

**NTNU**  
Norwegian University of Science and Technology

**PhD Defense Presentation**

**Managing Shared Resources in Chip Multiprocessor Memory Systems**


12. October 2010

Magnus Jahre

www.ntnu.no

**Outline**

- Chip Multiprocessors (CMPs)
- CMP Resource Management
- Miss Bandwidth Management
  - Greedy Miss Bandwidth Management
  - Interference Measurement
  - Model-Based Miss Bandwidth Management
- Off-Chip Bandwidth Management
- Conclusion



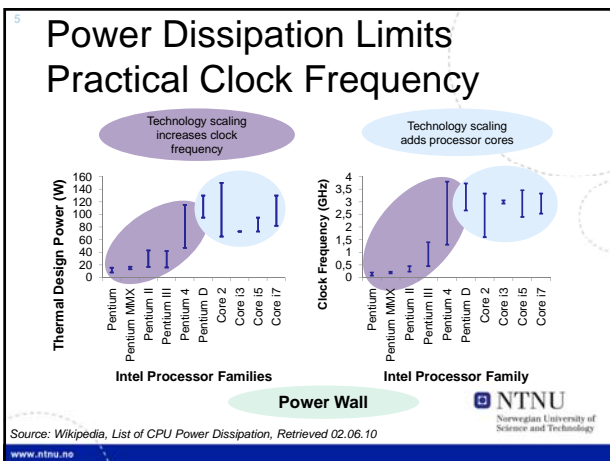
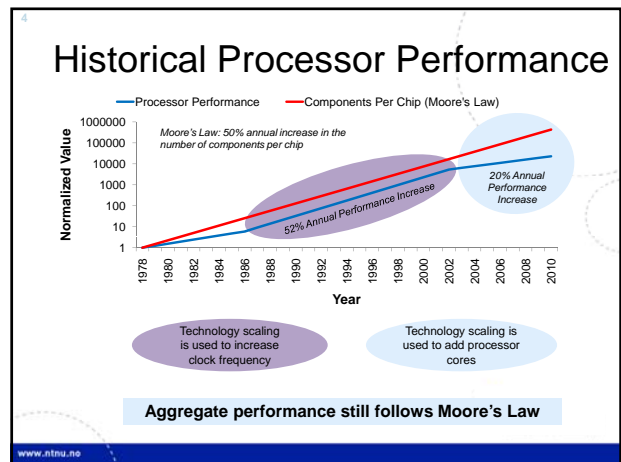
NTNU  
Norwegian University of Science and Technology

www.ntnu.no

**CHIP MULTIPROCESSORS**

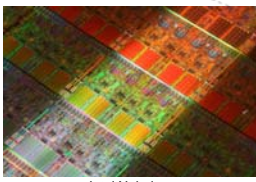
NTNU  
Norwegian University of Science and Technology

www.ntnu.no



**Chip Multiprocessors (CMPs)**

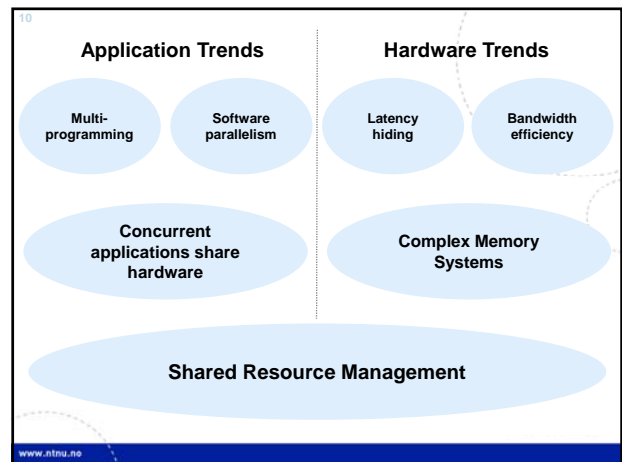
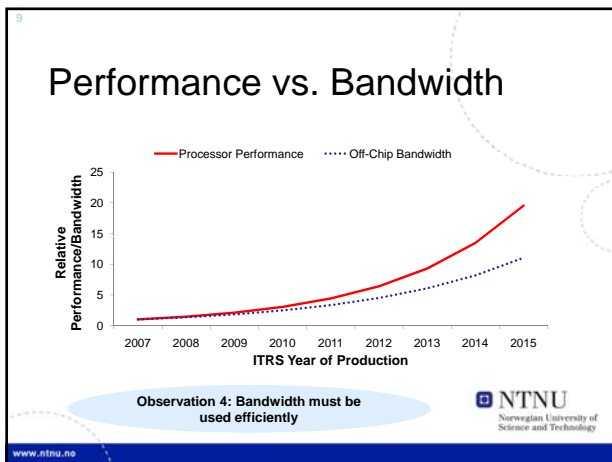
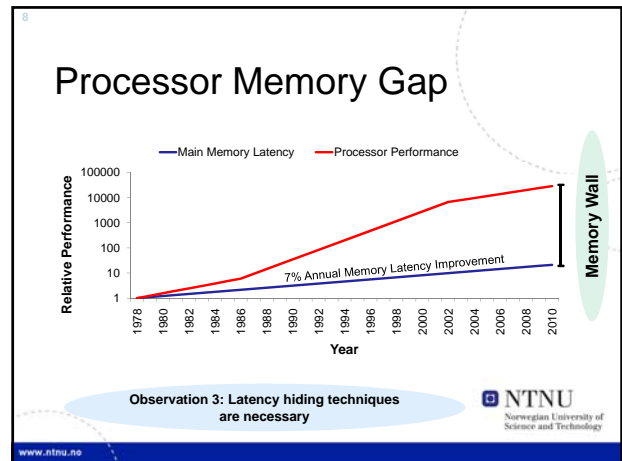
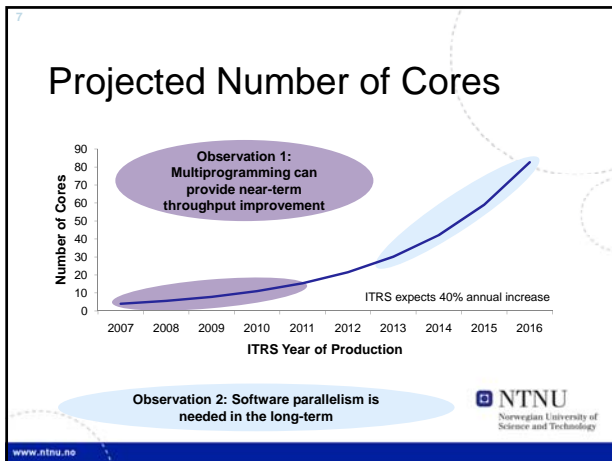
- CMPs utilize chip resources with a constant power budget
- How does technology scaling impact CMPs?



Intel Nehalem

NTNU  
Norwegian University of Science and Technology

www.ntnu.no



### 11 CMP RESOURCE MANAGEMENT

NTNU Norwegian University of Science and Technology

www.ntnu.no

- ### 12 Why Manage Shared Resources?
- Provide predictable performance
  - Support OS scheduler assumptions
  - Cloud: Fulfill Service Level Agreement
- NTNU Norwegian University of Science and Technology
- www.ntnu.no

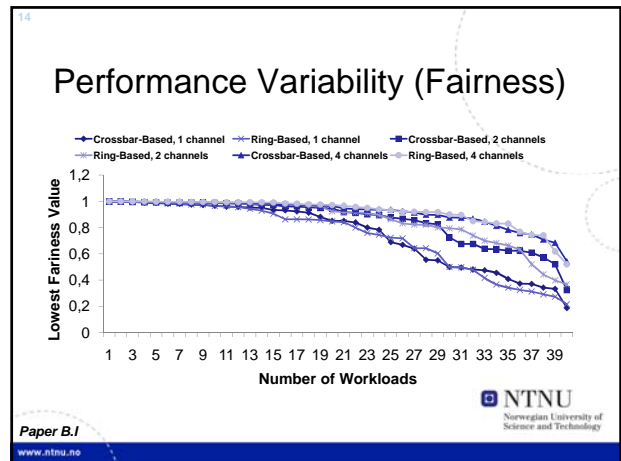
13

## Performance Variability Metrics

- Fairness
  - The performance reduction due to interference between processes is distributed across all processes in proportion to their priorities
  - *Equal priorities*: Performance reduction from sharing affects all processes equally
- Quality of Service
  - The performance of a process is never drops below a certain limit regardless of the behavior of co-scheduled processes

NTNU  
Norwegian University of Science and Technology

www.ntnu.no



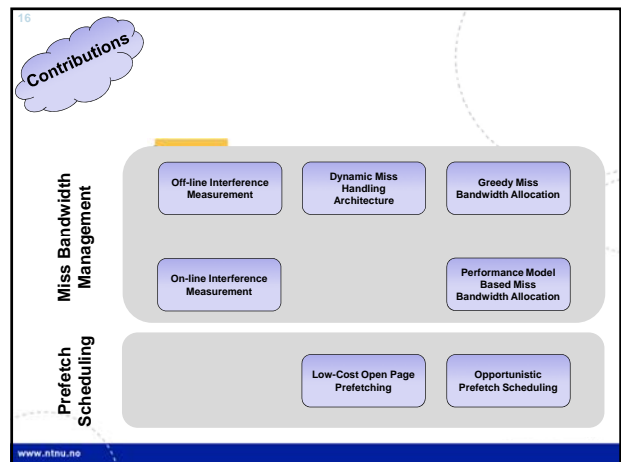
15

## Resource Management Tasks

- Measurement
- Allocation (Policy)
- Enforcement (Mechanism)

NTNU  
Norwegian University of Science and Technology

www.ntnu.no



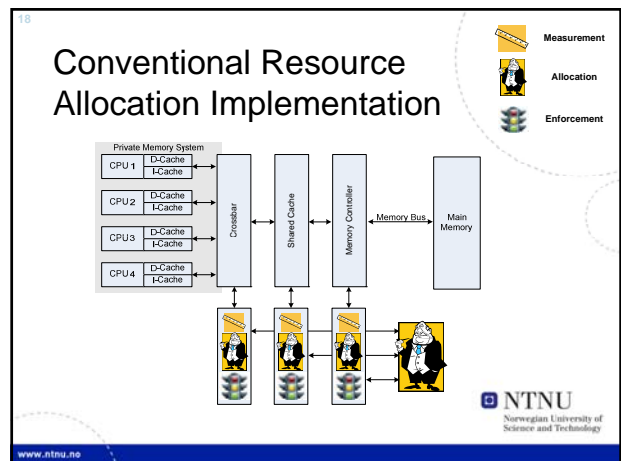
17

Miss Bandwidth Management

## GREEDY MISS BANDWIDTH MANAGEMENT

NTNU  
Norwegian University of Science and Technology

www.ntnu.no



19

## Alternative Resource Allocation Implementation

Measurement  
Allocation  
Enforcement

Private Memory System

CPU1 D-Cache I-Cache  
CPU2 D-Cache I-Cache  
CPU3 D-Cache I-Cache  
CPU4 D-Cache I-Cache

Crossbar

Shared Cache

Memory Controller

Memory Bus

Main Memory

Dynamic Miss Handling Architecture

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

20

## Dynamic Miss Handling Architecture

Accesses

Cache

Miss Handling Architecture (MHA)

D		1
E		1
B		1
		0
		0
Address	Target Info.	U

A DMHA controls the number of concurrent shared memory system requests that are allowed for each processor

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

21

## Greedy Miss Bandwidth Management

- Idea: Reduce the number of MSHRs if a metric exceeds a certain threshold
- Metrics:
  - Paper A.II: Memory bus utilization
  - Paper A.III: Simple interference counters (Interference Points)
- Performance feedback avoids excessive performance degradations

Paper A.II and A.III

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

22

## Miss Bandwidth Management

# INTERFERENCE MEASUREMENT

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

23

## Resource Allocation Baselines

Baseline = Interference-free configuration

Quantify performance impact from interference

Private Mode and Shared Mode

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

24

## Interference Definition

$$I_i = L_i - \mathcal{L}_i$$

Interference

Shared Mode Latency

Private Mode Latency

$$E_i = \hat{\mathcal{L}}_i - \mathcal{L}_i$$

Estimate Error

Private Mode Latency Estimate

Private Mode Latency Measurement

NTNU  
Norwegian University of Science and Technology

www.ntnu.no

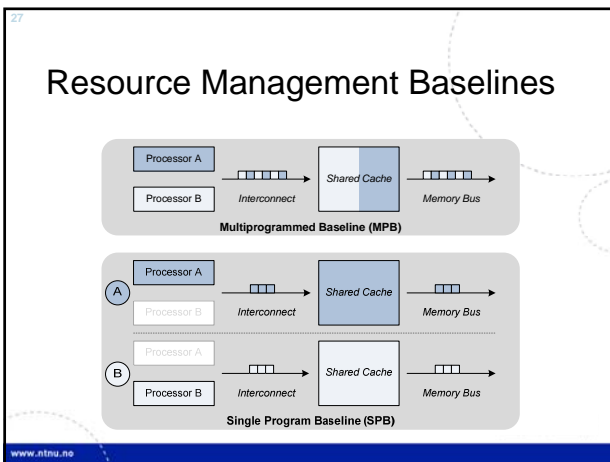
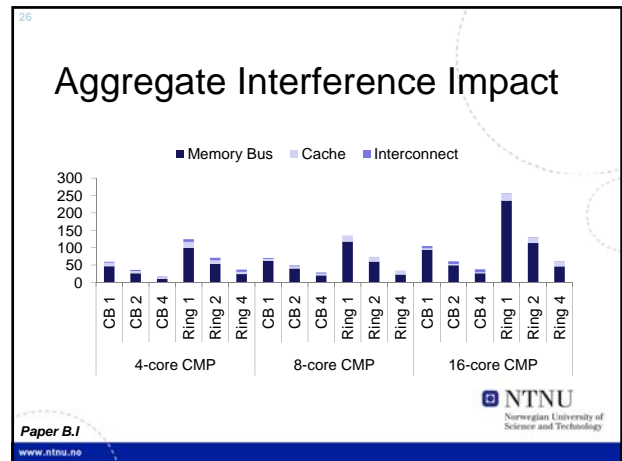
## Offline Interference Measurement

**Interference Penalty Frequency (IPF)** counts the number requests that experienced an interference latency of  $i$  cycles

**Interference Impact Factor (IIF)** is the interference latency times the probability of it arising, i.e.  $IIF(i) = i \cdot P(i)$

**Paper B.I**  
www.ntnu.no

NTNU  
Norwegian University of Science and Technology



## Baseline Weaknesses

- Multiprogrammed Baseline
  - Only accounts for interference of shared resources
  - Static and equal division of DRAM bandwidth does not give equal latency
  - Complex relationship between resource allocation and performance
- Single Program Baseline
  - ~~Does not exist in shared mode~~

**Online Interference Measurement  
Dynamic Interference Estimation  
Framework (DIEF)**

**Paper B.II**  
www.ntnu.no

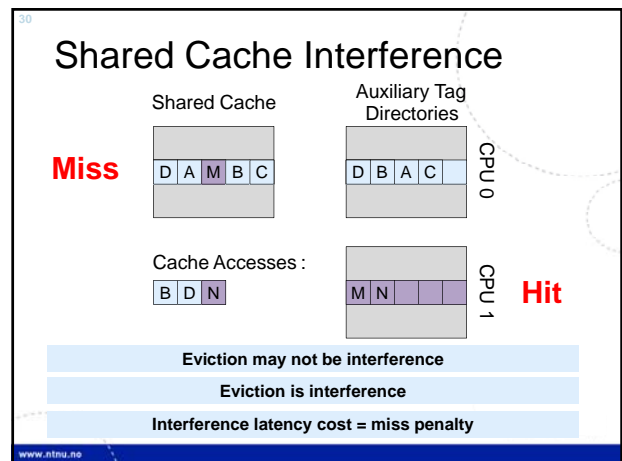
NTNU  
Norwegian University of Science and Technology

## Online Interference Measurement

- Dynamic Interference Estimation Framework (DIEF)
- Estimates private mode average memory latency
- General, component-based framework

**Paper B.II**  
www.ntnu.no

NTNU  
Norwegian University of Science and Technology



### Bus Interference Requirements

- Out-of-order memory bus scheduling
- Shared mode only cache misses and cache hits
- Shared cache writebacks

Computing private latency based on shared mode queue contents is difficult

Emulate private scheduling in the shared mode

NTNU  
Norwegian University of Science and Technology

### Shared Bus Queue

E	D	C	B
---	---	---	---

Memory Latency Estimation Buffer

D 40	C 40	B 200	A 120
---------	---------	----------	----------

Open Page Emulation Registers

15	Bank 0
32	Bank 1
⋮	⋮

⋯→ Arrival Order  
 - - -> Head Pointer  
 → Execution Order

Bank/Page Mapping: A → (0,15), B → (0,19), C → (0,15), D → (1,32)

Latency Lookup Table

B	200
---	-----

**Estimated Queue Latency = 120 + 40 + 40**

NTNU  
Norwegian University of Science and Technology

### Miss Bandwidth Management

## MODEL-BASED MISS BANDWIDTH MANAGEMENT

NTNU  
Norwegian University of Science and Technology

### Model-Based Miss Bandwidth Allocation

DIEF provides accurate estimates of the average private mode memory latency

Can we use the estimates provided by DIEF to choose miss bandwidth allocations?

We need a model that relates **average memory latency** to **performance**

NTNU  
Norwegian University of Science and Technology

### Performance Model

*Observation: The memory latency performance impact depends on the parallelism of memory requests*

$$IPC_p = \frac{N_p}{C_p^{Compute} + C_p^{MemStall}}$$

$$C_p^{MemStall} = (M_p) \cdot L_p$$

Very similar in private and shared mode

Shared mode measurements can provide private mode performance estimates

NTNU  
Norwegian University of Science and Technology

### Bandwidth Management Flow

**Measurement**

- Shared Mode Memory Latency
- Private Mode Memory Latency
- CPU Stall Time
- Committed Instructions
- Number of Memory Requests

**Modeling**

Per-CPU Models

Perf. Metric Model

**Allocation**

Find MSHR allocation that maximizes the chosen performance metric

Set number of MSHRs for all last-level private caches

NTNU  
Norwegian University of Science and Technology

37

# OFF-CHIP BANDWIDTH MANAGEMENT



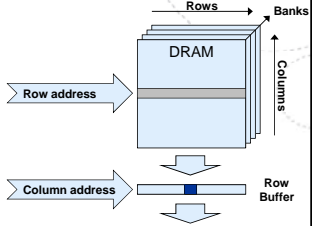

Norwegian University of Science and Technology

www.ntnu.no

38

# Modern DRAM Interfaces

- Maximize bandwidth with 3D organization
- Repeated requests to the row buffer are very efficient

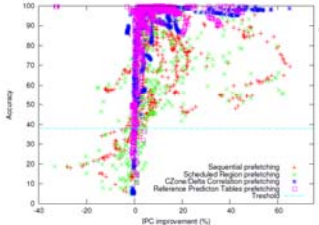

Norwegian University of Science and Technology

www.ntnu.no

39

# Low-Cost Open Page Prefetching

- Idea: Piggyback prefetches to open DRAM pages on demand reads
- Performance win if prefetcher accuracy is above ~40%

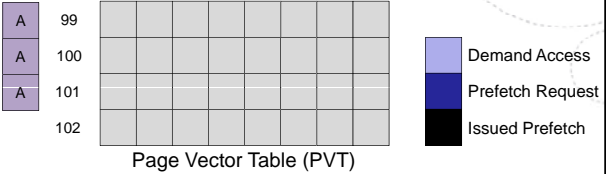
Norwegian University of Science and Technology

Paper C.I

www.ntnu.no


40

# Opportunistic Prefetch Scheduling



Page Vector Table (PVT)

Idea: Issue prefetches when a page is closed  
Increased efficiency: 8 transfers for 3 activations




Norwegian University of Science and Technology

Paper C.II

www.ntnu.no

41

# CONCLUSION




Norwegian University of Science and Technology

www.ntnu.no

42

# Conclusion

- Managing bandwidth allocations can improve CMP system performance
- Miss bandwidth management
  - Greedy allocations
  - Management guided by accurate measurements and performance models
- Off-chip bandwidth management with prefetching



Norwegian University of Science and Technology

www.ntnu.no

43

## Thank You



Visit our website:  
<http://research.idi.ntnu.no/multicore/>

**NTNU**  
Norwegian University of  
Science and Technology

www.ntnu.no

44

## EXTRA SLIDES

**NTNU**  
Norwegian University of  
Science and Technology

www.ntnu.no

45

## Future Work

- Performance-directed management of shared caches and the memory bus
- Improving OS and system software with dynamic measurements
- Combining dynamic MHAs with prefetching to improve system performance
- Managing workloads of single-threaded and multi-threaded benchmarks

**NTNU**  
Norwegian University of  
Science and Technology

www.ntnu.no

