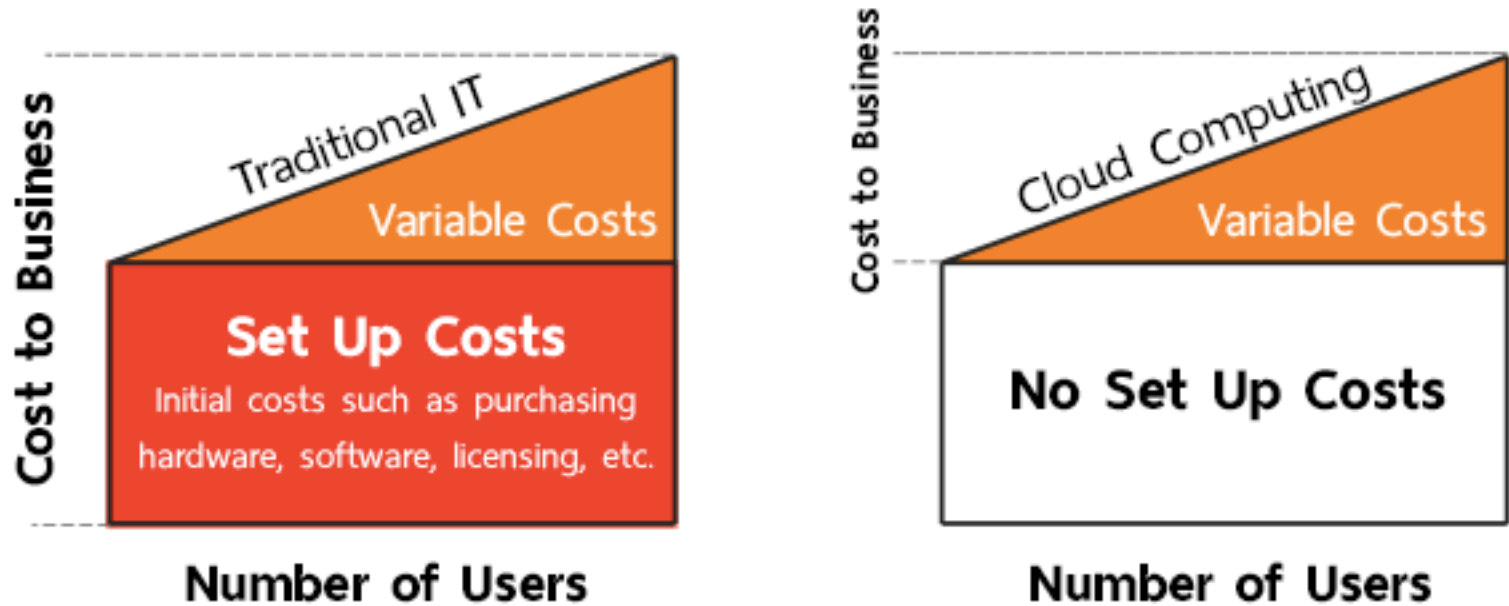
High Availability



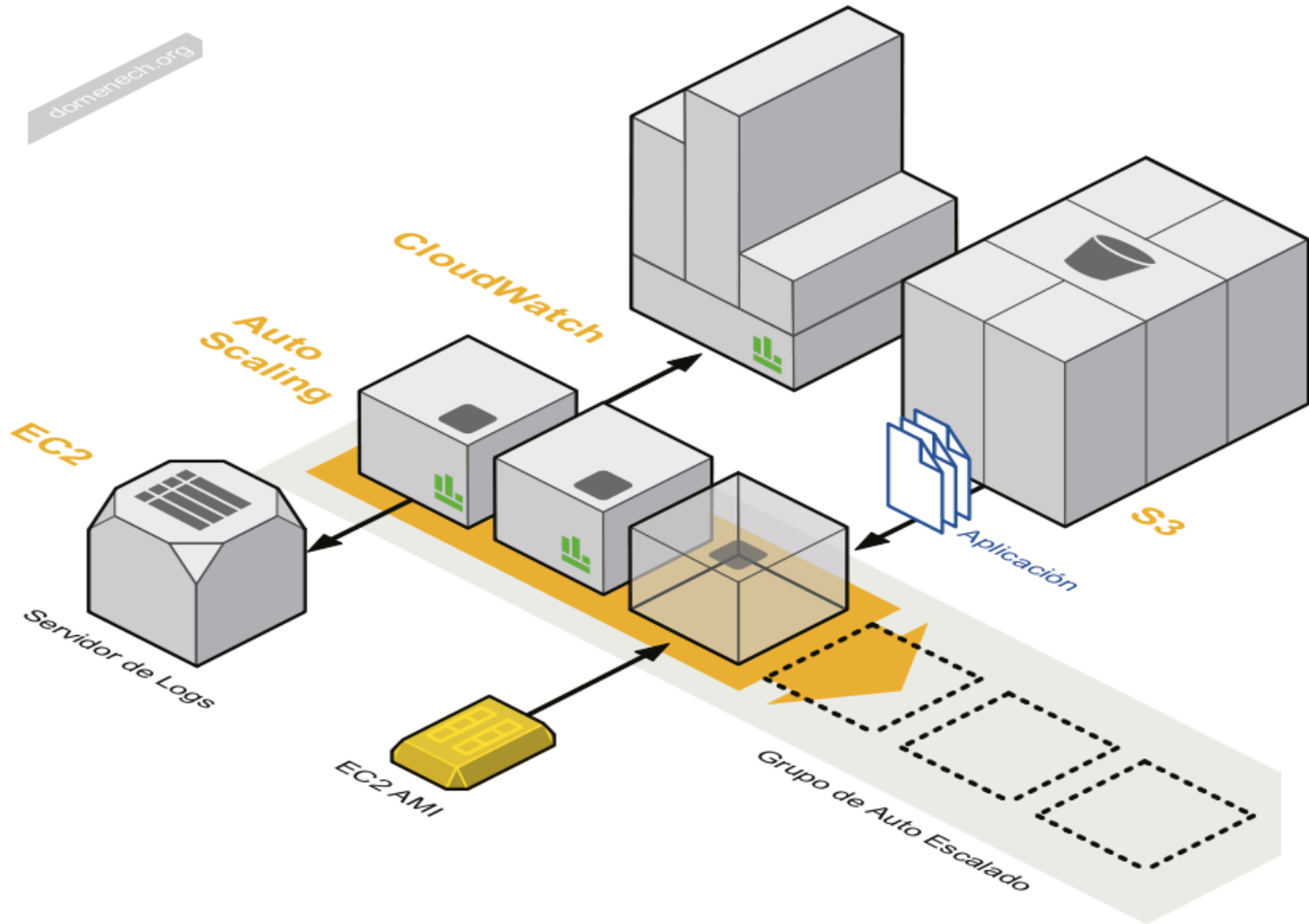
High Availability: AWS Data Centers



Cost-Effective



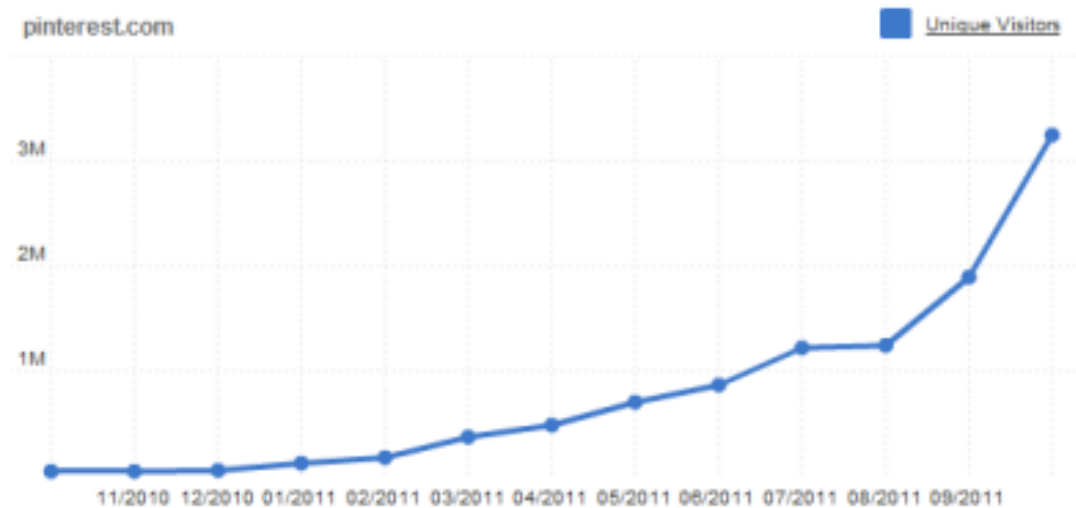
Scalability: Auto-Scaling



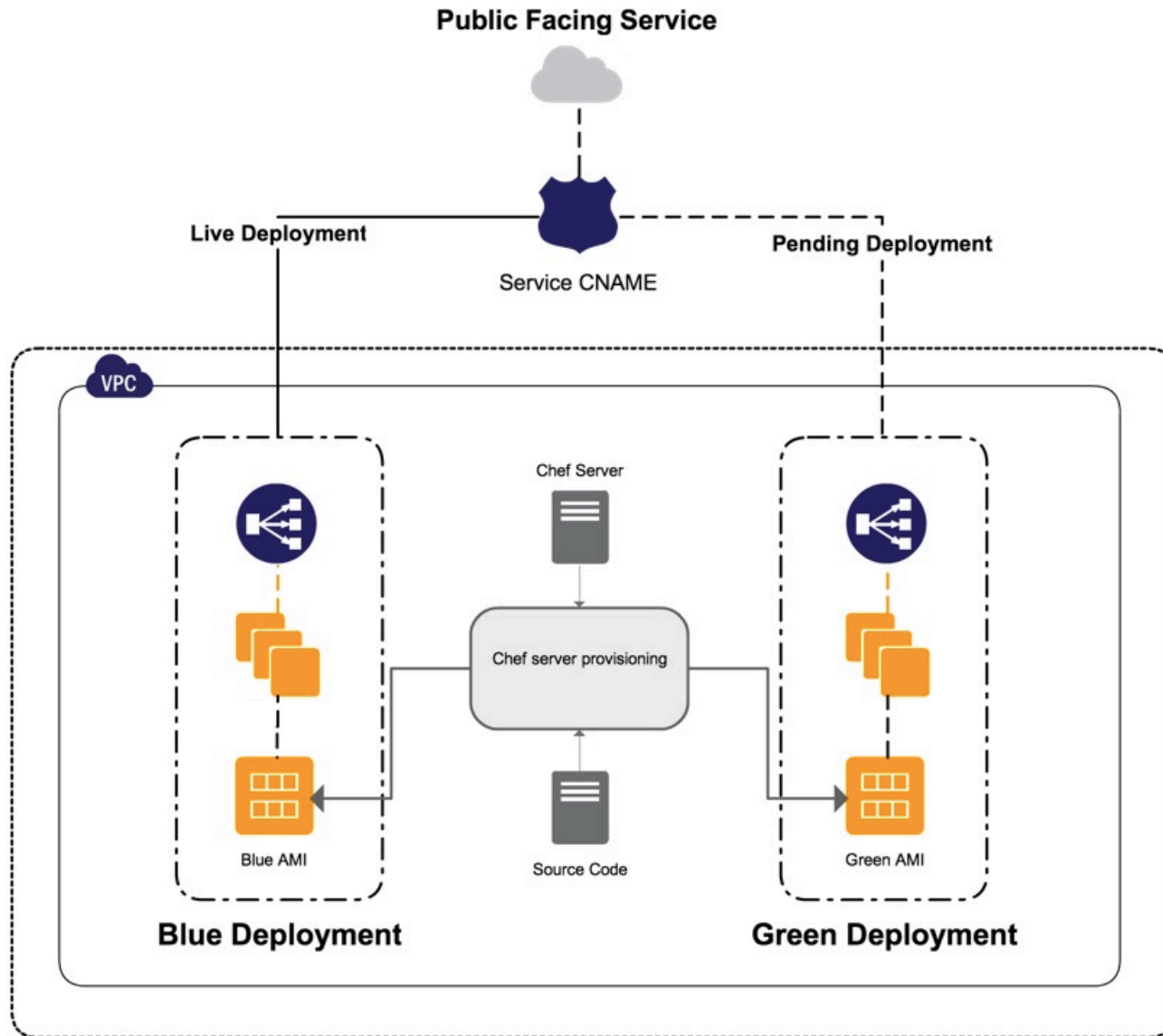
Focus on the Applications



Startups Made Easy

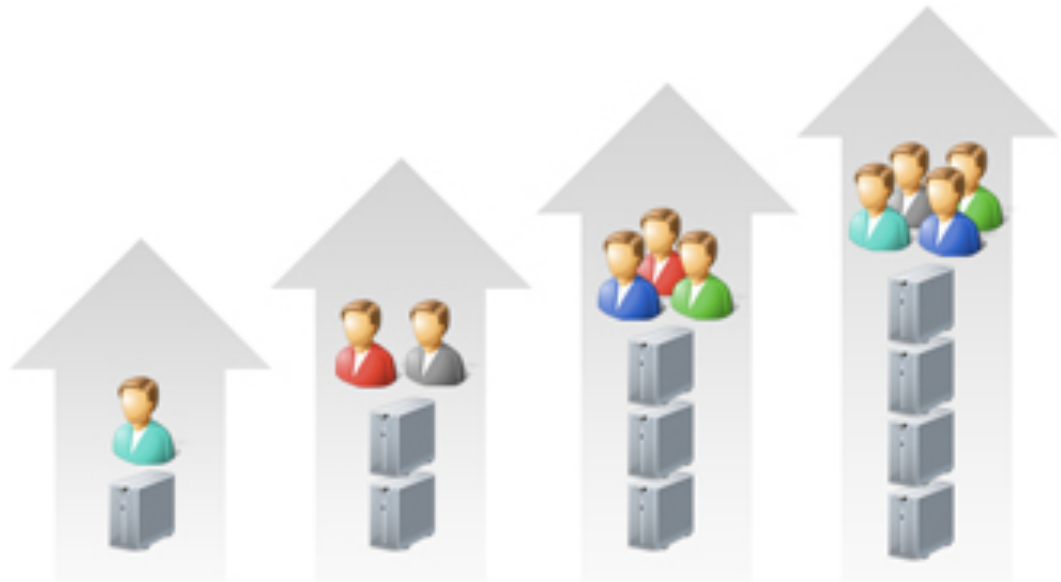


Deployment with Ease

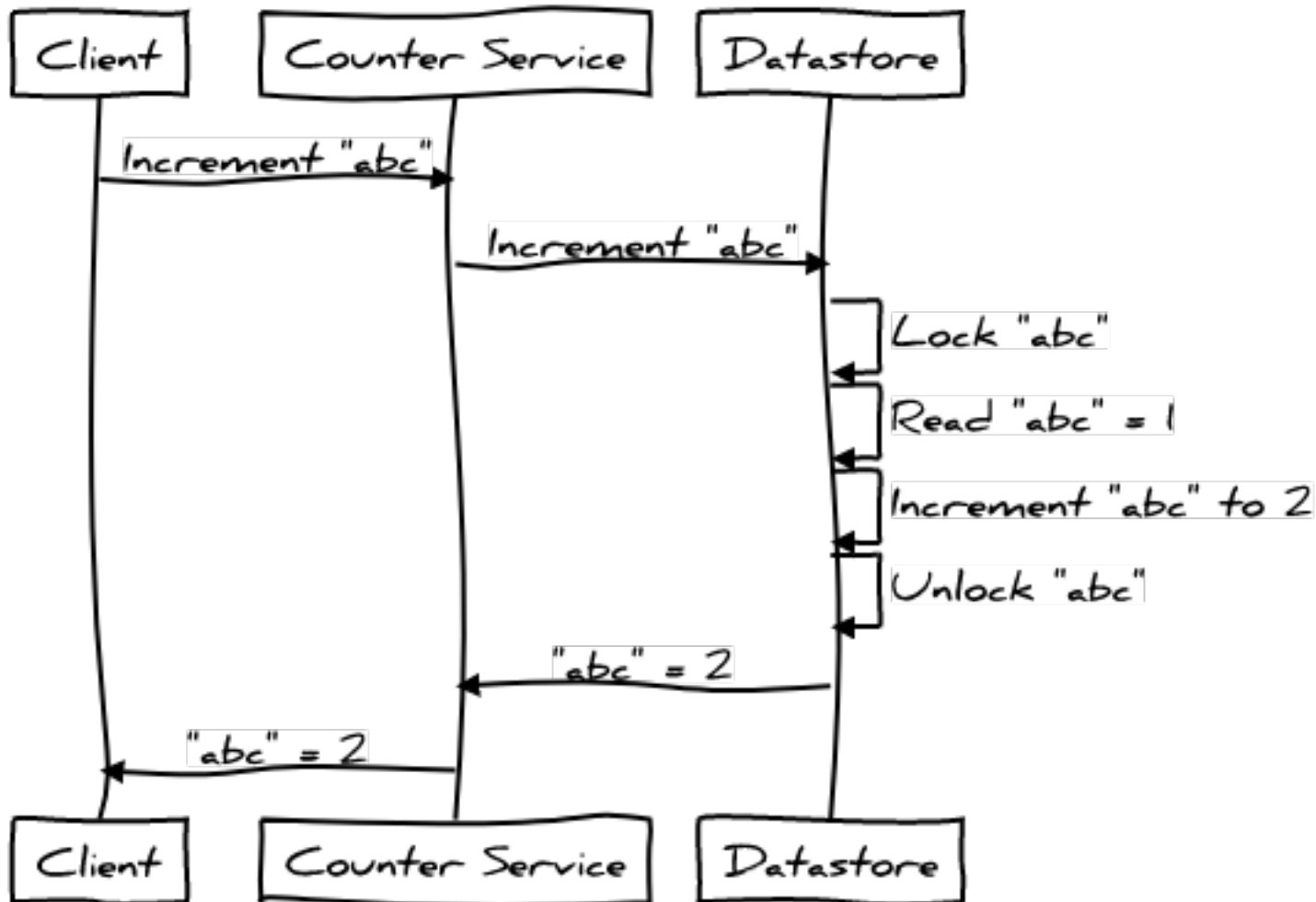


Always Keep Scalability in Mind

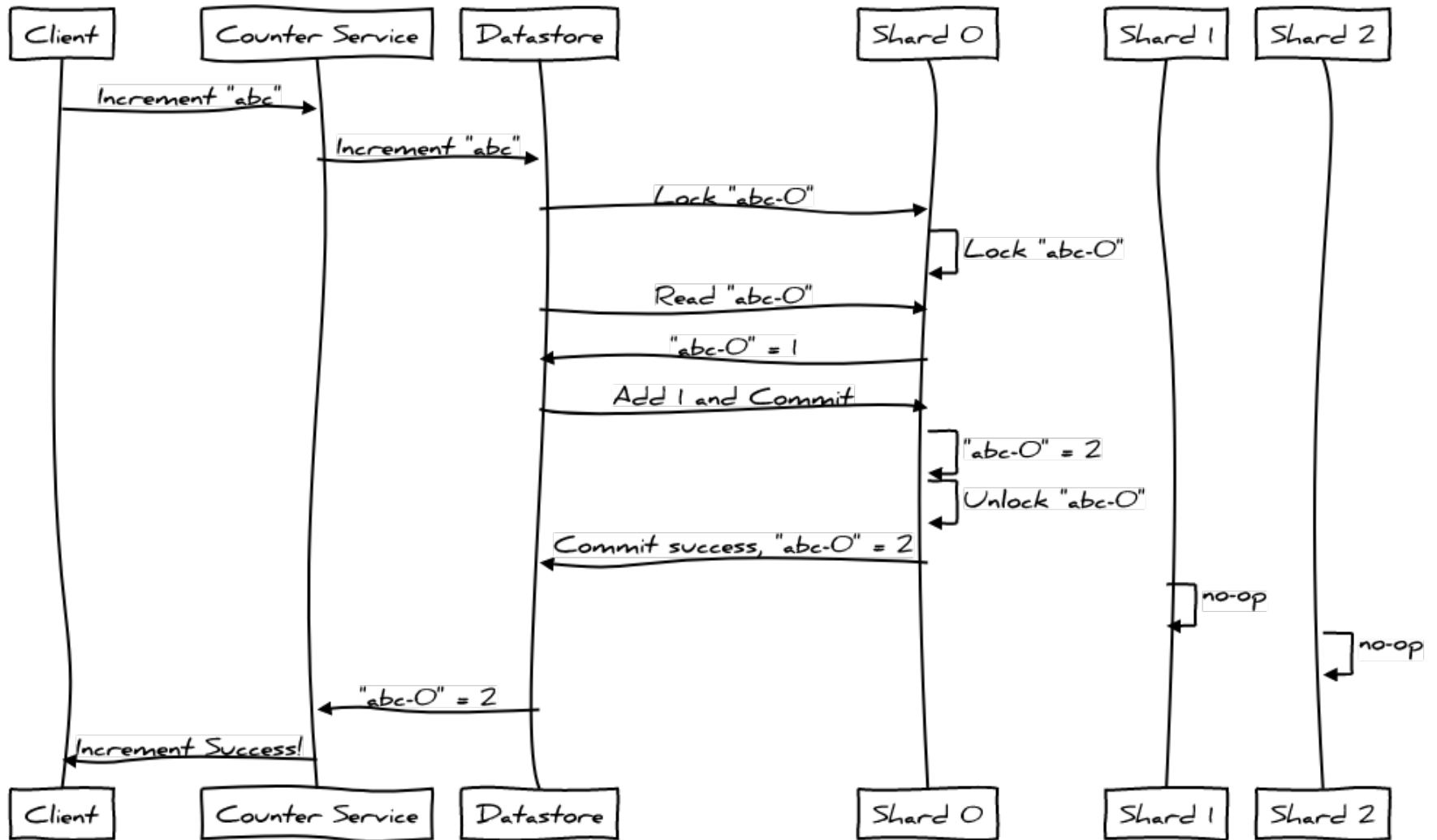
- ◆ Minimizing work
- ◆ Paging through large datasets
- ◆ Avoiding datastore contention
- ◆ Sharding counters
- ◆ Effective memcache



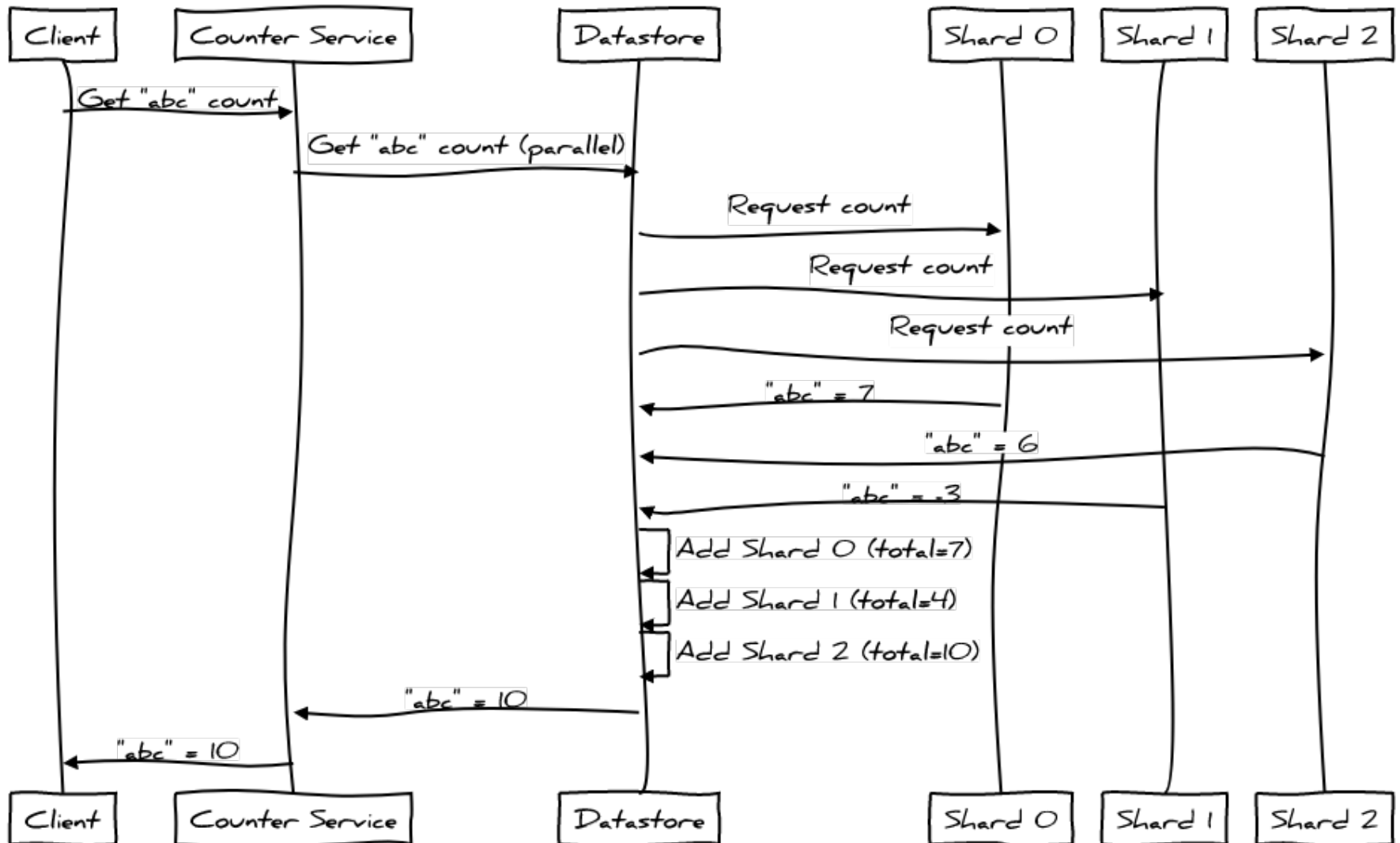
Sharding Counters



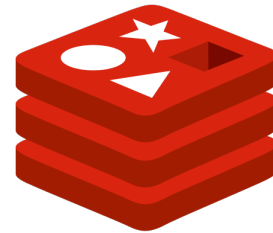
Sharding Counters



Sharding Counters



Effective Memcache



redis