

Thermostat heterogeneous
information selective interconnect
technologies placement hardware
overhead DDR
use CPU rate page systems
application coherence using high
ratio two memories pages without policy GDDR
while
TLB shown capacity each GPU support virtual
threshold OS client only between
across
hot any run cold GPUs
Online per slow shared
software cloud some Linux
huge KB request number access caches
Section table cost first range CPU-GPU accessed
higher fault workloads runtime applications
because data cache address requests
latency more slowdown footprint shows
accesses mechanism throughput Computer
total fraction performance Conference
degradation Architecture expansion DRAM
system bandwidth
BW-AWARE migration