

IEEE INTERNATIONAL CONFERENCE ON COMMUNICATIONS

COMMUNICATIONS: CENTREPOINT OF THE DIGITAL ECONOMY

### h-DDSS: Heterogeneous Dynamic Dedicated Servers Scheduling in Cloud Computing

CONNECT WITH IEEE ICC:

lin

Husnu Saner Narman Md. Shohrab Hossain Mohammed Atiquzzaman

School of Computer Science University of Oklahoma, USA. <u>atiq@ou.edu</u> <u>www.cs.ou.edu/~atiq</u>

June 2014



#### What is Cloud Computing





### Why Cloud Computing

- Simplicity
  - No need to set up software/hardware
- Flexibility
  - Easily extending memory/CPU capacity
- Maintenance
  - IT services
- Time and energy
  - No time or extra effort for desired environment
- Pay as you go

No need to pay for unused hardware or software



# What is Cloud Scheduling





#### **Customer Type**

- Different customers classes?
  - Paid and non-paid
- Customer requirements
  - Desired Platform based on Service Level Agreement
- How to satisfy different customer classes?
  - Reserve servers for each customer types
    - Dedicated Servers Scheduling
  - Priority
    - High or Low



#### **Customer Priority**







Without priority level in queuing theory

With priority level in cloud computing



#### **Reserved Servers**





#### **Dedicated Servers Scheduling**



## Q

#### **Dedicated Servers Scheduling**





#### Problems with DSS

- Does not dynamically update number of servers for each group
  - If arrival rate changes
  - If priority level changes
- Servers are homogeneous (Unrealistic)



The University of Oklahoma



#### **Dynamic Dedicated Servers Scheduling**



The University of Oklahoma



#### **Dynamic Dedicated Servers Scheduling**





#### **Problems with DDSS**

Servers are homogeneous (Unrealistic)





### Why Heterogeneous

 Failed or misbehaved servers of a multiserver system are replaced by new and more powerful ones



#### Heterogeneous Servers





#### Objective

- Improve performance of cloud systems for heterogeneous servers
  - Allowing heterogeneous servers to be dynamically allocated to customer classes based on
    - Priority level.
    - Arrival rate.



#### Contribution

- Propose Heterogeneous Dynamic Dedicated Servers Scheduling.
- Develop Analytical Model to evaluate performance
  - Average occupancy
  - Drop rate
  - Average delay
  - Throughput
- Comparing performance of
  - Heterogeneous Dynamic Dedicated Servers Scheduling
  - Dynamic Dedicated Servers Scheduling.



#### Heterogeneous Dynamic Dedicated Servers Scheduling



The University of Oklahoma



#### Heterogeneous Dynamic Dedicated Servers Scheduling





#### **Dynamic Approach**





#### **Modeling Assumptions**

- System is under heavy traffic flows.
- Arrivals follow Poisson distribution, and service times for customers are exponentially distributed.
- Type of queue discipline used in the analysis is FIFO.
- Service rate of all servers can be different.

 $\lambda_1$ : Arrival rate

of  $C_1$  customers



#### **Analytical Model**

- Only  $C_1$  customers performance metric developed.
- Markov Chain Model :





#### Performance

• Drop Probability :  $D = p_0 \frac{\mu_{tm}^{m+N} \rho^{m+N}}{\frac{m}{m}}$ 

$$p_0 rac{\mu_{tm}^{m+N} \rho^{m+N}}{\prod\limits_{j=1}^m \mu_{tj}}$$

#### Drop probability

Rate of dropped customers from the systems buffer.

Number of customers served in the systems.

Average waiting time

of a customer in the

systems buffer.

• Throughput:  $\gamma = \lambda_1(1 - D)$ 

• Occupancy: 
$$n = \begin{cases} p_0 \frac{\mu_{tm}^m}{\prod \mu_{ti}} \rho^{m+1} \left( \frac{1 - (N+1)\rho^N + N\rho^{N+1}}{(1-\rho)^2} \right) & \rho \neq 1 \\ \prod_{i=1}^{m} \mu_{ti}} p_0 \frac{\mu_{tm}^m}{\prod_{i=1}^m \mu_{ti}} \left( \frac{N(N+1)}{2} \right) & \rho = 1 \\ p_0 \frac{\mu_{tm}^m}{\prod_{i=1}^m \mu_{ti}} \left( \frac{M(N+1)}{2} \right) & \rho = 1 \end{cases}$$

• Delay: 
$$\delta = \frac{n}{\gamma}$$



#### Results

- We have used discrete event simulation to implement by following  $M/M_i/N/N$  and proposed scheduling.
- Each queue holds 30 customers.
- We ran simulation with 20000 customers for each arrival rate.
- We show h-DDSS with Fastest Server First (FSF) and Slowest Server First (SSF) to compare best and worst performance.



#### Traffic Arrival Rates

- Simulations were carried out with increased arrival rates of all types of customers to observe the impact of heavy traffic on the system.
- Customer arrival rates at different trials:

$$\lambda_1 = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\},\ \lambda_2 = \{2, 4, 6, 8, 10, 12, 14, 16, 18, 20\},\ \Psi_1 = \{2, 3\}, \Psi_2 = \{1\}$$
 and  
 $\mu = 1, 2, ... 7$  for heterogeneous servers and  
 $\mu = 4$ , for homogeneous servers with 7 servers

The University of Oklahoma



#### Validation of Analytic Formulas: Occupancy



Occupancy model matches with simulation.

The University of Oklahoma



#### Validation of Analytic Formulas: Throughput



Throughput model matches with simulation.



h-DDSS is heterogeneous.

#### h-DDSS vs DDSS DDSS is homogeneous.



DDSS shows better occupancy than h-DDSS for these priority levels.

The University of Oklahoma



h-DDSS is heterogeneous.

#### h-DDSS vs DDSS DDSS is homogeneous.



h-DDSS shows better occupancy than DDSS for these priority levels.



h-DDSS is heterogeneous.

#### h-DDSS vs DDSS is homogeneous.



DDSS shows better throughput than h-DDSS for these priority levels.



h-DDSS is heterogeneous.

#### h-DDSS vs DDSS is homogeneous.



h-DDSS shows better throughput than DDSS for these priority levels.



#### Summary of Results

- Priority levels do not affect the performance of DDSS and h-DDSS under low traffic.
- Under heavy traffic, priority levels have a significant impact on the class performances of DDSS.
- Under heavy traffic, performances of FSF and SSF in h-DDSS are same while FSF is better for low traffic arrivals.
- h-DDSS can be more efficient than DDSS for selected class priority levels



#### Conclusion

- We have proposed a novel scheduling algorithm for cloud computing considering priority, arrival rate and heterogeneous servers.
- Performance metrics of the proposed cloud computing system are presented through different cases.
- h-DDSS and DDSS are compared under different priority levels.
- Proposed scheduling algorithm can help Cloud Computing with homogenous and heterogeneous servers systems have higher throughput and be more balanced.





#### Thank You



http://cs.ou.edu/~atiq atiq@ou.edu