

Optimizing Shared Resource Contention in HPC Clusters

Sergey Blagodurov
Alexandra Fedorova

sergey_blagodurov@sfu.ca
alexandra_fedorova@sfu.ca



The problem: contention for shared multicore resources (shared caches, memory controllers, NUMA domains, etc.) within cluster nodes incurs up to 40% severe degradation to job performance.
HPC clusters are not contention-aware, with no virtualization to migrate jobs around.

The solution: Clavis-HPC, a novel contention-aware virtualized HPC framework. Here is how:

- 1) We monitor job processes on-the-fly and classify them with 2 parameters:
 - a) a process is a **devil** if it is memory intensive, has high last-level cache missrate, otherwise - a **turtle**.
 - b) if a given process is communicating with other processes.
- 2) We develop a multi-objective scheduling algorithm Clavis-Cluster that simultaneously:
 - a) **minimizes** the number of devils on each node;
 - b) **maximizes** the number of communicating processes on each node;
 - c) **minimizes** the number of powered up nodes in the cluster.
- 3) After the new schedule is found, we enforce it by introducing a low-overhead live migration into cluster: the job scheduler places processes into OpenVZ containers, Clavis-Cluster migrates containers.



A typical HPC job management cycle with our modifications highlighted in red: (state-of-the-art → **Clavis-HPC**)

