

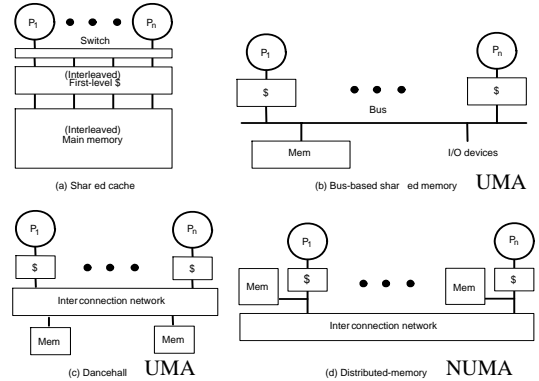


Sun's E6000 Server Family

Erik Hagersten
Uppsala University
Sweden



What Approach to Shared Memory

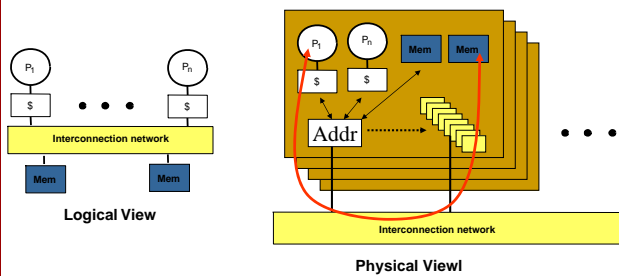


DARK 2008

Dept of Information Technology | www.it.uu.se © Erik Hagersten | user.it.uu.se/~eh



Looks like a NUMA but drives like a UMA



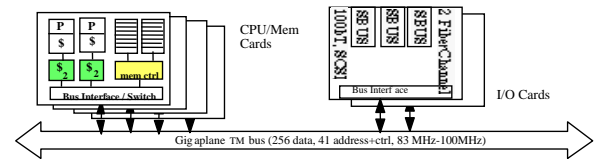
- Memory bandwidth scales with the processor count
- One "interconnect load" per (2xCPU + 2xMem)
- Optimize for the dancehall case (no memory shortcut)

DARK 2008

Dept of Information Technology | www.it.uu.se © Erik Hagersten | user.it.uu.se/~eh



SUN Enterprise Overview



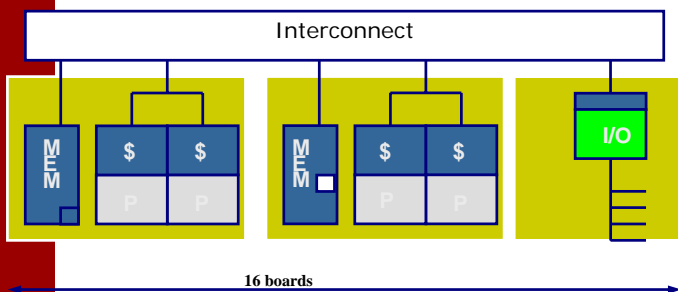
- 16 slots with either CPUs or IO
- Up to 30 UltraSPARC processors (peak 9 GFLOPs)
- Gigaplane™ bus has peak bw 2.67 GB/s; up to 30GB memory
- 16 bus slots, for processing or I/O boards

DARK 2008

Dept of Information Technology | www.it.uu.se © Erik Hagersten | user.it.uu.se/~eh



Enterprise Server E6000

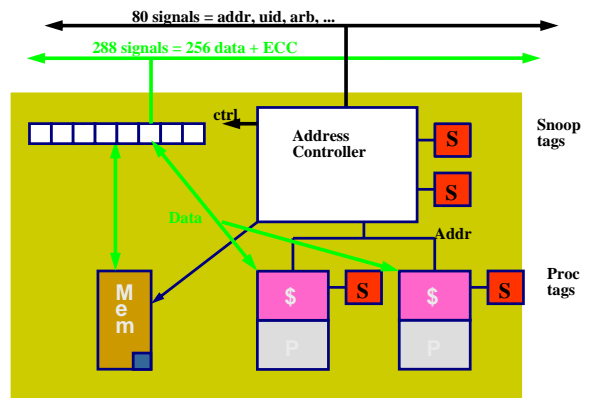


DARK 2008

Dept of Information Technology | www.it.uu.se © Erik Hagersten | user.it.uu.se/~eh



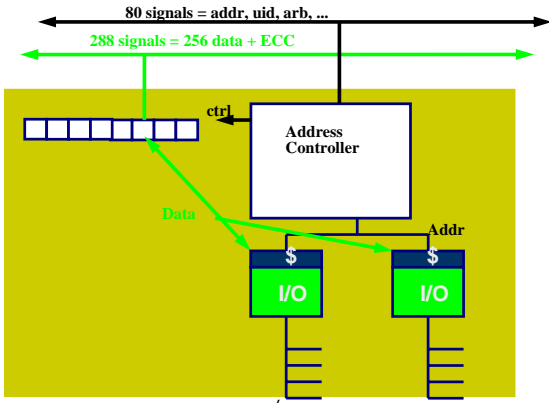
An E6000 Proc Board



DARK 2008

Dept of Information Technology | www.it.uu.se © Erik Hagersten | user.it.uu.se/~eh

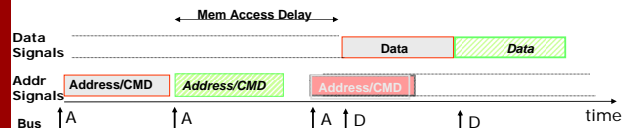
An I/O Board



Dept of Information Technology | www.it.uu.se | © Erik Hagersten | user.it.uu.se/~eh

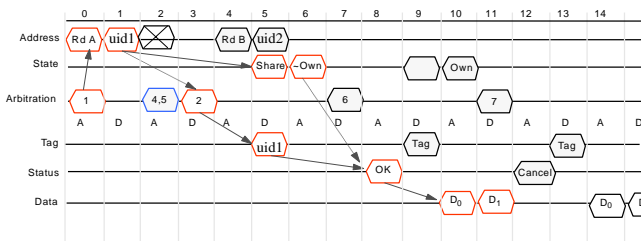
Split-Transaction Bus

- Split bus transaction into request and response sub-transactions
 - Separate arbitration for each phase
- Other transactions may intervene
 - Improves bandwidth dramatically
 - Response is matched to request
 - Buffering between bus and cache controllers



Dept of Information Technology | www.it.uu.se | 8 | © Erik Hagersten | user.it.uu.se/~eh

Gigaplane Bus Timing



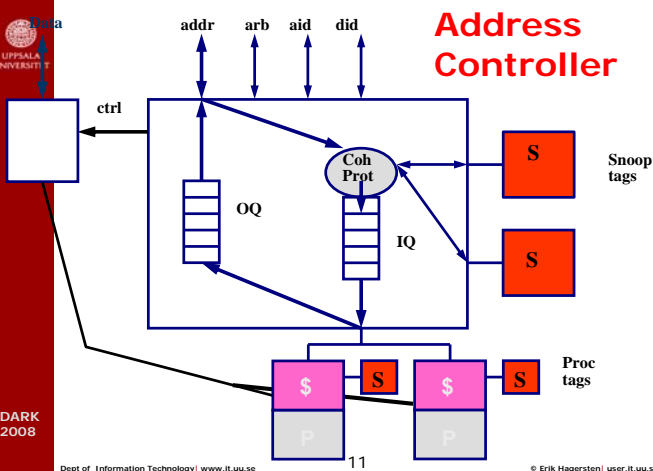
Dept of Information Technology | www.it.uu.se | 9 | © Erik Hagersten | user.it.uu.se/~eh

Electrical Characteristics of the Bus

- At most 16 electrical loads per signal
- 8 boards from each side (ex. 15 CPU + 1 I/O)
- 20.5 inches "centerplane"
- Well controlled impedance
- ~350-400 signals
- Runs at 90/100 MHz

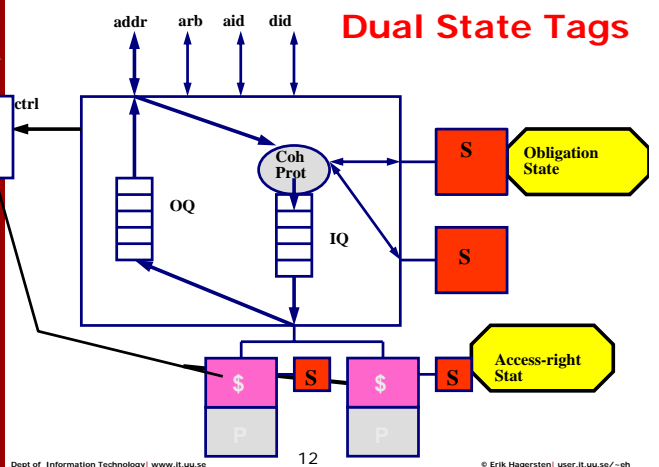
Dept of Information Technology | www.it.uu.se | 10 | © Erik Hagersten | user.it.uu.se/~eh

Address Controller



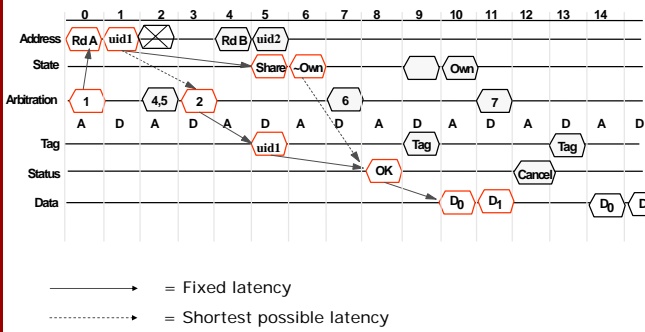
Dept of Information Technology | www.it.uu.se | 11 | © Erik Hagersten | user.it.uu.se/~eh

Dual State Tags

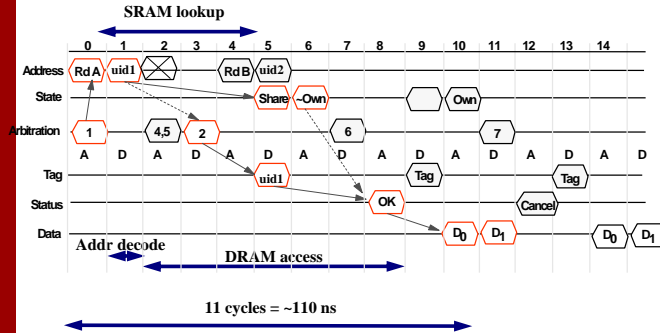


Dept of Information Technology | www.it.uu.se | 12 | © Erik Hagersten | user.it.uu.se/~eh

Timing of a single read trans Board 1 reading from mem 2



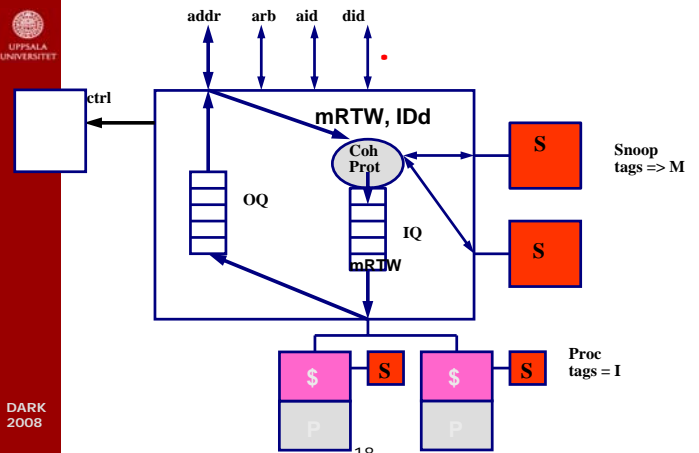
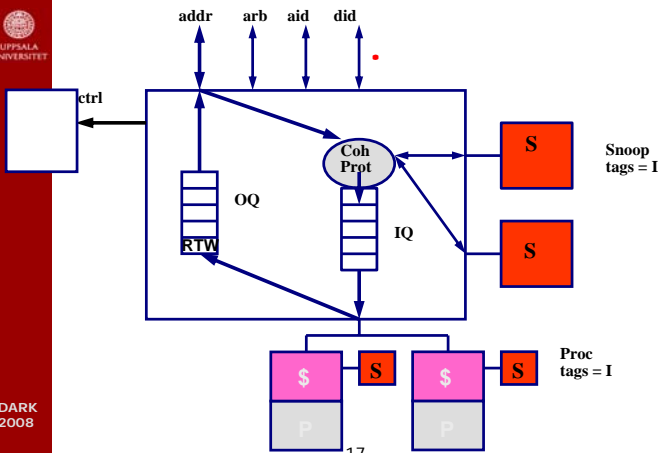
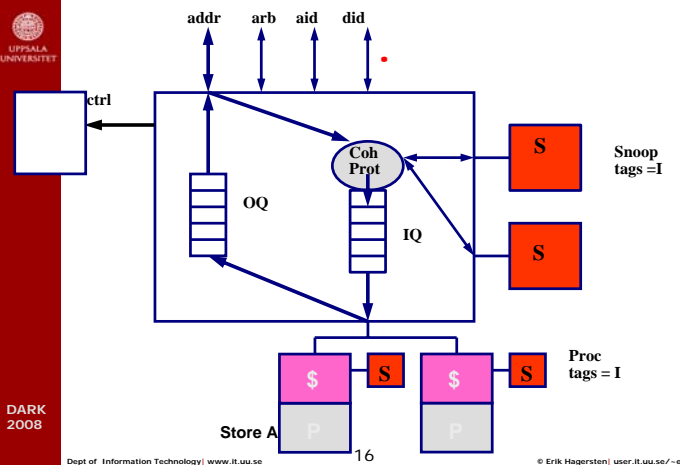
Protocol tuned for timing

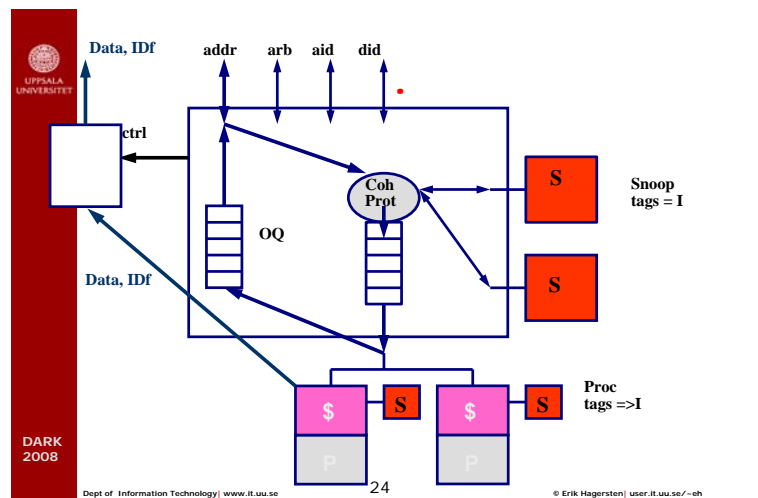
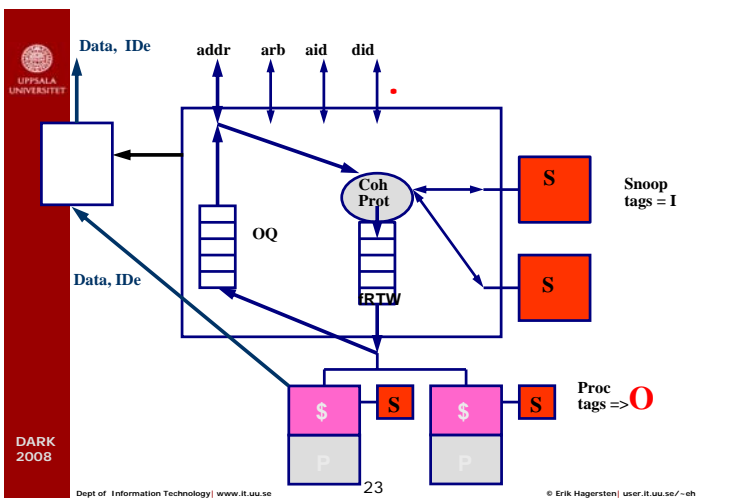
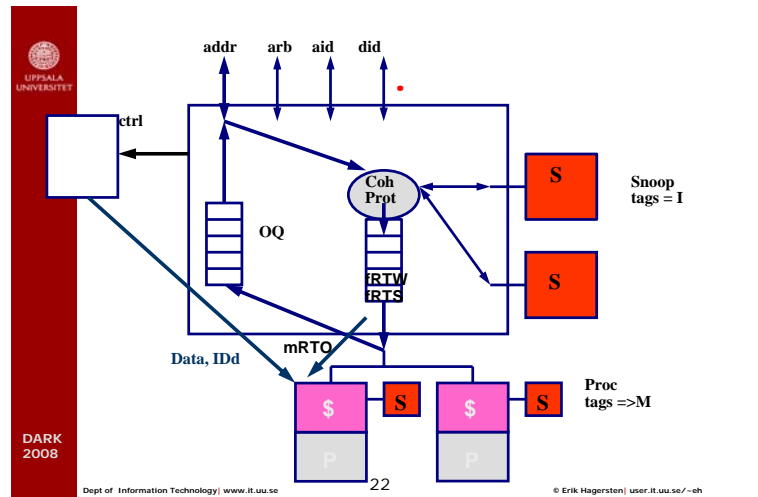
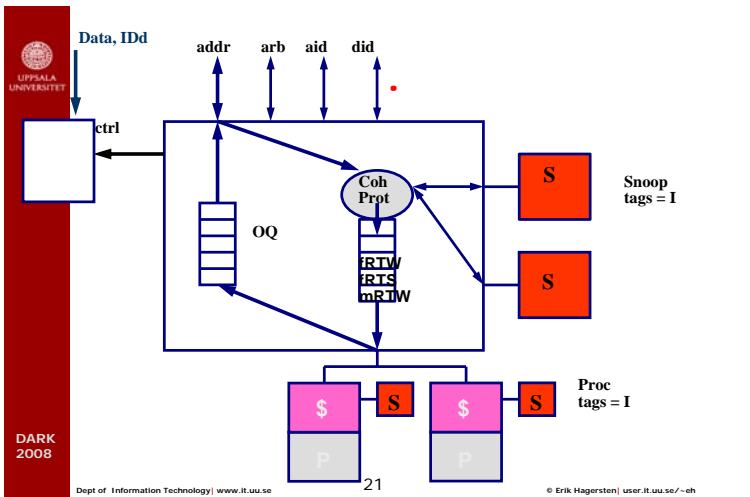
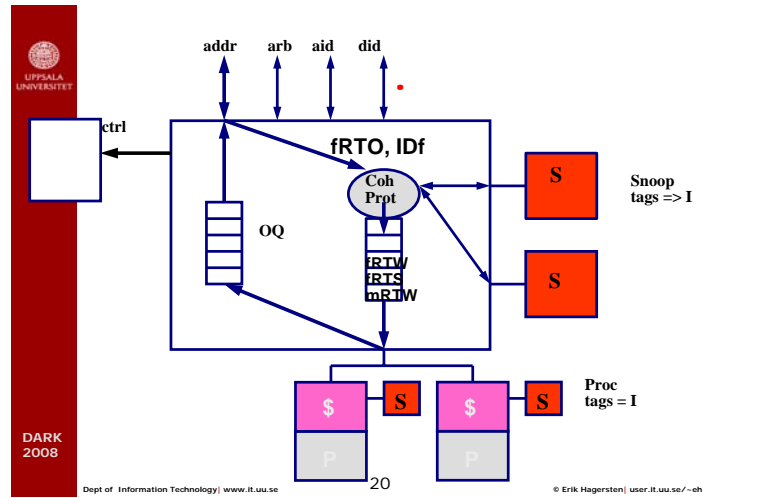
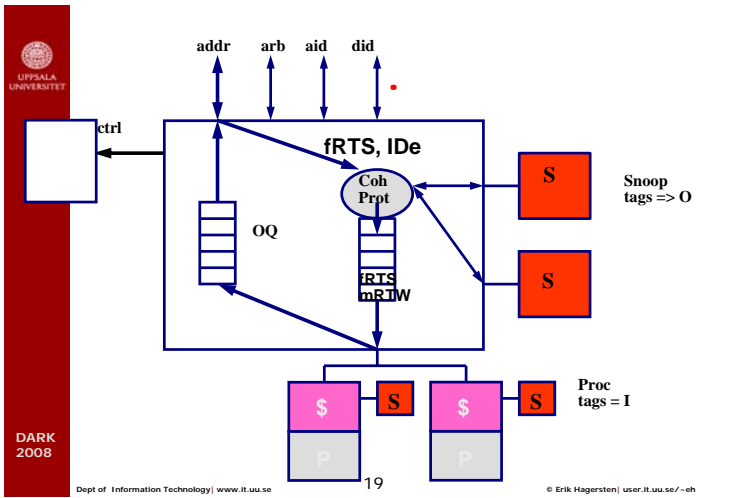


Foreign and own transactions queue in IQ State Change on Address Packet

- Data "A" initially resides in CPU7's cache
- CPU1: Issues a store request to "A"
- CPU1: Read-To-Write req, ID=d, (i.e., "write request")
- CPU13: LD "A" -> Read-To-Shared req, ID=e
- CPU15: ST "A" -> RTW req, ID=f

mRTO stored in IQ_{CPU1}
Own read IQtrans retired when data arrives
Later requests for A queued in IQ_{CPU1} behind mRTO
IQ_{CPU1} will eventually store: <mRTW_{Id}, fRTS_{Id}, fRTW_{Id}>





A cascade of "write requests"

- A initially resides in CPU7's cache
- CPU1: RTW, ID=a
- CPU2: RTW, ID=b
- ...
- CPU5: RTW, ID=f

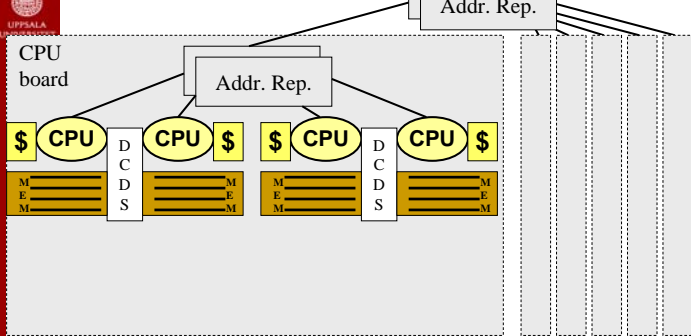
CPU tags		Snoop tags
I	IQ1 = <mRTW_{1Da}, fRTW_{1Db}>	I
I	IQ2 = <mRTW_{2Db}, fRTW_{2Dc}>	I
...
I	IQ5 = <mRTW_{5Df}>	M
...
S	IQ7 = <fRTW_{7Da}>	I



Implementing Sun's SunFire 6800

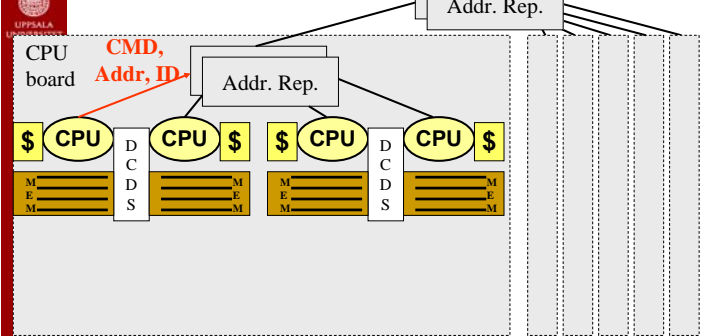
Erik Hagersten
Uppsala University
Sweden

FirePlane, 24 CPUs



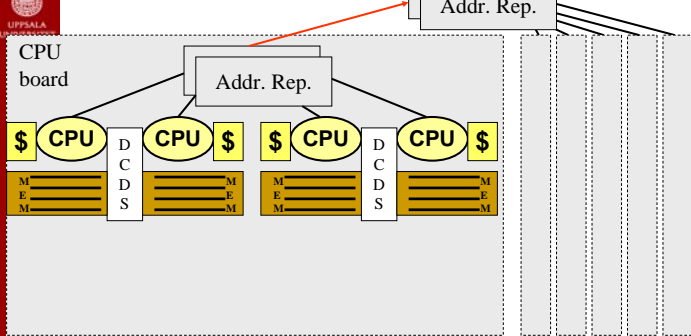
L2 cache = 8MB, snoop tags on-chip
CPU 1+GHz UltraSPARC III
Mem= 4+GB/CPU

FirePlane, 24 CPUs

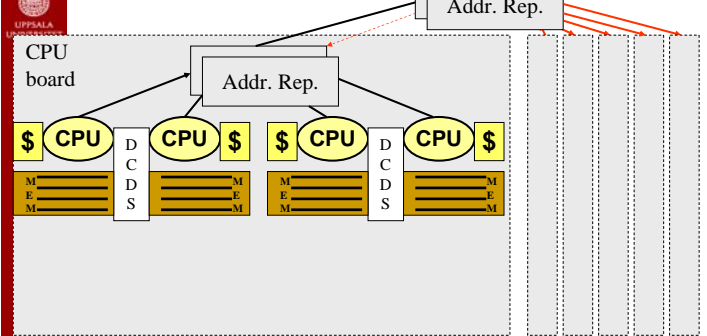


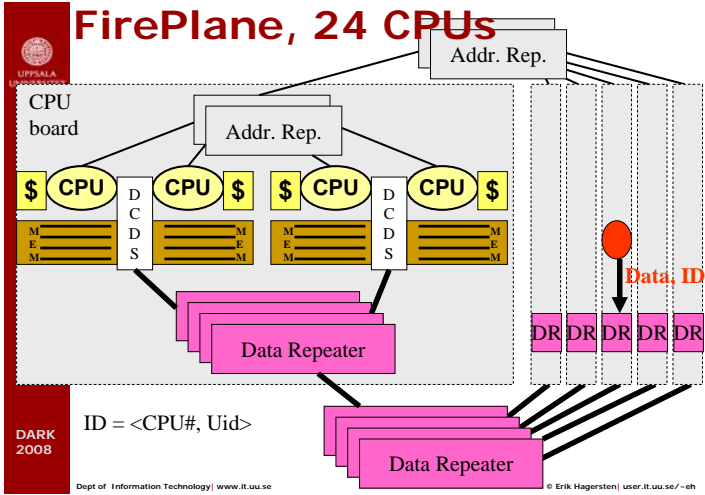
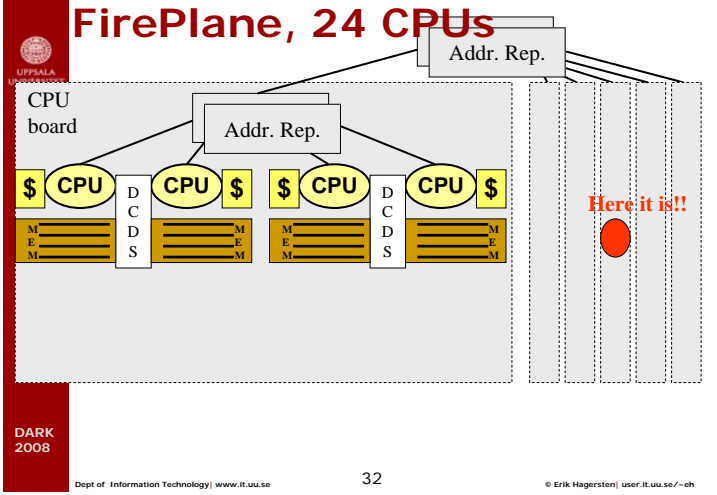
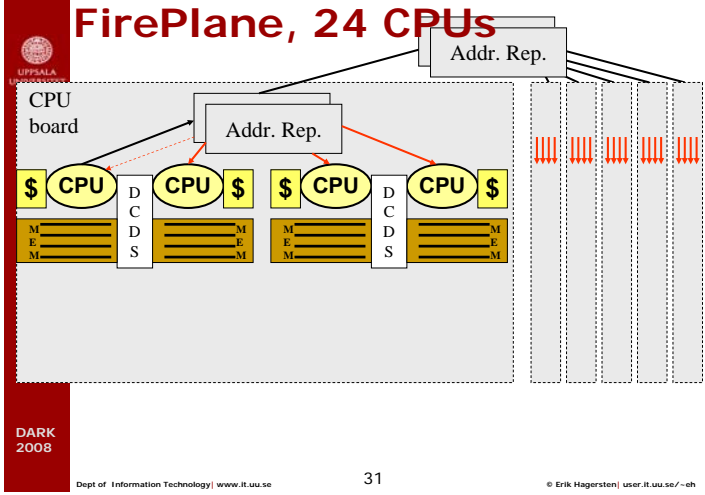
ID = <CPU#, Uid>

FirePlane, 24 CPUs



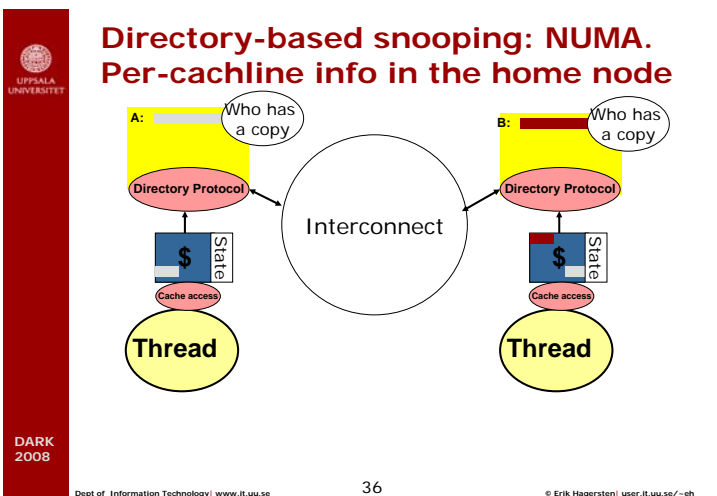
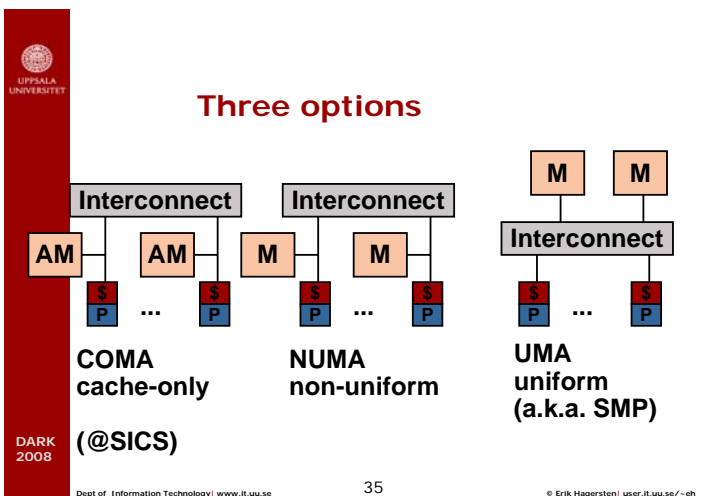
FirePlane, 24 CPUs





Scalable Shared-Memory Implementations

Erik Hagersten
Uppsala University
Sweden



"Upgrade" in dir-based

Who has a copy

Who has a copy

INV ACKV INV

ACK ACK

Thread Thread Thread

Read A ... Read B
Read A ... Read A
...
Write A

DARK 2008

Dept of Information Technology | www.it.uu.se

37

© Erik Hagersten | user.it.uu.se/~eh

Cache-to-cache in dir-based

Who has a copy

Who has a copy

ReadRequest ReadDemand

Ack Forward

Thread Thread Thread

Read A ... Read B
Read A ... Read A
...
Read A

DARK 2008

Dept of Information Technology | www.it.uu.se

38

© Erik Hagersten | user.it.uu.se/~eh

Fully mapped directory

Memory Directory

presence bits dirty bit

- k Nodes
- Each node is the "home" for $1/k$ of the memory
- Dir entry per cacheline in home memory: k presence-bits + 1 dirty-bit
- Requests are first sent to the home node's CA

DARK 2008

Dept of Information Technology | www.it.uu.se

39

© Erik Hagersten | user.it.uu.se/~eh

Reducing the Memory Overhead: SCI

--- Scalable Coherence Interface (SCI)

- home only holds pointer to rest of the directory info [$\log(N)$ bits]
- distributed linked list of copies, weaves through caches
 - cache tag has pointer, points to next cache with a copy
- on read, add yourself to head of the list (comm. needed)
- on write, propagate chain of invalidations down the list
- on replacement: remove yourself from the list

Main Memory (Home)

Node 0 Node 1 Node 2

Cache Cache Cache

DARK 2008

Dept of Information Technology | www.it.uu.se

40

© Erik Hagersten | user.it.uu.se/~eh

Cache Invalidation Patterns

Barnes-Hut Invalidation Patterns

Radiosity Invalidation Patterns

of invalidations

DARK 2008

Dept of Information Technology | www.it.uu.se

41

© Erik Hagersten | user.it.uu.se/~eh

Overflow Schemes for Limited Pointers

- Broadcast (Dir_B)
 - broadcast bit turned on upon overflow
 - bad for widely-shared invalidated data
- No-broadcast (Dir_{NB})
 - on overflow, new sharer replaces one of the old ones (invalidated)
 - bad for widely read data
- Coarse vector (Dir_{CV})
 - change representation to a coarse vector, 1 bit per k nodes
 - on a write, invalidate all nodes that a bit corresponds to

Mem: Overflow bit 2 Pointers

(a) No overflow

Overflow bit 8-bit coarse vector

(b) Overflow bit

(c) Coarse vector

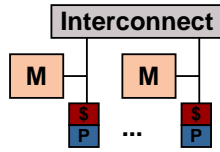
DARK 2008

Dept of Information Technology | www.it.uu.se

42

© Erik Hagersten | user.it.uu.se/~eh

cc-NUMA issues



- Memory placement is key!
- Gotta' migrate data to where it's being used
- Gotta' have cache affinity
 - Long time between process switches in the OS
 - Reschedule processor on the CPU it ran last
- Origin 2000's migration always turned off☹



Sun's WildFire System

Erik Hagersten
Uppsala University
Sweden

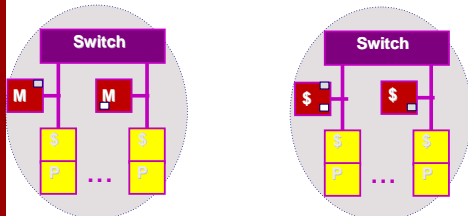
Sun's WildFire System

- Runs unmodified SMP apps in a more scalable way than E6000
- Minor modifications to E6000 snooping required
- CPUs generate local address OR global address
- Global address --> no replication (NUMA)
- Coherent Memory Replication (~Simple COMA@ SICS)
- Hardware support for detecting migration/replication pages
- Directory cache + address translation cache backed by memory
- Deterministic directory implementation (easy to verify)

WildFire: One Solaris spanning four nodes



COMA: self-optimizing DSM



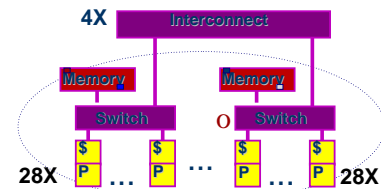
ccNUMA

COMA

COMA:

- Self-optimizing architecture
- Problem at high memory pressure
- Complex hardware and coherence protocol

Adaptive S-COMA of Large SMPs



- A page may have space allocated in many nodes
- HW maintains memory coherence per cache line
- Replication under SW control --> simple HW (S-COMA)
- Adaptive replication algorithm in OS (R-NUMA)
- Coherent Memory Replication (CMR)
- Hierarchical affinity scheduler (HAS)
- Few large nodes --> simple interconnect and coherence protocol

A WildFire Node

- 16 slots with either CPUs, IO or... WildFire extension board →
 - Up to **28** UltraSPARC processors
 - Gigaplane™ bus has peak bw 2.67 GB/s
 - Local access time of 330ns (Imbench)

DARK 2008 Dept of Information Technology | www.it.uu.se 49 © Erik Hagersten | user.it.uu.se/~eh

Sun WildFire Interface Board

DARK 2008 Dept of Information Technology | www.it.uu.se 50 © Erik Hagersten | user.it.uu.se/~eh

Sun WildFire Interface Board

DARK 2008 Dept of Information Technology | www.it.uu.se 51 © Erik Hagersten | user.it.uu.se/~eh

WildFire as a vanilla "NUMA"

DARK 2008 Dept of Information Technology | www.it.uu.se 52 © Erik Hagersten | user.it.uu.se/~eh

NUMA -- local memory access

DARK 2008 Dept of Information Technology | www.it.uu.se 53 © Erik Hagersten | user.it.uu.se/~eh

NUMA -- remote memory access

DARK 2008 SRAM overhead = $10/512 = 2\%$ (lower bound $2/512 = 0.4\%$)
 Dept of Information Technology | www.it.uu.se 54 © Erik Hagersten | user.it.uu.se/~eh

Global Cache Coherence Prot.

Mod dir entry

Access right changes

Mem, I/F, Dir\$, Mtag, Cache, Proc

DARK 2008

Dept of Information Technology | www.it.uu.se

55

© Erik Hagersten | user.it.uu.se/~eh

NUMA -- local memory access

Interconnect

Access right OK? NO!!

Mem, I/F, Dir\$, Mtag, Cache, Proc

DARK 2008

Dept of Information Technology | www.it.uu.se

56

© Erik Hagersten | user.it.uu.se/~eh

Gigaplane Bus Timing

Address	Rd A	uidA	XXX	XXX	Rd B	uidB	Rd C	uidC						
State	A	D	A	D	A	D	A	D	A	D	A	D	A	D
Arbitration	1	4.5	2				6				7			
Tag						uidA					uidB			
Status						OK					Cancel			
Data						D ₀	D ₁							

DARK 2008

Dept of Information Technology | www.it.uu.se

57

© Erik Hagersten | user.it.uu.se/~eh

WildFire Bus Extensions

Address	Rd A	uidA	XXX	XXX	Rd B	uidB						Rd A	uidA	
State	A	D	A	D	A	D	Ignore					A	D	A
Arbitration	1	4.5	2									7	9	
Tag						uidA						uidB		
Status						OK						Cancel		
Data						D ₀	D ₁							

Asserted by WildFire, Resent by WildFire

Ignore transaction squashes an ongoing transaction => not put in IQ

WildFire eventually reissues the same transaction

RTSF -- a new transaction sends data to CPU and memory

DARK 2008

Dept of Information Technology | www.it.uu.se

58

© Erik Hagersten | user.it.uu.se/~eh

WildFire Directory -- only 4 nodes!!

WildFire Directory

- k nodes (with one or more procs).
- With each cache-block in memory: k presence-bits, 1 dirty-bit
- With each cache-block in cache: 1 valid bit, and 1 dirty (owner) bit

• ReadRequest from main memory by processor i:

- If dirty-bit OFF then { read from main memory; turn p[i] ON }
- if dirty-bit ON then { recall line from dirty proc (cache state to shared); update memory; turn dirty-bit OFF; turn p[i] ON; supply recalled data to i; }

.....

DARK 2008

Dept of Information Technology | www.it.uu.se

59

© Erik Hagersten | user.it.uu.se/~eh

NUMA "detecting excess misses"

I thought you had the data!!*?@

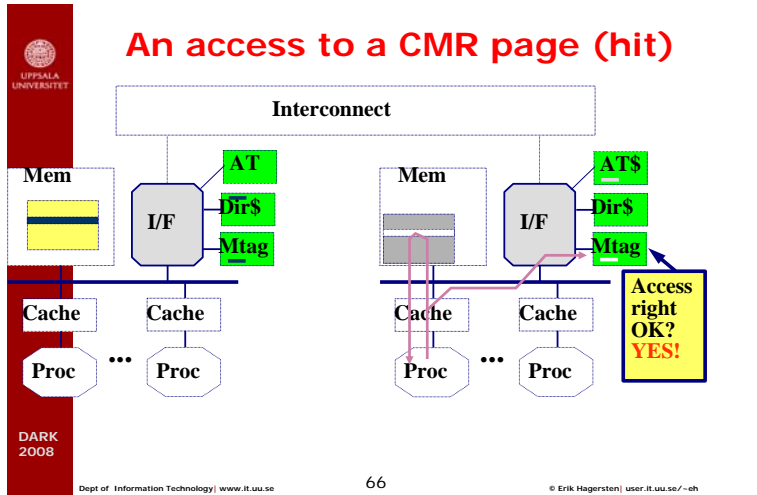
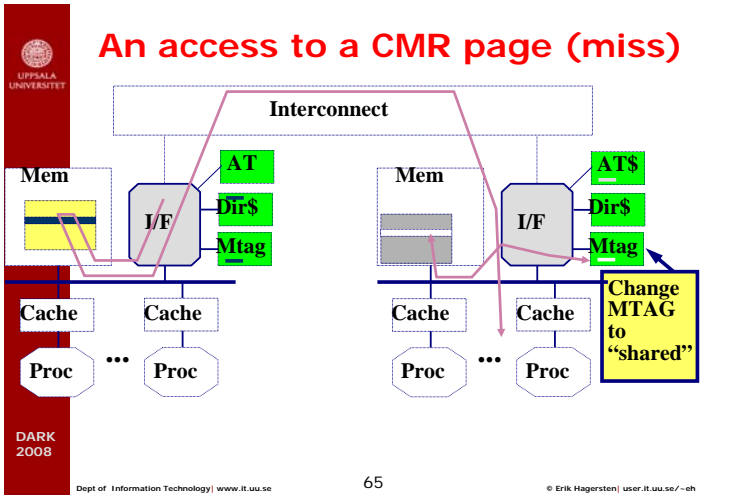
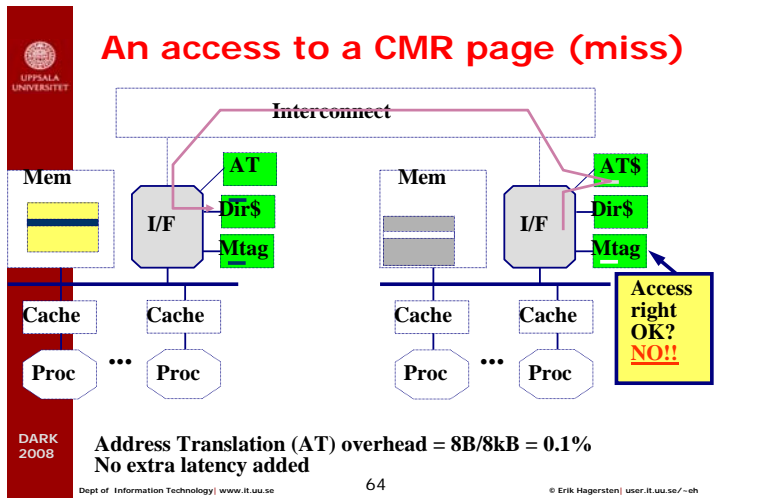
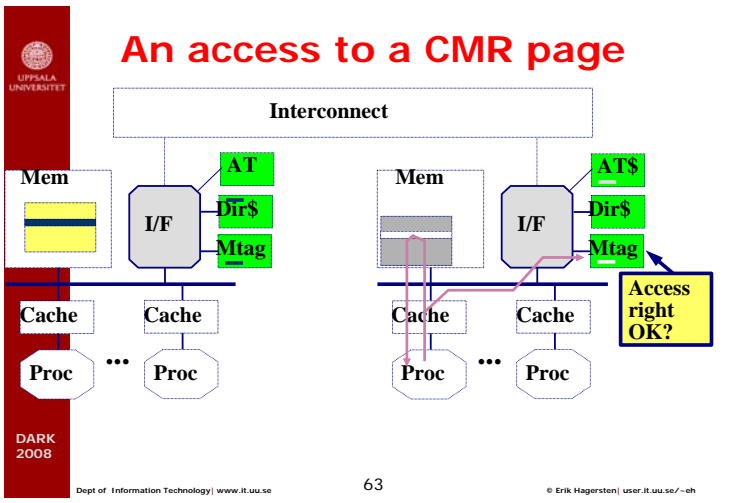
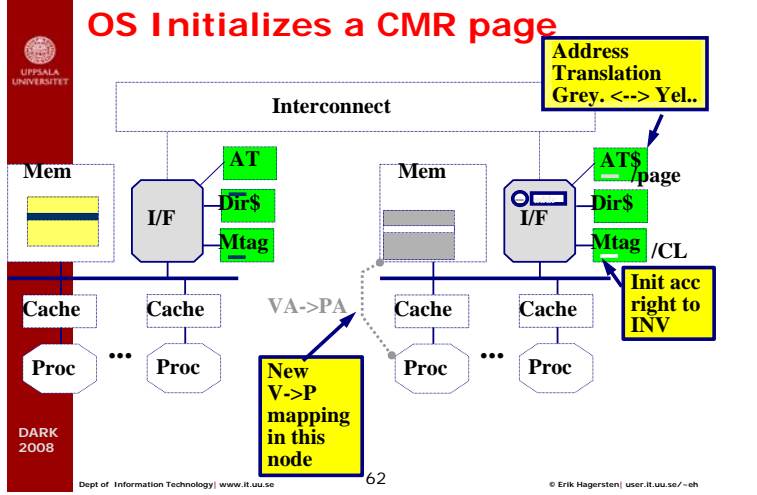
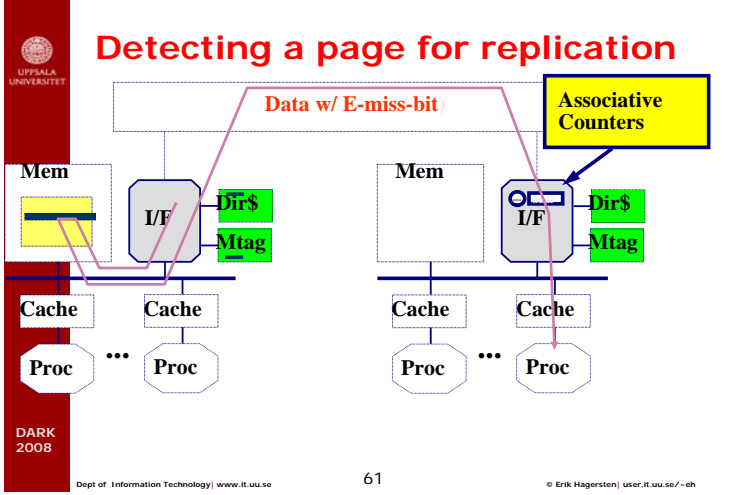
Mem, I/F, Dir\$, Mtag, Cache, Proc

DARK 2008

Dept of Information Technology | www.it.uu.se

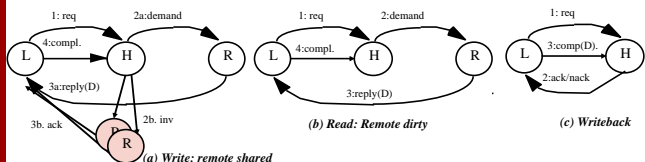
60

© Erik Hagersten | user.it.uu.se/~eh



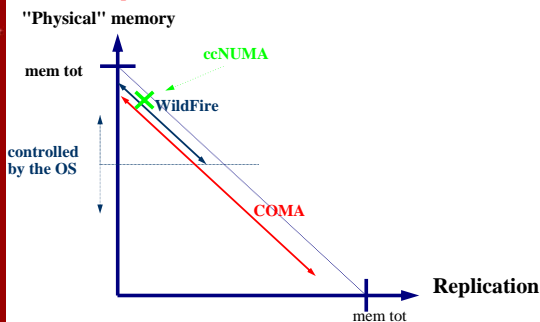
Deterministic Directory

- MOSI protocol, fully mapped directory (one bit/node)
- Directory blocking: one outstanding trans/cache line
- Directory blocks new requests until completion received
- The directory state and cache state always in agreement (except for silent replacement...)



67

Replication Issues Revisited



- Only "promising" pages are replicated
- OS dynamically limits the amount of replication
- Solaris CMR changes in the hat_layer (=port)

68

Advantages of Multiprocessor Nodes

Pros:

- amortization of fixed node costs over multiple processors
- can use commodity SMPs
- fewer nodes to keep track of in the directory
- much communication may stay within node (NUCA)
- can share "node caches" (WildFire: Coherent Memory Replication)

Cons:

- bandwidth shared among processors and interface
- bus may increase latency to local memory
- snoopy bus at remote node increases delays there too

69

Memory cost of replication

Example: Replicate 10% of data in all nodes

- 50 nodes, each with 2 CPUs
 ==> 490% overhead

- 4 nodes, each with 25 CPUs
 ==> 30% overhead

70

Does migration/replication help? NAS parallel Benchmark Study (Execution time in seconds)

[M. Bull, EPCC 2002]

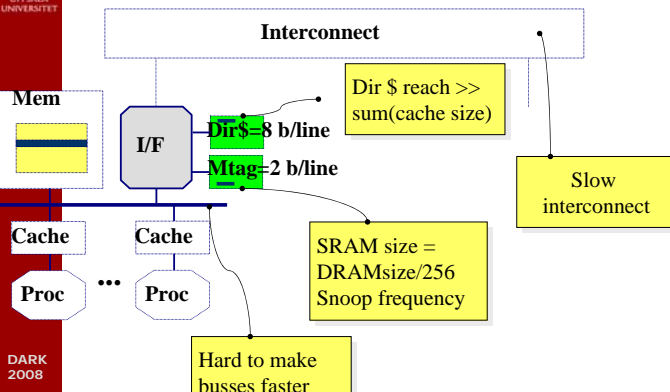
Shallow				BT			
No Initial Plac.		Initial Placement		No Initial Plac.		Initial Placement	
No migr	Migr	No Migr	Migr	No migr	Migr	No Migr	Migr
No Repl	26	5.9	6.1	960	610	620	600
Repl	7.2	6.2	6.1	590	580	580	580

Unopt. FT				HWopt. SWopt.			
No Initial Plac.		Initial Placement		No Initial Plac.		Initial Placement	
No migr	Migr	No Migr	Migr	No migr	Migr	No Migr	Migr
No Repl	520	330	380	1540	780	760	780
Repl	250	260	190	670	680	670	670

MG				CG			
No Initial Plac.		Initial Placement		No Initial Plac.		Initial Placement	
No migr	Migr	No Migr	Migr	No migr	Migr	No Migr	Migr
No Repl	230	230	240	1060	700	940	700
Repl	220	220	220	300	280	300	290

71

WildFire's Technology Limits



72



Sun's SunFire 15k/25k

Erik Hagersten
Uppsala University
Sweden



StarCat Sun Fire 15k/25k (used at Lab2)



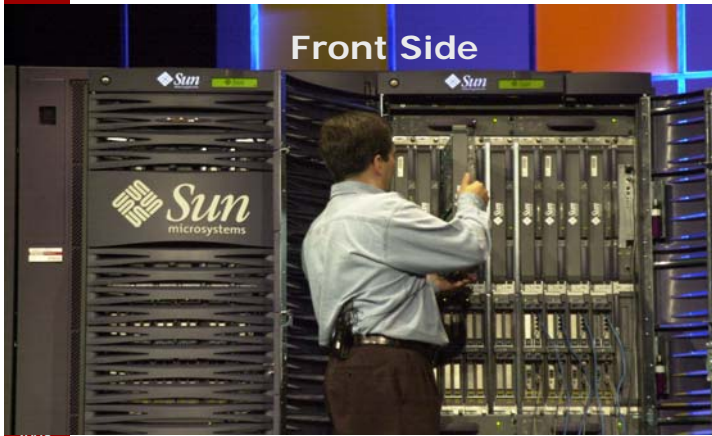
DARK
2008

Dept of Information Technology| www.it.uu.se

74

© Erik Hagersten| user.it.uu.se/~eh

Front Side



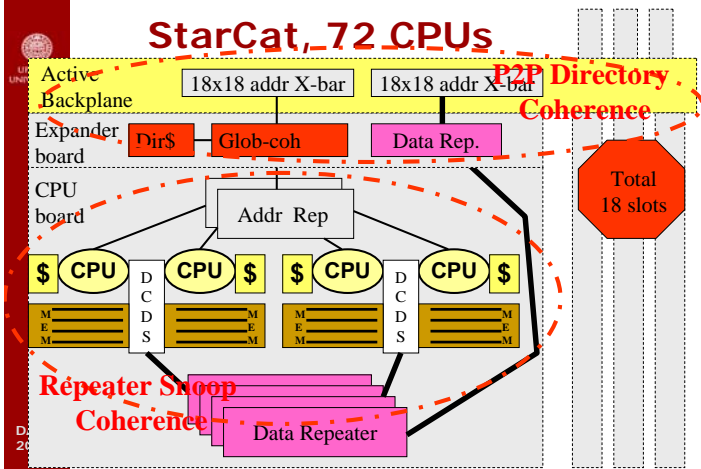
75

Back Side



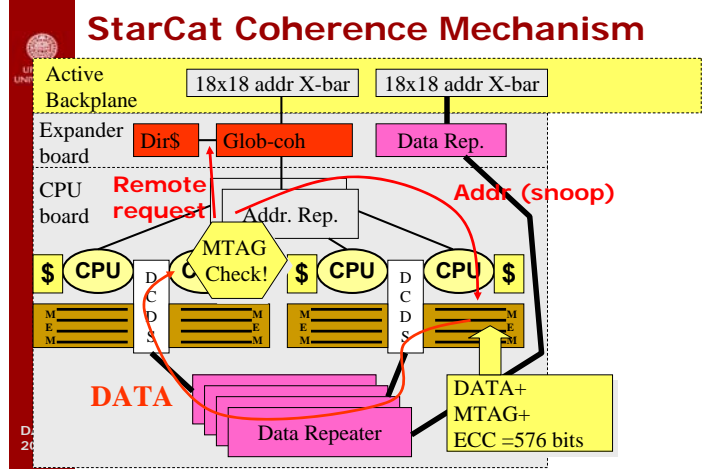
76

StarCat, 72 CPUs



77

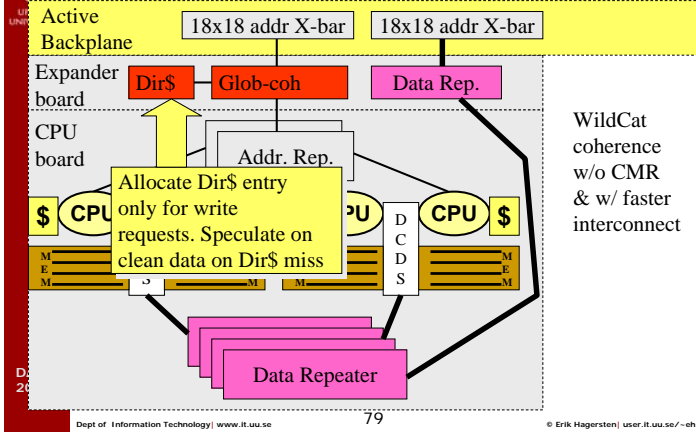
StarCat Coherence Mechanism



78

© Erik Hagersten| user.it.uu.se/~eh

StarCat, 72 CPUs



StarCat Performance Data

