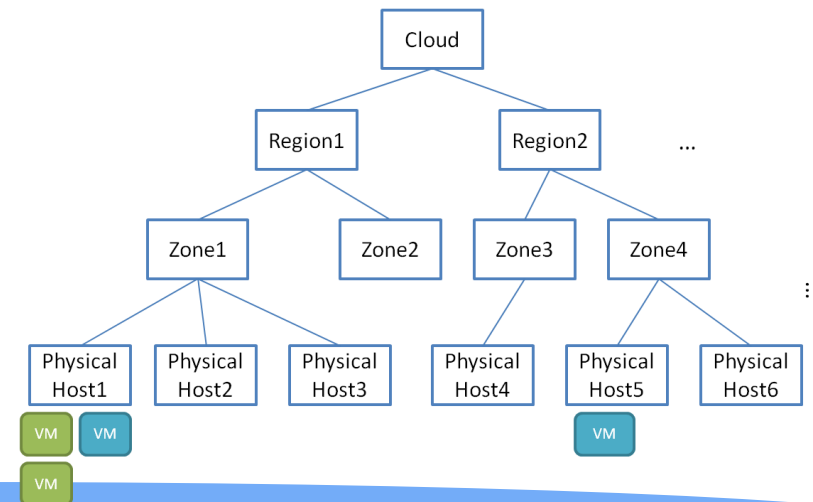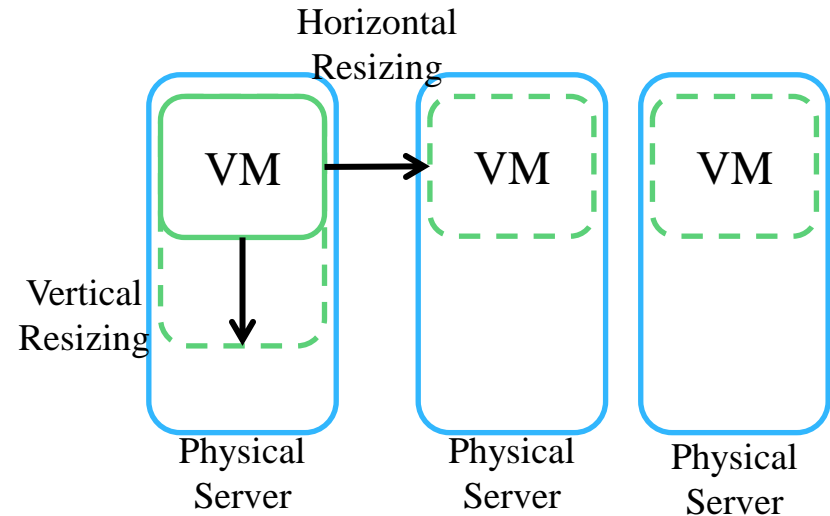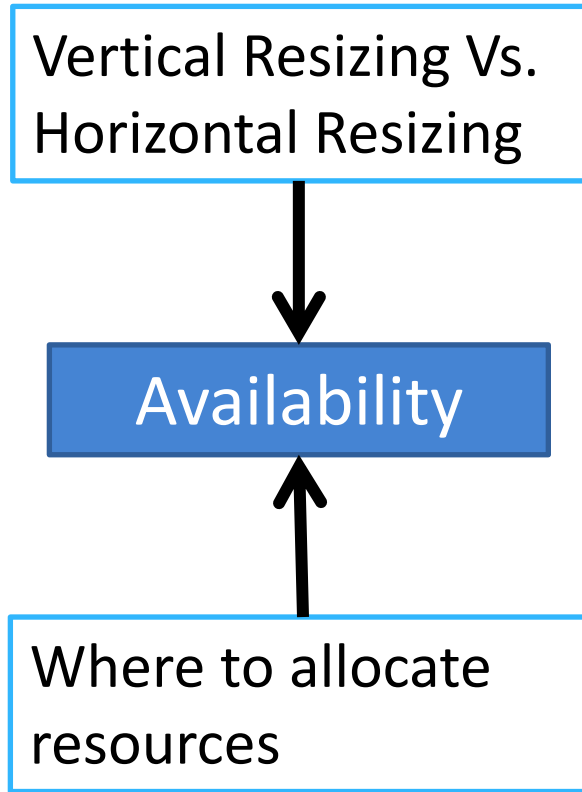# An Availability-aware Approach to Resource Placement of Dynamic Scaling in Clouds

Wenting Wang, Haopeng Chen, Xi Chen

Shanghai Jiaotong University, China

2012-06-24

# Agenda

- Problem Statement

- Modeling and Approach

- Evaluation

# Problem Statement

REliable, INtelligent & Scalable Systems

- Two Problems in scaling resources

# Agenda

- Problem Statement

- Modeling and Approach

- Evaluation

# Modeling and Approach

- Availability Modeling

- ## The Availability of One VM

$$P_i(VM) = P_i \prod_{j \in PP(i)} P_j$$



- ## The Failure Probability of Two VMs

$$\left(\overline{P_u(VM)}\right) \bigwedge \left(\overline{P_v(VM)}\right) = \overline{\prod_{n \in C(u,v)} P_n}$$

$$+ \prod_{n \in C(u,v)} P_n * \left(\overline{\prod_{x \in N(u), x \notin C(u,v)} P_x}\right)\left(\overline{\prod_{y \in N(v), y \notin C(u,v)} P_y}\right)$$
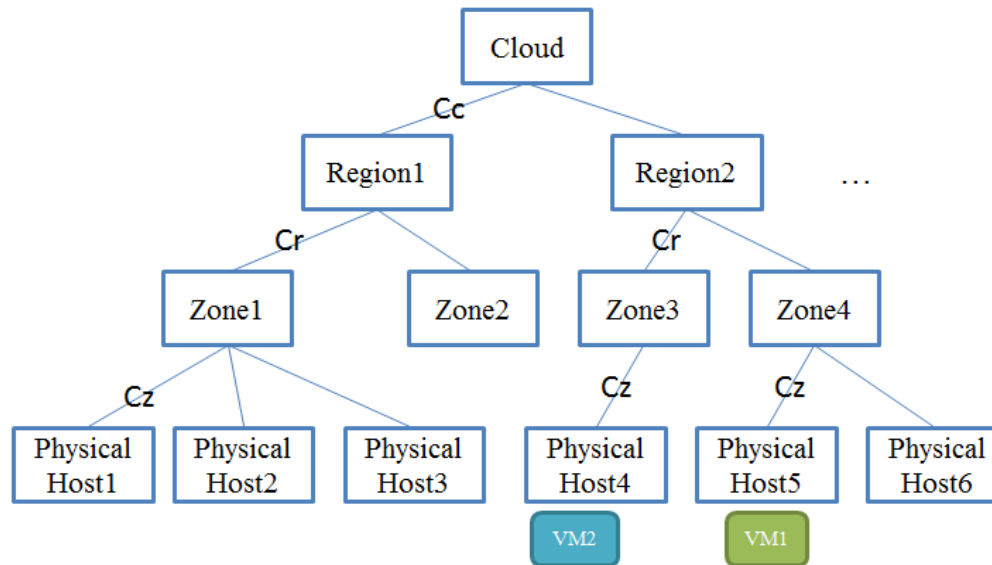
- The Failure Probability of multiple VMs

$$\overline{P'_m} = \overline{P_m} + P_m \prod_{n \in children(m)} \overline{P'_n}$$

> The application availability(denoted as A) is based on the sub-tree generated by multiple VMs:

$$A = 1 - P(\bigcap_{n=1}^{k} \left(\overline{P_n(VM)}\right)$$

- Communication Cost Modeling
  - Let $cc(v_1,v_2)$ donates the communication cost between $VM_1$ and $VM_2$.
  - Then the communication cost from one VM v to the other VMs in an application (where S is the set of all VMs composing the application) is $cc(v, S - \{v\}) = \sum_{x!=v} cc(v, x)$

# Modeling and Approach

REliable, INtelligent & Scalable Systems

```
1: Ac=calculateAvailability();
2: t =1;
3 :For(k=1;k<=Quantity;k++){
4:  If(scale == up){
5:    //the current availability is met
6:    If(Ac>=Ar){
7:       VerticalResizeUp(S, 1);
8:    }
9:    Else{// Ac<Ar
10:       HorizontalResizeUp(S, 1, t);
11       t++;
12:    }
13:  }
14:  Else{ //scale down
15:    VerticalResizeDown(S, 1);
16:  }
17:  Ac=calculateAvailability();
18: } //end for
```

```
19: while(Ac<Ar && relocatedTimes >0){
20:   //rebalance overall application
21:   Relocate(S);
22:   Ac=calculateAvailability();
23:   relocatedTimes --;
24:}
```

Input: **Quantity** denotes demanded unitized resource quantity need to be resized up or down

Input: **Scale**=up/down shows scaling flag

Input: **relocatedTimes** is the max times of relocation

# Agenda

- Problem Statement

- Modeling and Approach

- Evaluation

# Simulation

- Evaluation of Availability Model



(a)Demand satisfaction     (b) Averaged availability     (c)Communication cost Increase ratio

Figure3 availability enhancement and performance change when scaling up



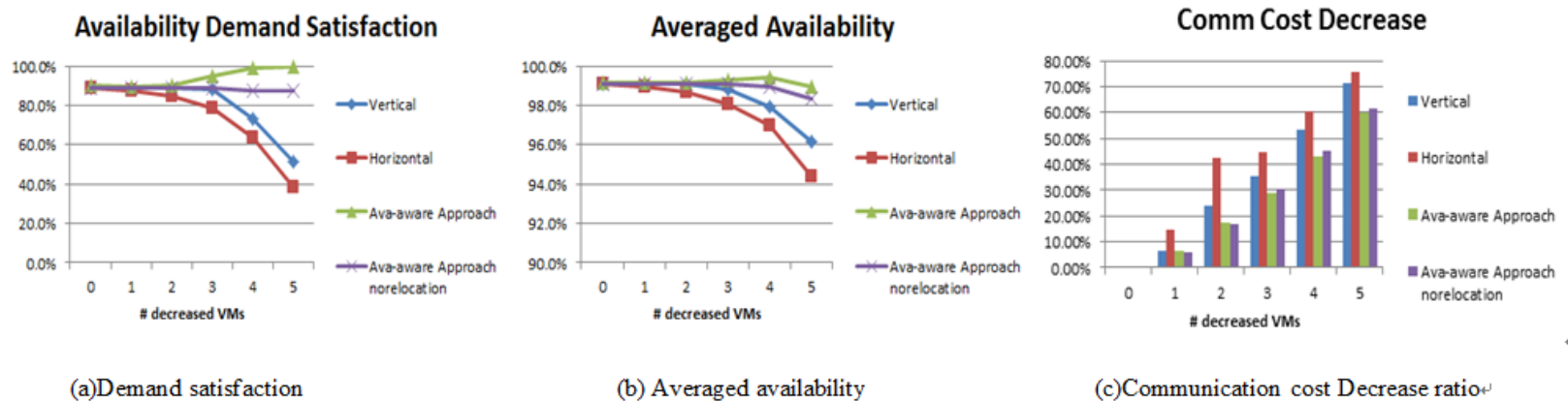(a)Demand satisfaction     (b) Averaged availability     (c)Communication cost Decrease ratio

Figure4 availability declination and performance change when scaling down

- Q & A

- Thank you~