# Survivable and Bandwidth-Guaranteed Embedding of Virtual Clusters in Cloud Data Centers

**Ruozhou Yu**, Guoliang Xue, and Xiang Zhang Arizona State University Dan Li Tsinghua University



## Outlines

#### Introduction and Motivation

System Model and Algorithm Design

Performance Evaluation

Conclusions



## **The Cloud Shift**

Cloud computing: seems an omnipotent solution to all kinds of performance requirements



But is it as mighty as it seems?



## **Inside the Cloud**

An illusion of infinite computing resources created by large clusters of interconnected machines in data centers



#### Performance bottleneck: Cloud network!



#### VM & Bandwidth

Traditional approach: Network-agnostic VM allocation
 Recent advance: Bandwidth-guaranteed VM allocation
 Or Virtual Cluster Embedding (VCE)!



Existing algorithms can allocate bandwidth-guaranteed VMs with minimum bandwidth, migration costs, etc.

But we know that Cloud machines do fail, quite often...



#### **Survivable VCE**

Question: How can we ensure VM availability even when its host machine could fail?

**Answer**: We prepare extra VMs and bandwidth just in case!

**Question**: And how much will that cost us?

**Answer**: No problem! We can minimize that!

**Question**: How are we going to achieve that?

**Answer**: Dynamic programming!



## Outline

Introduction and Motivation

#### **System Model and Algorithm Design**

- Performance Evaluation
- Conclusions



## **Network Topology**

#### Assumption: the DCN has a tree structure

Abstracts many common DCN topologies (FatTree, VL2, etc)





## **VM Survivability Model**

Primary VMs: VMs that are active during normal operations;
 Backup VMs: VMs in standby mode, activated when a primary VM's PM fails

Each backup VM synchronizes the states of multiple primary VMs



Question: Can we find a bandwidth-guaranteed allocation of both primary and backup VMs to cover an arbitrary single-PM failure, with the minimum number of backup VMs?



## **Dynamic Programming for SVCE**

**Given**: topology tree *T*, request  $J = \langle N, B \rangle$ 

#### **Assumption**: single PM failure

Interpretation: a failure can be either within a subtree, or outside a subtree, but cannot be both.

Key observation: each subtree's ability to provide VMs is independent from the rest of the tree, both during normal operations and during an arbitrary failure

Two layers of Dynamic Programming

- Outer DP: DP for entire subtrees
- Inner DP: DP for the first k sub-subtrees of each subtree



## **DP in Details**

- **Outer DP**:  $N_{\nu}[n_0, n_1]$  as the **minimum number of total VMs** needed in subtree  $T_{\nu}$ , to ensure that
  - $T_{v}$  can provide at least  $n_0$  VMs when no failure is in  $T_{v}$ ;
  - $rac{1}{2}$   $T_{v}$  can provide at least  $n_{1}$  VMs when any PM fails in  $T_{v}$ .
- Inner DP: N<sub>v</sub>'[n<sub>0</sub>, n<sub>1</sub>, k] as the minimum number of total VMs needed in the first k subtrees of v, to ensure that
   The k subtrees can provide n<sub>0</sub> VMs when no failure is in them;
   The k subtrees can provide n<sub>1</sub> VMs when any PM fails in them.
- Alternately update the two tables:
  - ✤ N<sub>v</sub>[n<sub>0</sub>, n<sub>1</sub>] depends on N<sub>v</sub>'[n<sub>0</sub>', n<sub>1</sub>', d<sub>v</sub>] (d<sub>v</sub> is the # subtrees under v);
    ♠ N<sub>v</sub>'[n<sub>0</sub>, n<sub>1</sub>, k] depends on N<sub>v</sub>[n<sub>0</sub>'', n<sub>1</sub>''] of lower-layer nodes.























## **Heuristic SVCE**

#### **Optimal DP time complexity:** $O(|V| N^6)$

 $\clubsuit$  where |V| is # tree nodes, N is # requested VMs.

**Question**: Can we find a near-optimal solution with less time?

- ❑ Observation: if we find a normal VCE with N+N' VMs, such that each PM hosts at most N' VMs, then we can always recover from any single PM failure.
- Algorithm: search from N'=1 to N, each time using an existing VCE algorithm to find a VCE with N' extra VMs, and each PM's # VMs is bounded by N'.

#### **Time Complexity**: $O(N \cdot |V| \log |V|)$



## Outline

Introduction and Motivation

System Model and Algorithm Design

#### Performance Evaluation

Conclusions



## **Simulation Setups**

- Tree-structured DCN
  - ✤ 4-layer 8-ary (512 PMs, 73 switches)
  - ✤ 5 VM slots / PM
  - ToR bandwidth: 1 Gbps | Aggr/Core bandwidth: 10 Gbps
- Tenant VCs
  - ✤ 1000 requests
  - ✤ 15 VMs and 300 Mbps per VM, on average
  - Poisson arrivals
- Comparison:
  - ✤ OPT: Optimal DP SVCE algorithm
  - ✤ HEU: Heuristic SVCE algorithm
  - SBS: Shadow-based solution (dedicated VC backup)



#### **Simulation Results: Average VM Usage**





#### Simulation Results: Acceptance Ratio



Average # requested VMs



#### **Simulation Results: Running Time**





## Outline

- Introduction and Motivation
- System Model and Algorithm Design
- Performance Evaluation
- Conclusions



#### Conclusions

#### A first study on Survivable VCE

- ✤ A two-layer optimal DP algorithm
- ✤ A faster near-optimal heuristic algorithm

Discussions

- Extension to tree-like topologies (FatTree, VL2, etc.)
- Extension to cover a constant number of simultaneous failures

#### Future work

- SVCE on generic data center topologies (BCube, JellyFish, etc.)
- Covering link failures in addition to PM failures





# Q&A? THANK YOU VERY MUCH!

#### **Hose Model Bandwidth Guarantee**

 $\Box Request J = <N, B>$ 

♦ N = 7, B = 100 Mbps



Number of VMs  $T_c$  can offer (bandwidth constrained):  $n_c \in [0, 2] \cap [5, 7]$ 



## DP in Details /2

Bandwidth feasible VMs

Outer DP update:
♦ PM level: 
$$N_h[n_0, n_1] = \begin{cases} n_0 & \text{if } n_1 = 0, n_0 \in [0, c_h] \cap \Lambda_h \\ \lambda_h & \text{if } n_1 = 0, n_0 \in \overline{\Lambda}_h, \lambda_h \leq c_h \quad (3) \\ \lambda_h & \text{if } n_1 = 0, n_0 \in \overline{\Lambda}_h, \lambda_h \leq c_h \quad (3) \\ \infty & \text{otherwise} & Bandwidth \\ \text{infeasible VMS} \end{cases}$$
♦ Switch level:  $N_v[n_0, n_1] = \begin{cases} N_v'[n_0, n_1, d_v] & \text{if } n_0, n_1 \in \Lambda_h \\ N_v'[\lambda_v, n_1, d_v] & \text{if } n_0 \in \overline{\Lambda}_h, n_1 \in \overline{\Lambda}_h \\ N_v'[\lambda_v, \lambda_v, d_v] & \text{if } n_0 \in \Lambda_h, n_1 \in \overline{\Lambda}_h \\ N_v'[\lambda_v, \lambda_v, d_v] & \text{if } n_0, n_1 \in \overline{\Lambda}_h \end{cases}$  (4)

Inner DP update: No subtree:  $N'_v[n_0, n_1, 0] = \begin{cases} 0 & \text{if } n_0 = n_1 = 0 \\ \infty & \text{otherwise} \end{cases}$ 



