

TS-Bat: Leveraging Temporal-Spatial Batching for Data Center Energy Optimization

GLOBECOM 2017, Singapore

Motivation

Server Farms are generally over-provisioned

CPUs and other hardware are configured to satisfy peak load
 On average, typical server farm utilization is < 30%

Current server farms lack efficient power management

CPUs are kept active more than necessary

Only shallow power saving states (e.g., Core C1 state) are used

A more effective energy optimization requires

A better characterization of multi-core power consumption
 Energy-aware scheduling and control of processor low power states
 Workload adaptivity (QoS constraints, job arrival rate & service time)



11/20/21

Understanding Multi-core Processor Power







Temporal Batching – Front End



Job Batch Processing – Back End



Temporal Batching Problem Formulation

Determine the max number of jobs to be batch

For each job i

Batching delay: B_i = (K–i)/λ + σ
 Queuing delay: U_i =S⁹⁵*(i–1)/C+W
 Total Delay: D_i = B_i + U_i

■ Find max *K* so that that $\forall i \in [1, K]$: $D_i + S^{95} \leq Q$

- λ : job arrival rate
- C : number of cores per server
- Q : target tail latency constraint
- S : service time distribution,
- S⁹⁵: 95th percentile service time
- σ :processor wakeup latency

Integration with Spatial Batching

Group jobs onto specific servers Instead of evenly distributing load... Keep a subset of CPUs in active state Create opportunities for other servers to enter deep sleep states Cuts down power significantly

Challenges: How to choose idle and active servers?



8

Temporal + Spatial Batching – Scale-up

9





Experimental Setup

Testbed

- Temporal Batching single server
 - Intel Xeon E5 10 core process
- Temporal Batching + Spatial Batching
 - Cluster of 16 servers
 - Each server equipped with Intel Xeon 5650 6-core server
 - Network : Netgear 24-port Gigabit switch (start topology)

Software Framework

- Load generator: *httperf*
- Application server: apache httpd
- Benchmark: CGI script running PARSEC-2.1 benchmarks at back-end

The targeted QoS constraint

- 95th percentile latency* is less than 2× job service time *

*95th Percentile latencies are adopted as measurements of QoS constraints in Latency critical applications – Rubik, Micro 2014; PEGASUS, ISCA 2014; Tailbench, ASPLOS 2016



Temporal Batching – C State Residency





Temporal Batching – Response Time



- \checkmark For QoS 10x: 95th response time 925 ms < target latency 1080 ms
- Quality of Service constraints are effectively maintained



Energy Savings



GW

Conclusion

We characterized the power and performance of multicore processor in data centers.

• We proposed TS-Bat, a QoS-aware energy optimization framework that combines temporal and spatial batching.

TS-Bat achieved upto 68% processor energy saving under various QoS constraints.



11/20/21