

Power-Efficient Server Utilization in Compute Clouds

Overview

1. Motivation
2. SPECpower benchmark
3. Load distribution strategies
4. Cloud configuration
5. Results
6. Conclusion

2/14



1. Motivation

- Huge power consumption of data centers
- Often underutilized resources
- Smart load distribution (of the virtual machines) should minimize the total power consumption

3/14



2. SPECpower Benchmark

- Server Side Java (SSJ): SPECpower_ssj2008
- Java Apps with discrete load levels (0,10%,..., 100%)
- Measured power consumption of commonly used servers

SPEC Standard Performance Evaluation Corporation

4/14



SPEC data for a server

Performance			Power	Performance to Power Ratio
Target Load	Actual Load	ssj_ops	Average Active Power (W)	
100%	99.7%	1,435,697	315	4,564
90%	90.1%	1,297,557	290	4,482
80%	79.9%	1,150,879	240	4,793
70%	70.0%	1,007,940	218	4,618
60%	60.1%	865,402	201	4,296
50%	50.0%	720,975	182	3,971
40%	39.9%	575,426	159	3,630
30%	30.0%	432,447	137	3,148
20%	20.0%	288,246	117	2,454
10%	10.1%	145,611	103	1,419
Active Idle		0	69.4	0
			Σ ssj_ops / Σ power =	3,900

5/14



3. Load distribution strategies

Four strategies for *divisible* loads:

1. Relative load balancing
2. Absolute load balancing
3. Best performance to power first
4. Adaptive load distribution

6/14



Relative load balancing

- Compute the percentage of requested load to total maximum performance

$$\alpha = \frac{o}{\sum_{i=1}^n O_i^{max}}$$

- Assign to each server

$$O_i = \alpha \cdot O_i^{max}$$



Absolute load balancing (alb)

- Assigns the same amount o/n to each server
- If $o_i > o_i^{max}$: excess load $o_i^{max} - o_i$ of server i will be equally distributed to more powerful servers

8/14



Best performance to power first (bpppf)

- Sort servers regarding their performance to power ratio (100% load) in decreasing order
- Assign to each server its o_i^{max} until the requested load o is completely distributed
- Remaining servers will be idle

9/14



Adaptive Load Distribution (ald)

- Find best scaling for each server by minimizing

$$p_{total}(o) = \sum_{i=1}^n p_i(o \cdot \alpha_i)$$

- Constrains

$$\alpha_i \cdot o \leq o_i^{max}$$

$$\sum_i \alpha_i = 1$$

$$o \leq \sum_i o_i^{max}$$

10/14



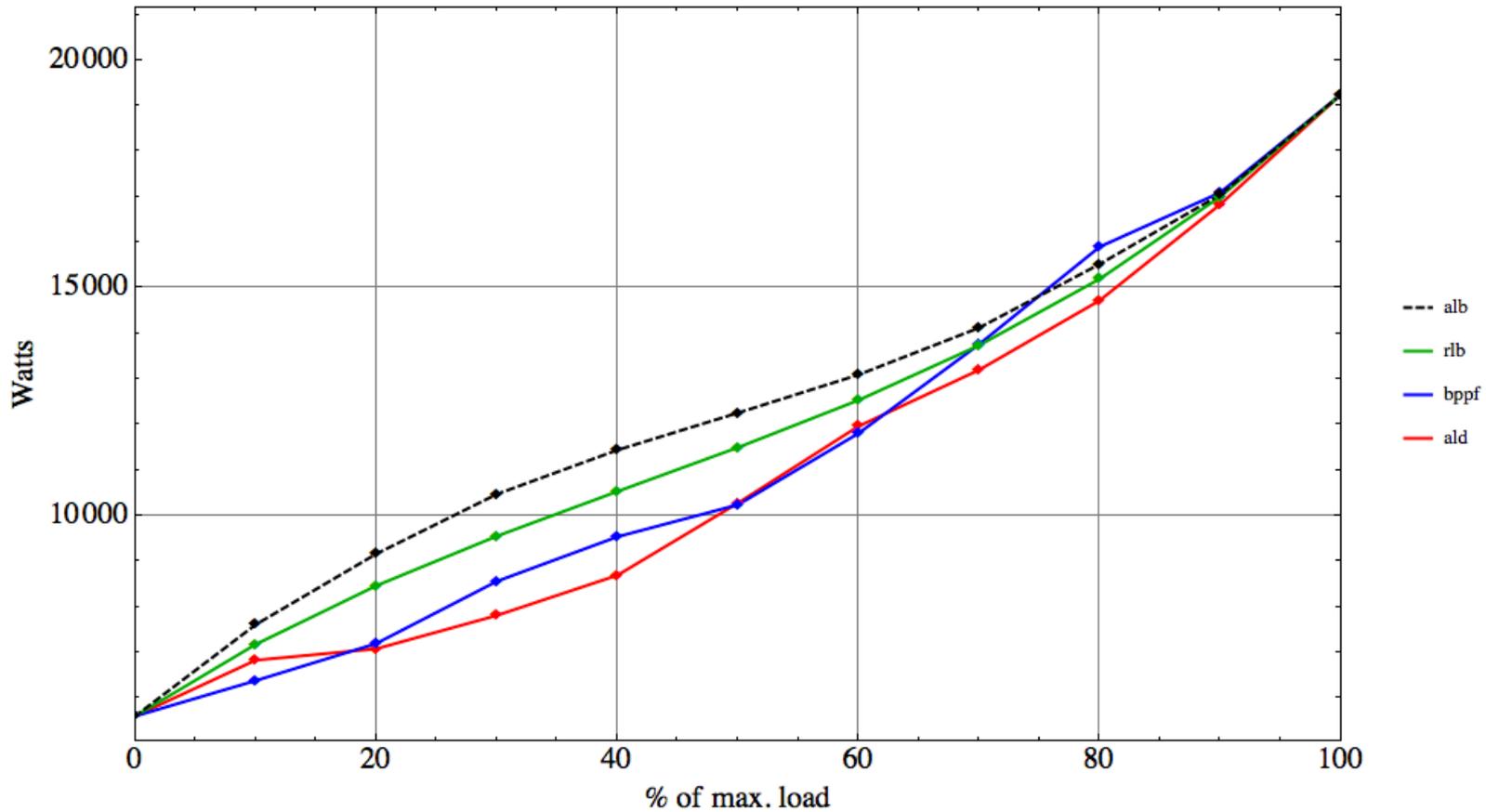
4. Cloud configuration

- Heterogeneous Server-Hardware Nov. 2006 – Okt. 2012
- 12 Single Server (mainly Intel Xeon-based)
- 4 Cluster System (1x4, 1x16, 1x18, 1x32)
- SPEC Performance/Power values: 268 – 5521
- Power consumption: 5,3 – 19,1 kW

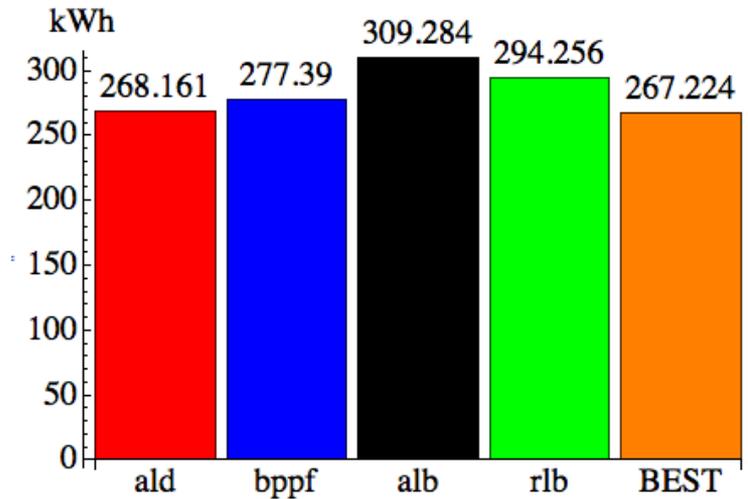
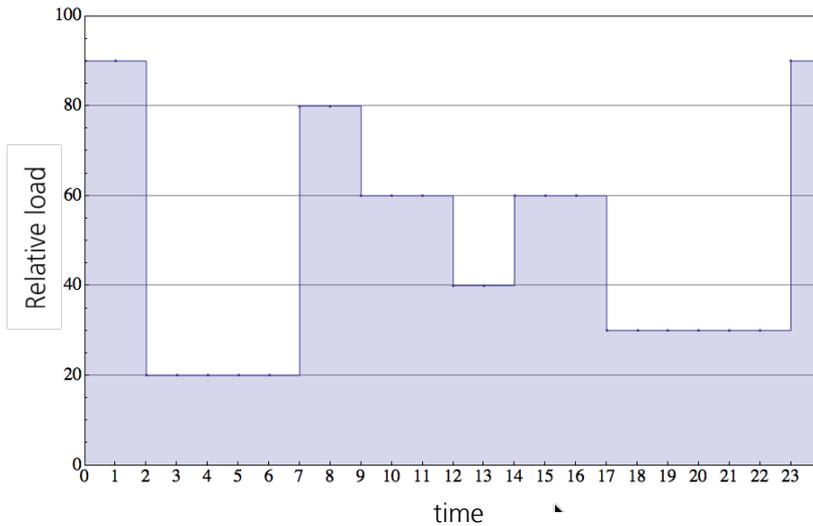
11/14



5. Results



Energy consumption for one day



-15% saving to alb

If best solution for each load used

13/14



6. Conclusion

- Assuming divisible loads in cloud environments described by SPECpower data
- Four different load distribution strategies
- Most sophisticated one is close to the best solution improves power-efficiency just by approx. 9%
- Power consumption of idle nodes
- Load-controlled switching of the servers could help to further improve power-efficiency

14/14

