

ELASTIC SCHEDULING FOR MICROSERVICE APPLICATIONS IN CLOUDS

ABSTRACT

Microservices are widely used for flexible software development. Recently, containers have become the preferred deployment technology for microservices because of fast start-up and low overhead. However, the container layer complicates task scheduling and auto-scaling in clouds. Existing algorithms do not adapt to the two-layer structure composed of virtual machines and containers, and they often ignore streaming workloads. To this end, this project proposes an Elastic Scheduling for Microservices (ESMS) that integrates task scheduling with auto-scaling. ESMS aims to minimize the cost of virtual machines while meeting deadline constraints. Specifically, we define the task scheduling problem of microservices as a cost optimization problem with deadline constraints and propose a statistics-based strategy to determine the configuration of containers under a streaming workload. Then, we propose an urgency-based workflow scheduling algorithm that assigns tasks and determines the type and quantity of instances for scale-up. Finally, we model the mapping of new containers to virtual machines as a variable-sized bin-packing problem and solve it to achieve integrated scaling of the virtual machines and containers. Via simulation-based experiments with well-known workflow applications, the ability of ESMS to improve the success ratio of meeting deadlines and reduce the cost is verified through comparison with existing algorithms.