# *When Hybrid Cloud Meets Flash Crowd: Towards Cost-Effective Service Provisioning*

Yipei Niu[1], Bin Luo[1], Fangming Liu[1],
Jiangchuan Liu[2], Bo Li[3]

Email: fmliu@hust.edu.cn

[1]*Huazhong University of Science & Technology*
[2]*Simon Fraser University*
[3]*The Hong Kong University of Science & Technology*

# E-commerce miracle in promotion seasons

**During promotion seasons**

- E-commerce websites offer attractive discounts.



**Alibaba claims:**

- The GMV (gross merchandise volume) on Nov. 11, 2013 is $5.8 billion
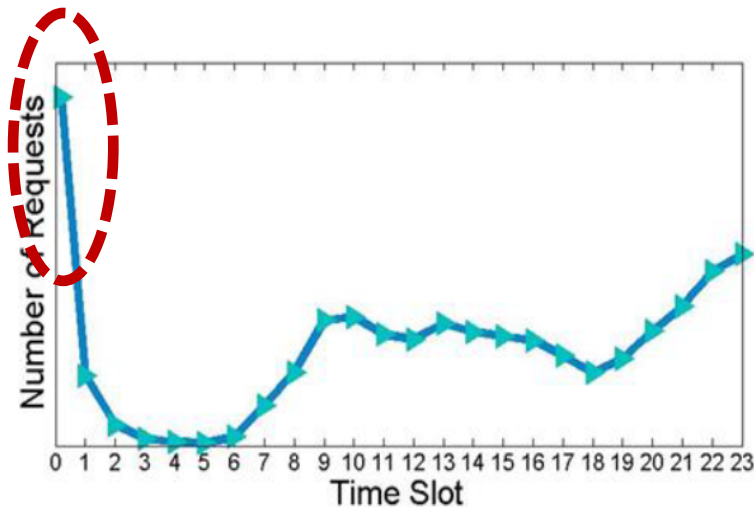- The GMV on Nov. 11 reaches $9.3 billion in 2014

**In the U.S.**

- Online sales exceeds $1 billion on Thanksgiving Day
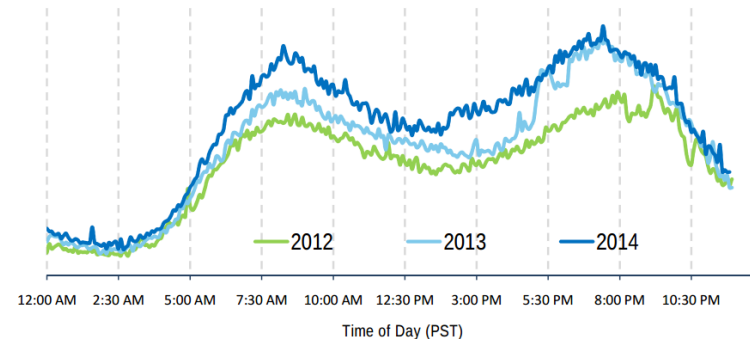- Online sales passes $1.5 billion on Black Friday

# Are e-commerce websites excited?



Thanksgiving Day 2014
U.S. Retail, 24-hr Real-time Sales Chart
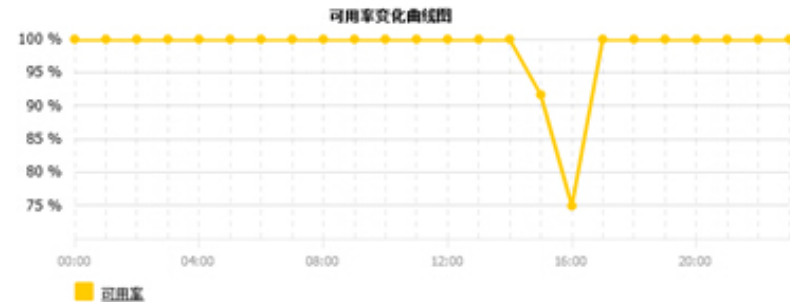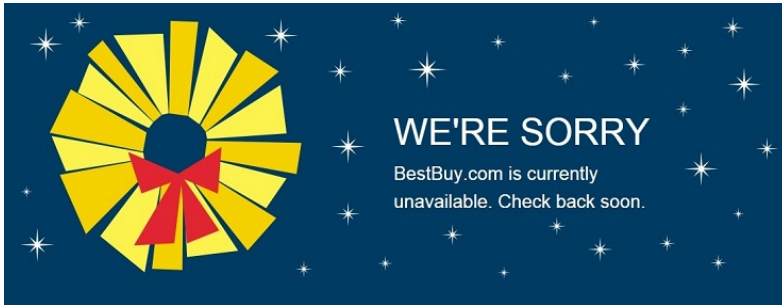
IBM Digital Analytics Benchmark

Online sales peaked in the morning and again in the evening on Thanksgiving. Online shopping began to pick up earlier (around 1:45 p.m. PST) than in 2013 leading up to the evening peak time of 7:30 p.m. PST. Online sales grew 14.3% over 2013.

**During Double Eleven Shopping Festival**

- 13.7 million buyers simultaneously visited Tmall in 2013
- 340,000 orders were placed during the first minute in 2013
- 15,000 online transactions per second at the peak in 2013
- 47,500 payment transactions per second at the peak in 2014

**E-commerce websites have to face bursty, immense, and unpredictable flash crowds brought by promotion seasons**

Source: http://www.alizila.com/1111-shopping-festival-fast-facts

# Why e-commerce websites headache?



- **BestBuy.com Crashes on Black Friday in 2014**
  - On Friday morning, BestBuy.com went offline.
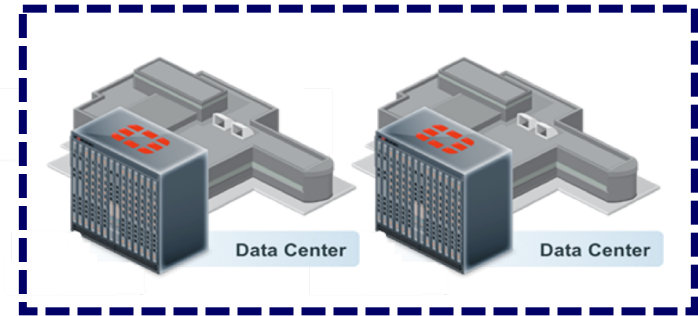  - Around 11:30 a.m., the site was back online, after approximately 1.5 hours offline

- **Vancl is unavailable on shopping festival**
  - On Nov.11, 2014, most e-commerce websites have the availability of 100%.
  - However, only Vancl outraged three times, 20 minutes unavailability in total.

Source: http://www.pcmag.com/article2/0,2817,2472895,00.asp

# Is private cloud OK?

- **Private cloud**
  - Dedicated datacenter or server cluster
  - Virtual resources provided by cloud providers



**Private cloud**

- **Private cloud solution**
  - **Advantages**
    - Enhanced security
    - Ultimate control
  - **Disadvantages**
    - Limited capacity
    - Low scalability
    - Complex to operate

- **Requirement of security**
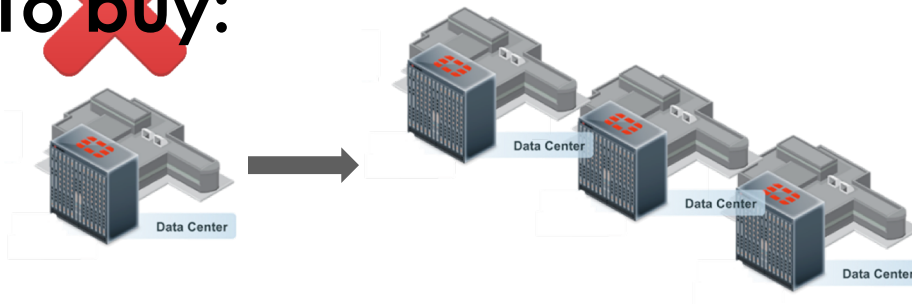  - Protect confidential data

- **Requirement of performance**
  - Maximum uptime
  - Fast page load time
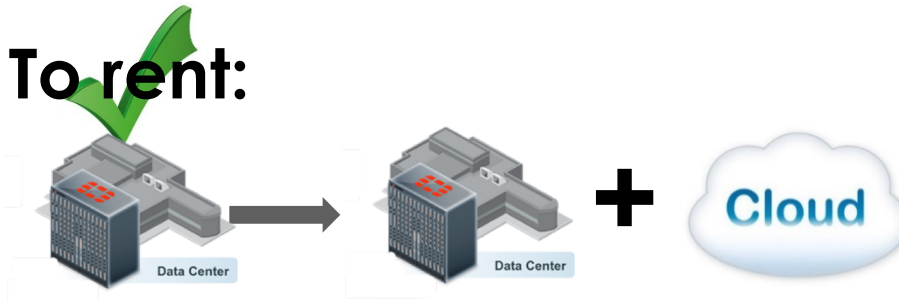
**How to increase capacity and improve scalability?**

# To buy or To rent?

**To buy:**

**To rent:**

- Cost Increases linearly
  - ❑ Infrastructure
- Unable to scale up or down based on workloads
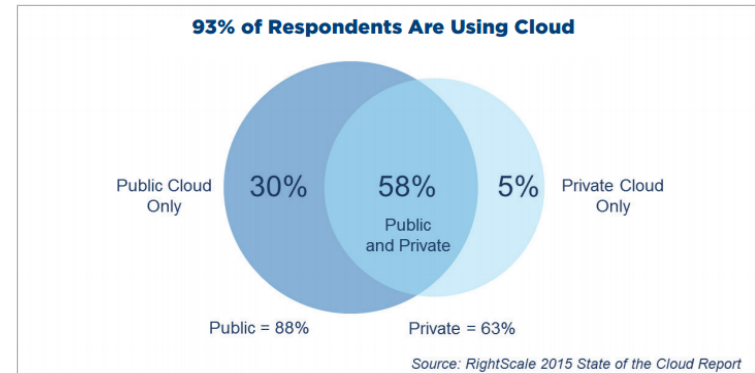  - ❑ Temporary use

- Low price
- Auto scaling
  - ❑ Scalable capacity
  - ❑ Easy to operate
- Potentially unlimited resources

**Hybrid cloud solution is a wise choice!**

# Castle in the air?

- **Hybrid cloud solution is promising and popular**

  - 82% of enterprises have a hybrid cloud strategy, up from 74 percent in 2014
  - 58% of respondents are using hybrid cloud



- **Hybrid cloud solution is already leveraged to handle peak or normal traffic.**
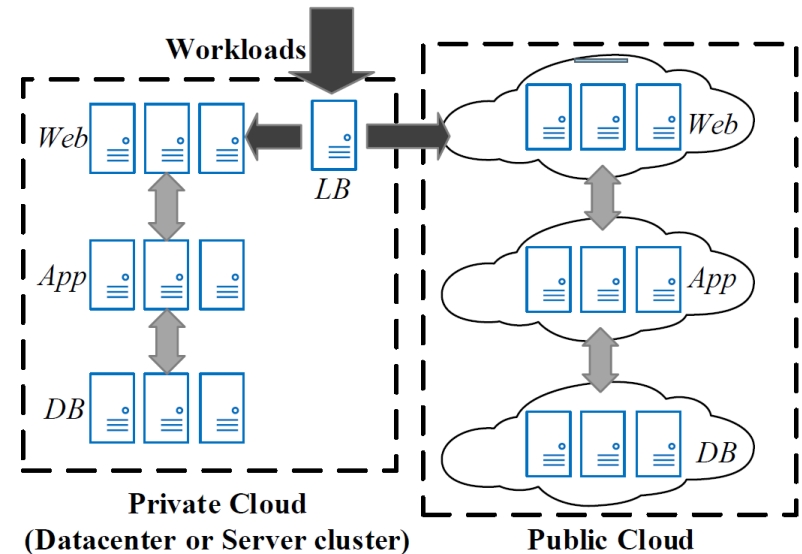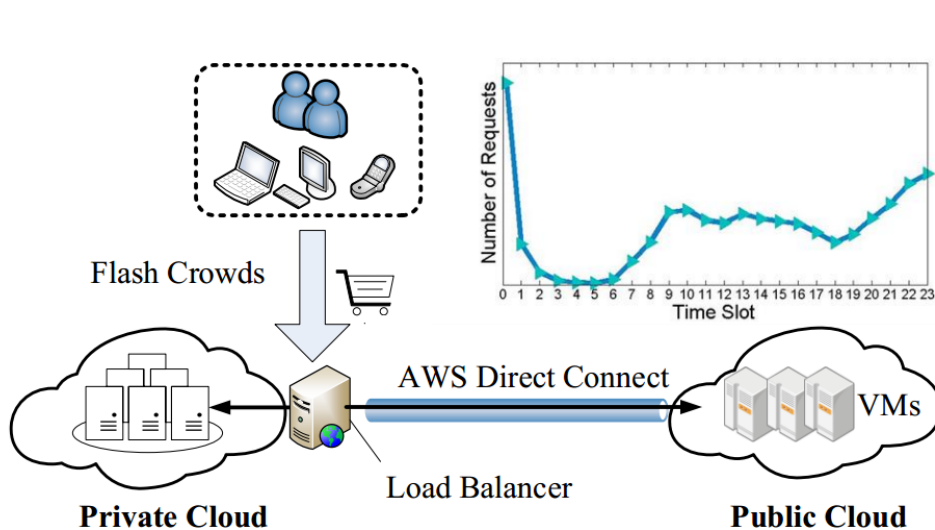
  - **Alibaba**
    - 75% orders were processed by Alibaba cloud on Nov. 11, 2013
    - 96% orders were processed by Alibaba cloud on Nov. 11, 2014
  - **Ebay**
    - 95% of eBay traffic is powered by its OpenStack cloud in 2014.
    - It was zero in 2011.

Source: http://www.computerweekly.com/news/2240222899/Case-study-How-eBay-uses-its-own-OpenStack-private-cloud

# Own the base & Rent the peak



Flash Crowds

Private Cloud

AWS Direct Connect

Load Balancer

Public Cloud

VMs

Workloads

Web

App

DB

LB

Private Cloud (Datacenter or Server cluster)

Web

App

DB

Public Cloud

- **Challenges**
  - Workload distributing
  - Public cloud scaling

- **Architecture**
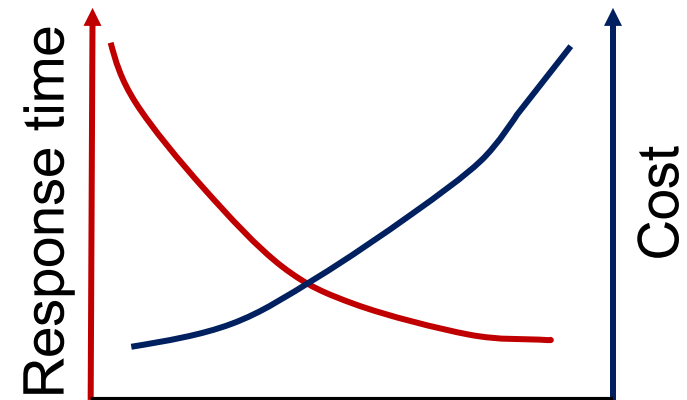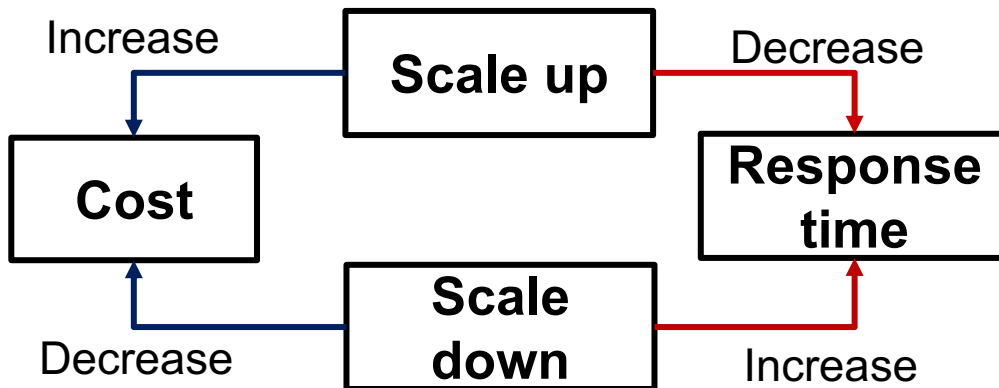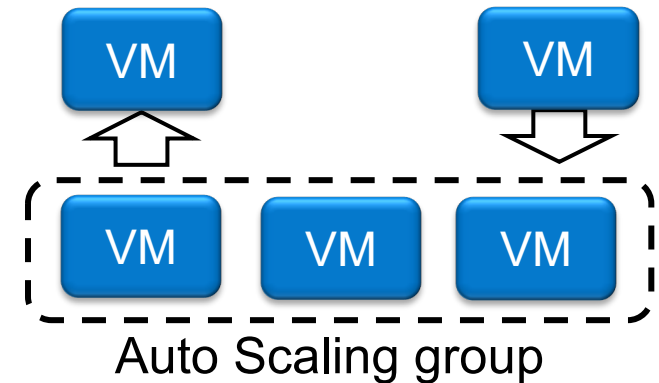  - Three tiers both in private and public clouds
- **Advantages**
  - High flexibility
  - Easy to adopt hybrid cloud

# Trade-off between performance & cost in hybrid cloud

- **Auto Scaling**
    - Enable public cloud to add or remove VMs automatically
    - Workloads assigned to the public cloud is closely related to both performance and cost



Auto Scaling group



**How to provision cost-effective services ?**

# Overview of system model



Flash Crowds

M/M/1 queue

Auto scaling

AWS Direct Connect

Single-tier

Single-tier

Load Balancer

**Private Cloud**

**Public Cloud**

Workload distributing

M/M/1 queue

M/M/1 queue

# Modeling auto scaling



Flash Crowds

Auto scaling

AWS Direct Connect

Load Balancer

**Private Cloud**

**Public Cloud**

VMs

# Auto scaling in public cloud
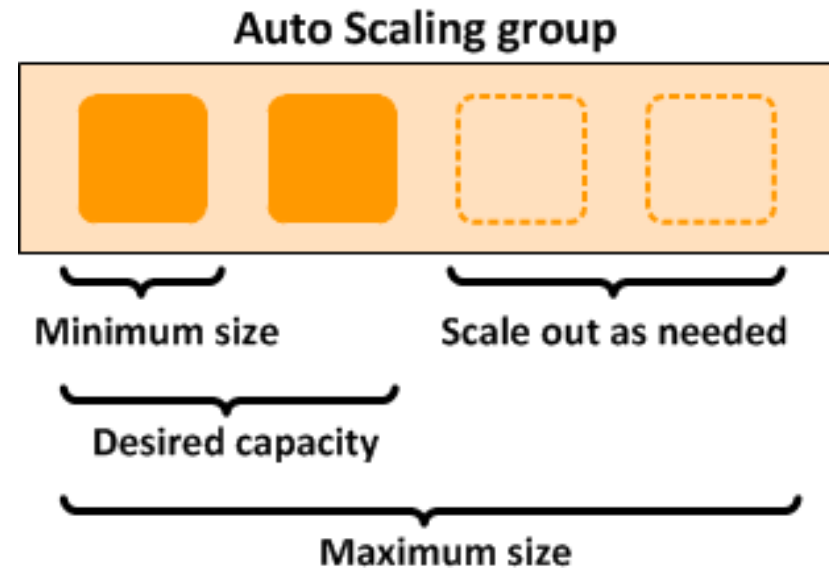
- Auto scaling
  - Allow users to scale public cloud capacity up or down automatically according to defined metric of performance.
  - Number of once scaling
    - Scaling up: m
    - Scaling down: n
  - Monitored metric
    - CPU utilization

**Auto Scaling group**

Minimum size

Scale out as needed

Desired capacity

Maximum size

$$S(t) = \begin{cases} m, & \alpha(t) \geq \alpha_u \\ 0, & \alpha_d \leq \alpha(t) \leq \alpha_u \\ -n, & \alpha(t) < \alpha_d \end{cases}$$

# Modeling websites



**Flash Crowds**

Number of Requests / Time Slot

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23

**AWS Direct Connect**

Single-tier

Single-tier
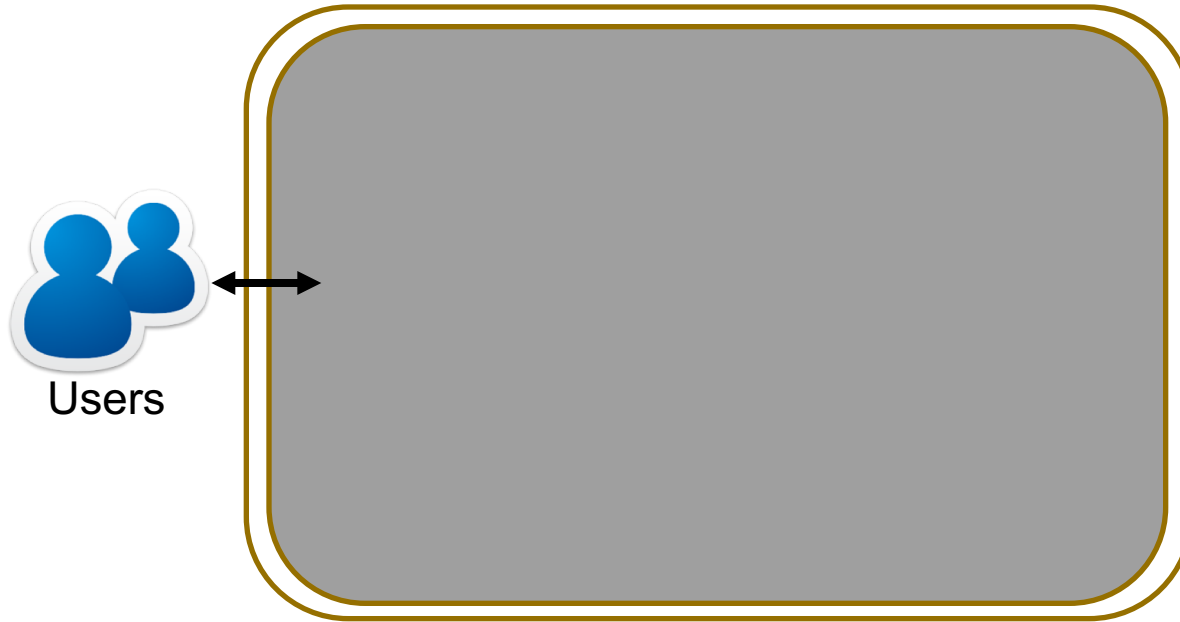
**Load Balancer**

**Private Cloud**

**Public Cloud**

# Single-tier architecture

- **Assumption**
    - The architecture of website is single-tier.



Users

- **Single-tier architecture [4]**
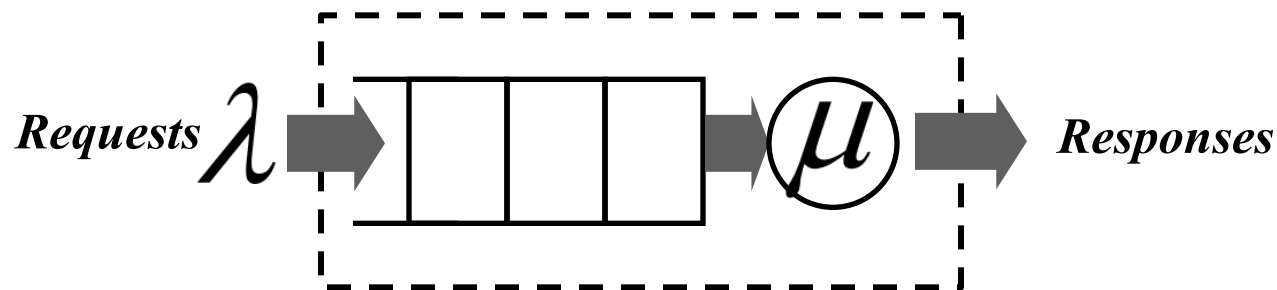    - Encapsulates all the functionalities in the website

# Single-tier architecture

- **Assumption**
  - The architecture of website is single-tier.

$Requests$ $\lambda$ ➤ [ | | ] ➤ $(\mu)$ ➤ $Responses$

- **Modelling**
  - Assumptions
    - Request arrival process is a Poisson process [1] [3]
    - Request serving time is exponential distribution [1] [2]
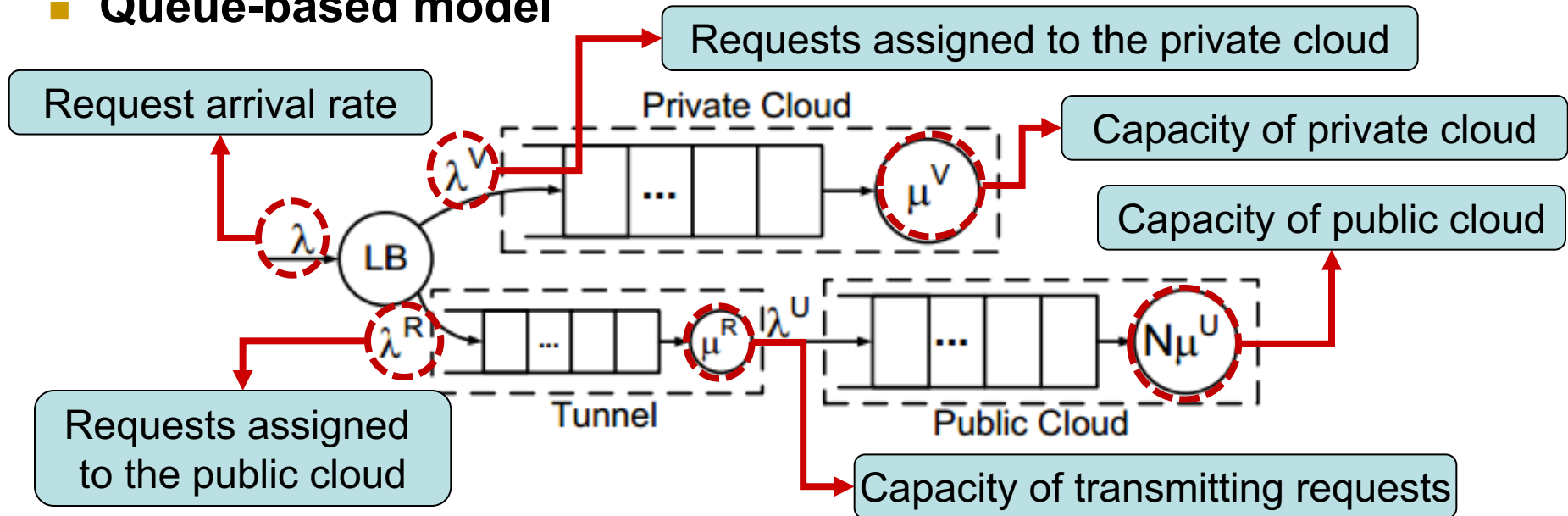  - M/M/1 queue [2]

# Evaluating Performance & Cost

# Evaluating response time in hybrid cloud

- **Queue-based model**

Requests assigned to the private cloud

Request arrival rate

Capacity of private cloud

Capacity of public cloud

Requests assigned to the public cloud

Capacity of transmitting requests



- **Little's Law**
  - Evaluating average response time

$$D^R(t) = \frac{1}{\mu^R(t) - \lambda^R(t)}$$

$$D^U(t) = \frac{1}{N(t)\mu^U - \lambda^U(t)}$$

$$D^V(t) = \frac{1}{\mu^V - \lambda^V(t)}$$

$$D(t) = \frac{\lambda^V(t)}{\lambda(t)}D^V(t) + \frac{\lambda^R(t)}{\lambda(t)}(D^R(t) + D^U(t))$$

# Calculating cost in hybrid cloud

- ## Cost = Cost of tunnel + Cost of VMs

**Levels of available tunnel**

| Level | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| **Bandwidth (Mbps)** | 50 | 100 | 200 | 300 | 400 | 500 | 10,000 |
| **Price ($/hour)** | 0.03 | 0.06 | 0.12 | 0.18 | 0.24 | 0.30 | 2.25 |

- ## Cost of tunnel $K(t) = \sum_{l=0}^{L} k_l x_l(t)$

**Price of AWS EC2 instance**

| Type | vCPU | ECU | Memory | Usage |
|---|---|---|---|---|
| M2.2xlarge | 8 | 26 | 30GB | $0.56/Hour |

- ## Cost of VMs $I(t) = AN(t)$

# How to distribute workloads?



Flash Crowds

AWS Direct Connect

VMs

Load Balancer

**Private Cloud**

**Public Cloud**

Workload distributing

# Problem formulation

- To control cost, we introduce a time-averaged budget M.

- Problem
  - Online, NP-hard and non-linear

$$\min \quad \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} D(t)$$

Minimize response time

$$\text{s.t.} \quad \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \left( AN(t) + AS(t) + \sum_{l=0}^{L} k_l x_l(t) \right)$$

Control cost under budget

$$\lambda(t) = \lambda^V(t) + \lambda^R(t)$$

Workload distributing
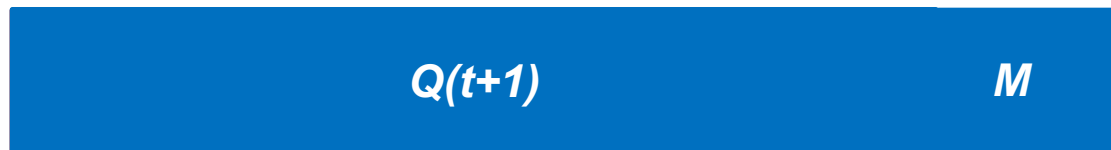
$$\lambda^U(t) < N(t)\mu^U$$

$$\lambda^R(t) < \mu^R(t)$$

Ensure the traffic intensity is below 1

$$\lambda^V(t) < \mu^V$$

# How to minimize response time & control cost?

- To measure how much the cost exceeds the budget, we introduce *Q(t)*.

| Q(t+1) | M |
|:---:|:---:|

- Like a queue, with "cost" coming and "budget" leaving.
- Lyapunov optimization approach can address it
    - Combines response time and cost.
    - Introduces a *V* can be used to determine which part we want to emphasize.

$$\min \ VD(t) + Q(t)\left(AS(t) + \sum_{l=0}^{L} k_l x_l(t) + AN(t) - M\right)$$

# Optimality analysis

- Based on the constraint, there must exist a positive $\epsilon$, which can transform the constraint to the following form:

$$\mathbb{E}\{AN(t) + AS(t) + \sum_{l=0}^{L} k_l x_l(t)\} \leq M - \epsilon$$

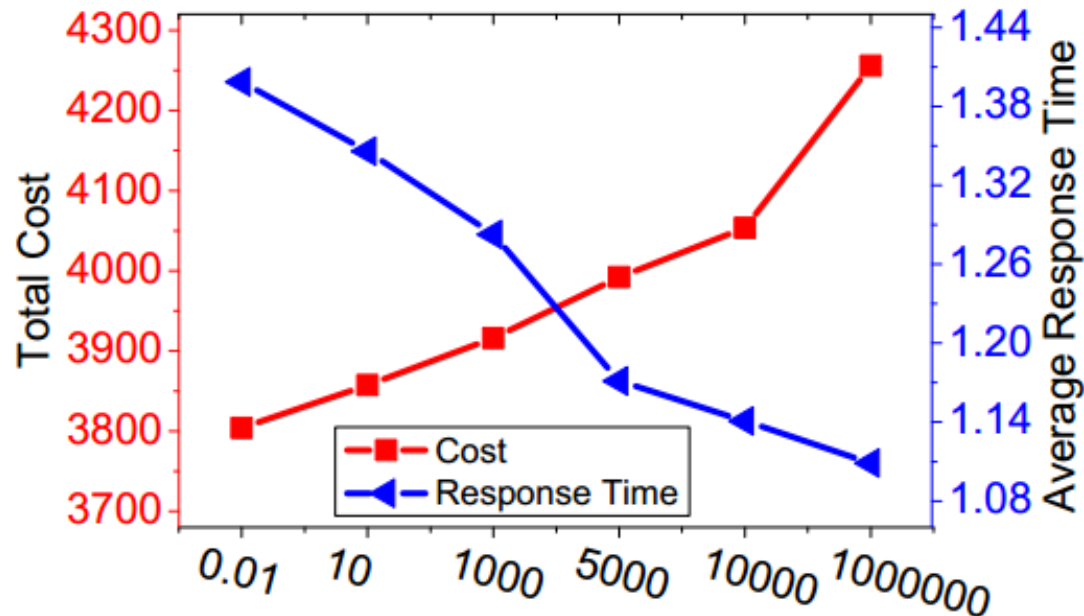- By applying *P\** to *Drift-Plus-Performance*, and summing up it over time slots, and then dividing both sides by *T*, we have

$$\frac{V}{T} \sum_{t=0}^{T-1} D(t) + \frac{L(Q(T)) - L(Q(0))}{T} \leq B + VP^* - \frac{\epsilon}{T} \sum_{t=0}^{T-1} Q(t)$$

- Making $T \to \infty$, as a result, $\frac{L(Q(T)) - L(Q(0))}{T} = 0$

- Note that $\frac{V}{T} \sum_{t=0}^{T-1} D(t) > 0$, we have $\lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} Q(t) \leq \frac{B + VP^*}{\epsilon}$
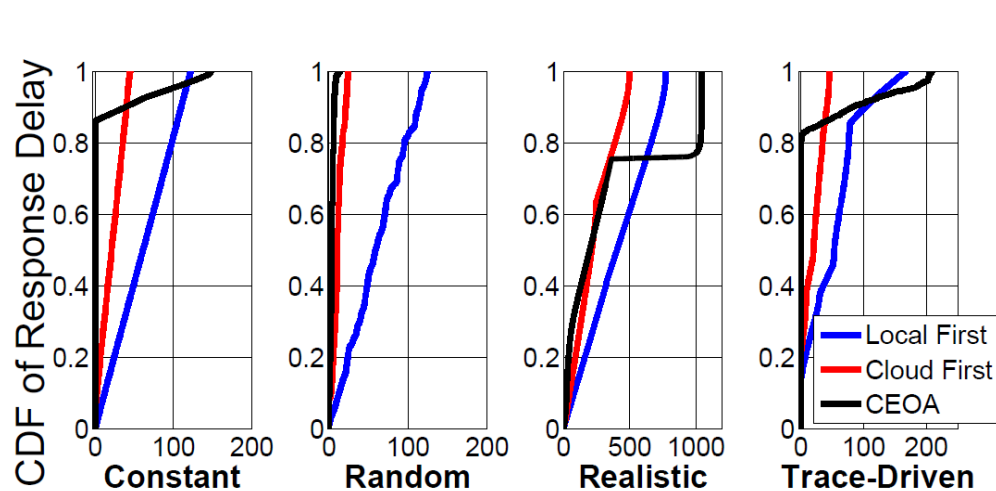
- Note that $-\frac{\epsilon}{T} \sum_{t=0}^{T-1} Q(t) < 0$, we have $\lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} D(t) \leq \frac{B}{V} + P^*$
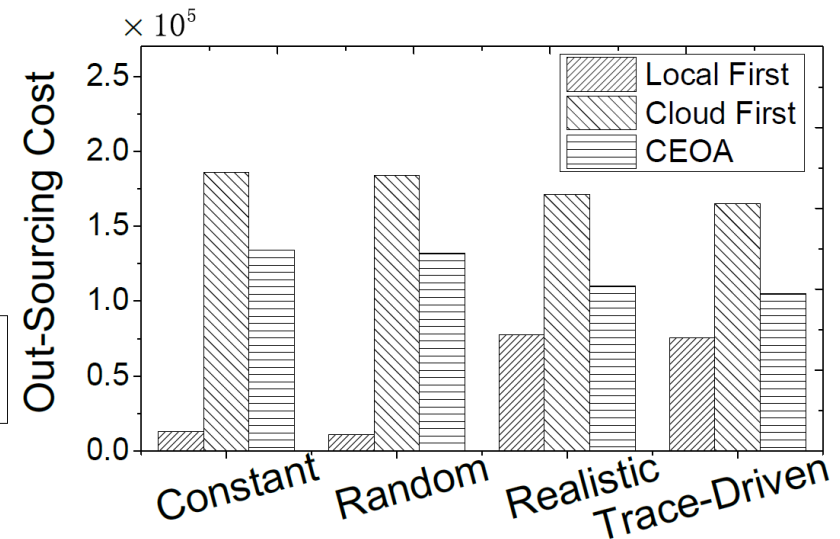
# Tradeoff between performance and cost



- By tuning the value of *V* to a small one, we observe that the average response time is large while the outsourcing cost is small.
- Meanwhile, by setting a large value of *V*, it brings markedly increase of the outsourcing cost and decrease of the average response time.
- When the average response time drops, the outsourcing cost grows remarkably.
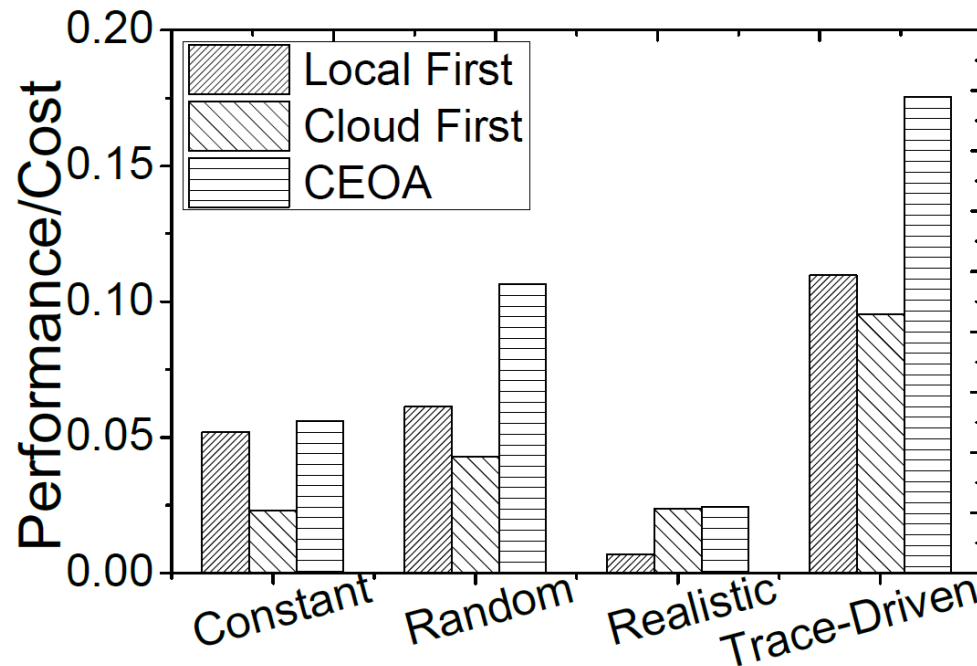
# Performance & Cost



- "Cloud first" strategy provisions the best services while "Local first" strategy provisions the worst services
- CEOA is searching the balancing point of service quality between those two baseline solutions

- The outsourcing cost of the "local first" strategy is the minimum
- The outsourcing cost of the "cloud first" strategy is the maximum
- The cost of CEOA is between the two baseline solutions.

# Performance-Cost ratio



- The "local first" strategy is the most economic
- The "cloud first" strategy helps the website provision the best services
- CEOA has the largest performance cost ratio.

**CEOA, i.e., our solution, enables the website to provision cost-effective services**

# Conclusion

- We design an online algorithm to help an e-commerce website provision cost-effective services.

- By applying Lyapunov optimization approach, our online algorithm can
    - make real time decision on how to offload workloads from a private cloud.
    - prove our online algorithm can approach a dedicated [O(1/V), O(V )] tradeoff between outsourcing cost and average response time.

- Through simulations with empirical real e-commerce PV trace, we demonstrated the effectiveness of our solution

# Reference

| No. | Paper | Source |
|-----|-------|--------|
| [1] | An analytical model for multi-tier internet services and its applications | SIGMETRICS '05 |
| [2] | Performance Guarantees for Web Server End-Systems: A Control-Theoretical Approach | TPDS'02 |
| [3] | Wide area traffic: the failure of Poisson modeling | ToN'95 |
| [4] | Provisioning Servers in the Application Tier for E-commerce Systems | IWQoS'04 |

# Q&A

## *Thank You!*

***"Cloud Datacenter & Green Computing"*** *Research Group*
*Huazhong University of Science & Technology*

http://grid.hust.edu.cn/fmliu/

fmliu@hust.edu.cn