

Harmony: Dynamic Heterogeneity-Aware Resource Provisioning in the Cloud

Abstract:

Data centers today consume tremendous amount of energy in terms of power distribution and cooling. Dynamic capacity provisioning is a promising approach for reducing energy consumption by dynamically adjusting the number of active machines to match resource demands. However, despite extensive studies of the problem, existing solutions for dynamic capacity provisioning have not fully considered the heterogeneity of both workload and machine hardware found in production environments. In particular, production data centers often comprise several generations of machines with different capacities, capabilities and energy consumption characteristics. Meanwhile, the workloads running in these data centers typically consist of a wide variety of applications with different priorities, performance objectives and resource requirements. Failure to consider heterogeneous characteristics will lead to both sub-optimal energy-savings and long scheduling delays, due to incompatibility between workload requirements and the resources offered by the provisioned machines. To address this limitation, in this paper we present HARMONY, a Heterogeneity-Aware Resource Management System for dynamic capacity provisioning in cloud computing environments. Specifically, we first use the K-means clustering algorithm to divide the workload into distinct task classes with similar characteristics in terms of resource and performance requirements. Then we present a novel technique for dynamically adjusting the number of machines of each type to minimize total energy consumption and performance penalty in terms of scheduling delay. Through simulations using real traces from Google's compute clusters, we found that our approach can improve data center energy efficiency by up to 28% compared to heterogeneity-oblivious solutions.