Workloads, Scalability and QoS Considerations in CMP Platforms

Presenter – Don Newell Sr. Principal Engineer Intel Corporation

© 2007 Intel Corporation



- Trends and research context
- Evolving Workload Scenarios
- Platform Scalability
- Platform QoS



Trends

- Cores, Threads, and Virtual Machines



Trends - Discussion Points

<u>Goal</u>: Scaling up to 100's of logical processors on a single CPU die

- Scaling hardware threads means provisioning and re-architecting other platform resources
- Scaling hardware threads means learning to share – QoS revisited



Trends - Pool of Virtual Resources





Trends - Decomposed OS



Platform Scaling



Tera-Scale Workload Scenarios



Tera-Scale Cache Hierarchy Design



Hierarchical sharing increases caching effectiveness significantly

OLTP Tera-Scale Case Study



Significant Performance Potential => Memory Scalability Challenges

Scalability Issues

Tera-scale Headroom requirements

- Start with 32 cores in 1st gen; Maybe grow to 48 cores in next gen?
- How much memory & interconnect bandwidth will we need?



On-Socket DRAM Caches (For Memory Scalability)

Enable Large Capacity L4s

- Low Latency
- High Bandwidth
- Technologies
- 3D Stacking
- Multi-chip Packages (MCP)
 Benefits
- Significant reduction in miss rate
- Avoids bandwidth wall





12

QoS and Performance Management



Background for QoS Discussion



Observations

Multi-core enables simultaneous execution of multiple workloads

But not all Workloads are equal -- users do have preferences

How well does the user-preferred application run? Should platforms optimize for the user-preferred application?

Resource Management

Capitalist

No management of resources Fair distribution of resources If you can generate more Give equal share of resources to requests, you will use more all executing threads resources Does not necessarily guarantee equal performance Grab as you will E.g. Partitioning resources for fairness and isolation • E.g. All of today's policies Elitist Utilitarian Focus on individual efficiency Focus on overall efficiency Provide more performance and Give more resource to those that resources to the VIP need it the most, less to others Limited resources to non-VIP E.g. Cache-friendly vs. Unfriendly, resource-aware scheduling E.g. Service Level Agreements, Foreground/Background

Communist/Fair

The Multi-Workload Problem



Execution Time of a foreground application Time Significant response time slowdown Execution 50%

Conroe core 2 Duo Measurements

Platform does not distinguish in resource allocation

Preferred (foreground) application can suffer significant slow down

Platform QoS can improve user-preferred application performance

Contention - Client side examples



Measured on a Dual core Conroe (4M cache, DDR2 667)

Similar data collected for server applications

Resource contention will impair the performance of important apps

(Performance Differentiation)

Application Behavior & Overall Performance



Potential to improve overall performance

Monitor resource usage and group/partition accordingly



Clovertown (8Core / 8App) Experiments (used destructive, normal and neutral apps)

Managing resource contention can improve overall throughput too

(Performance Management)

Service Level Agreements in the Enterprise

Server Consolidation



Disparate resource usage and contention hurts SLAs

(Need for Performance Isolation and SLA enforcement)

Shared Platform Resources



Problem Summary

CMP Amy heterogeneous threads, apps, VMs

Not all applications are equal

- Users have preferences
 - End-users (client) want to treat foreground preferentially
 - End-users (server) want service level differentiation (SLA) or isolation

Applications use resources differently

- Destructive vs. Constructive vs. Neutral Threads
- No performance management to protect from bad behavior

Priority-based OS scheduling no longer sufficient

- With more cores, OS will allow high and low priority applications to run simultaneously and contend for resources
- Low priority applications will steal platform resources from high priority apps → loss in performance & user experience

Platform has no support for application differentiation

- Platform has <u>no knowledge</u> of preferences or resource usage
- Platform has <u>no support</u> for fine-grained tracking of many shared resources that have significant performance value

21

Platform QoS





Goals – Preferential Treatment of VIP, Better Overall Throughput

Visible QoS Spectrum (Cache/Memory)



QoS Aware Architecture - Cache



24

Cache/Memory QoS Benefits

Client SPEC Case Study



Response Time -- Lower is better

Based on Measurements, Simulations and Analytical Projections

Significant Benefits of Cache/Memory QoS

QoS Aware Architecture: Power





Improves performance of the user-preferred application

Virtualization: From VMs to VPAs

(Managing Transparent Resources)



Summary

Large-scale CMP is going to happen

-Lots of work to be done to identify and remove platform and architectural limitations preventing applications and execution environments from scaling up to 100's of logical processors on a single CPU die

•Scalability concerns can be addressed

-Hierarchy of Shared Caches

-Large DRAM caches

QoS concerns can be addressed

-Dynamic Cache Allocation (Cache QoS)

-Dynamic Power Management (Power QoS)

•Smart performance management requires more visibility be available to the execution environment

-Resource utilization counters for schedulers, etc.