# (2)

- **Distributed Shared Memory: Concepts and Systems**

    Jelica Protic, Milo Tomasevic, and Veljko Milutinovic

    IEEE PARALLEL & DISTRIBUTED TECHNOLOGY,

    Vol. 4, No. 2; SUMMMER 1996, pp. 63-79


- **Distributed Shared Memory Home Pages**
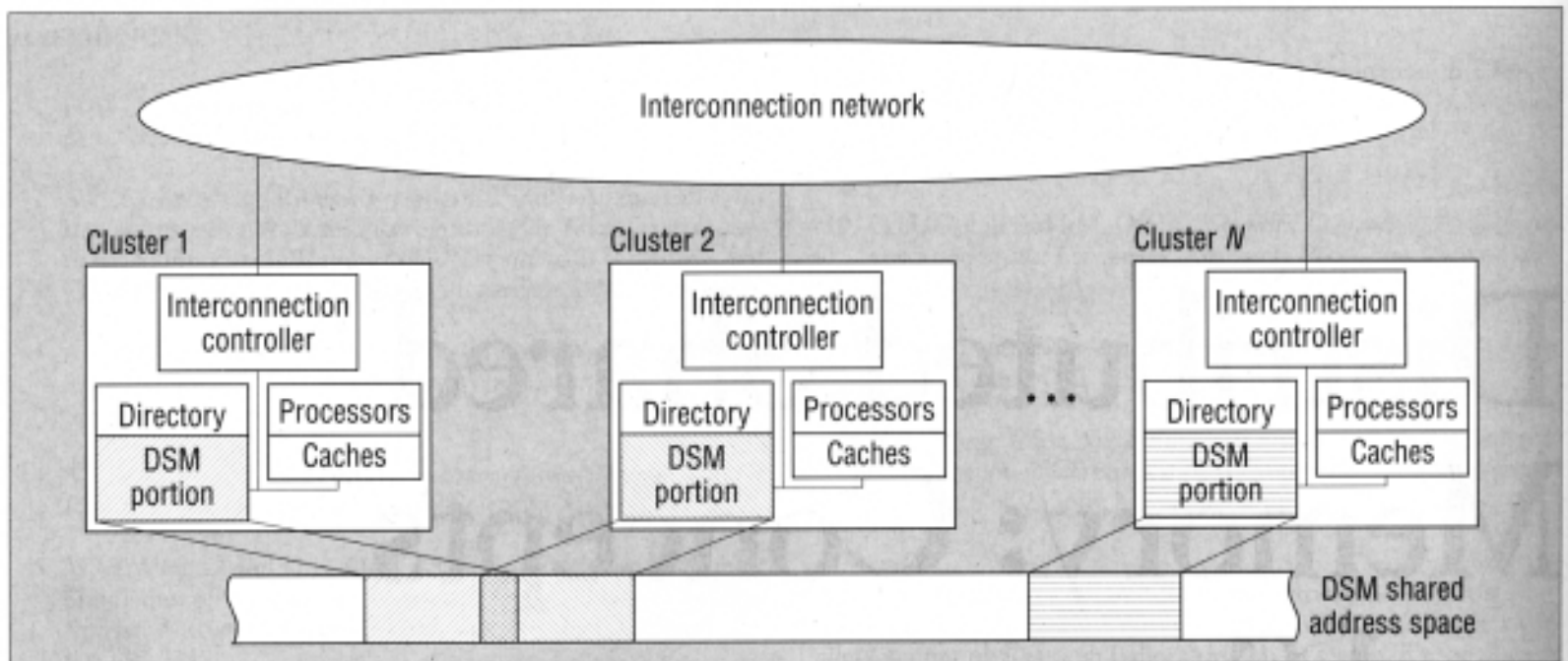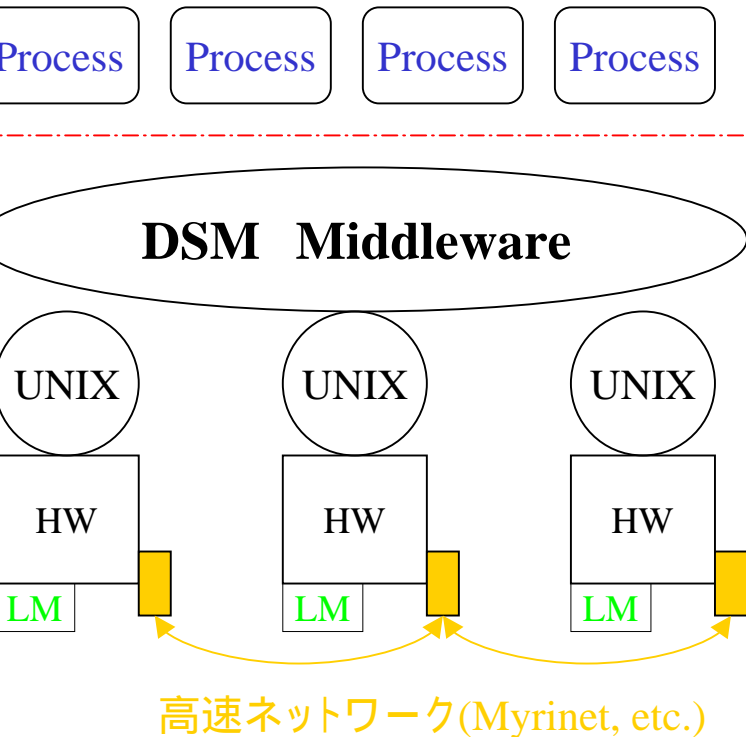  (TreadMarks          P.Keleher          )

    - **http://www.cs.umd.edu/~keleher/dsm.html**

Figure 1. Structure and organization of a DSM system.

| Process | Process | Process | Process |

**DSM  Middleware**

UNIX    UNIX    UNIX

| HW | HW | HW |

| LM | LM | LM |

(Myrinet, etc.)

**(Software) Distributed shared memory**

provides the programmer with the illusion of a single virtual address space, which is shared among a network of processors that do not share physical memory.
As local memory is updated, the modifications are propagated to the other processors, so that all maintain a consistent view.

[          ]
Read Only Data      replication

# Classifications of DSM systems

- How the access actually executes?

- Where the access is implemented?

- What the precise meaning of the word consistent is ?

Latency, Granularity, Availability

# DSM Algorithms

- **Single Reader/Single Writer algorithms**

  DSM

- **Multiple Reader/Single Writer algorithms**

  HW DSM,

- **Multiple Reader/Multiple Writer algorithms**

  Page-based SW DSM          False Sharing


Avenues for performance improvement

  Directory

  SW

# DSM mechanism

- Software DSM implementation
- Hardware DSM implementation
- Hybrid DSM implementation

# Software DSM

Table A. Software DSM implementations.

| IMPLEMENTATION | TYPE OF IMPLEMENTATION | TYPE OF ALGORITHM | CONSISTENCY MODEL | GRANULARITY UNIT | COHERENCE POLICY |
|---|---|---|---|---|---|
| IVY | User-level library + OS modification | MRSW | Sequential | 1 Kbyte | Invalidate |
| Mermaid | User-level library + OS modifications | MRSW | Sequential | 1 Kbyte, 8 Kbytes | Invalidate |
| Munin | Runtime system + linker + library + preprocessor + OS modifications | Type-specific (SRSW, MRSW, MRMW) | Release | Variable size objects | Type-specific (delayed update, invalidate) |
| Midway | Runtime system + compiler | MRMW | Entry, release, processor | 4 Kbytes | Update |
| TreadMarks | User-level | MRMW | Lazy release | 4 Kbytes | Update, invalidate |
| Blizzard | User-level + OS kernel modification | MRSW | Sequential | 32-128 bytes | Invalidate |
| Mirage | OS kernel | MRSW | Sequential | 512 bytes | Invalidate |
| Clouds | OS, out of kernel | MRSW | Inconsistent, sequential | 8 Kbytes | Discard segment when unlocked |
| Linda | Language | MRSW | Sequential | Variable (tuple size) | Implementation-dependent |
| Orca | Language | MRSW | Synchronization dependent | Shared data object size | Update |

# Software DSM implementation

- Write Detection
  - Page-Based  vs  Object-Based
- Coherence Enforcement

---

- IVY               OS level          -----> Shared Virtual Memory
- TreadMarks        User level        -----> Diff//LRC
- Midway            Compiler level-----> Entry Consistency
- Shasta            Compiler level-----> Any
- Linda             Language level----> content addressable
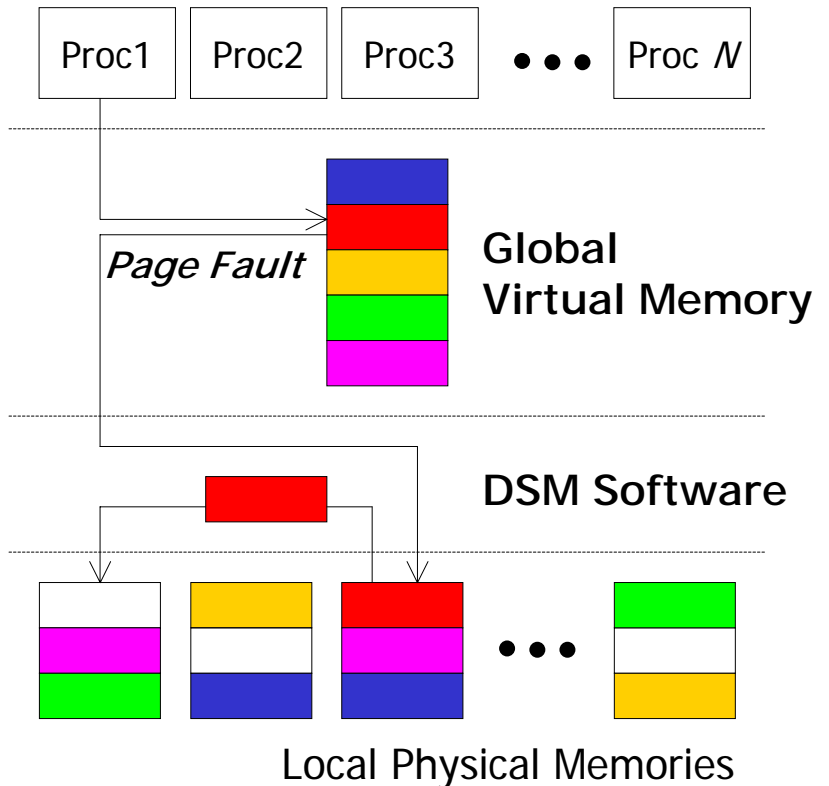                                                    "Tuple" space

# Princeton "IVY"

1KB Page based DSM

OS Modification+
     User-level Library

MRSW

Sequential

Invalidate

| Proc1 | Proc2 | Proc3 | • • • | Proc $N$ |

*Page Fault*

**Global Virtual Memory**

**DSM Software**

**Local Physical Memories**

# Rice     "TreadMarks"

4KB Page based DSM
User-level Library
    TLB    Write       Twin
     Release      Diff
     Acquire     Patch
MRMW
Lazy Release
Update/Invalidate

| False Sharing |
| --- |

X     1st Write

x:

x:     X writeable

"Twin"

Release

Twin:     "Diff"

Encode

x:     X

Non-writeable

# Midway

- **Entry Consistency**

  In Midway, there is <u>an explicit binding of locks to the data</u> that is logically guarded by each lock.
  - As the application acquires a lock for its own synchronization, Midway piggybacks the memory updates on the lock acquisition message. Thus Midway sends no extra messages.
  - Furthermore, the updates are sent only to the acquiring processor and only for the data explicitly guarded by the acquired lock. This serves to batch together updates and minimize the total amount of data transmitted.

- **Midway detects updates to shared memory via compiler and runtime support.**

- To provide high performance communication, Midway has its own application oriented protocols which reduce message counts, and it utilizes Mach's low-overhead network interfaces to reduce message latency.

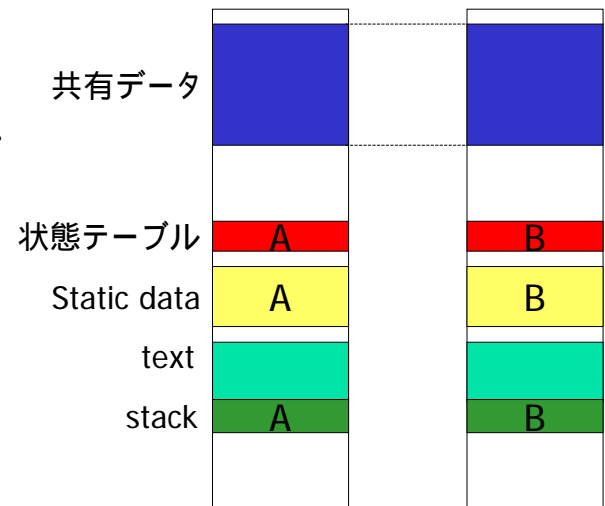# DEC WRL "Shasta"

Fine Grain Software DSM

Compiler Level

| Miss Handling and Protocol Code |
|---|

| Application Execututable for SMP | → | Shasta Compiler | → | Application Execututable with Miss Checks |
|---|---|---|---|---|

| Message Passing Library |
|---|

SW
miss check

MRSW
Any Consistency Model
Any Protocol

| | | |
|---|---|---|
| | A | B |
| Static data | A | B |
| text | | |
| stack | A | B |

**Shasta**

# Hybrid DSM implementation

- ## SHRIMP @Princeton Univ.
  - Virtual Memory Mapped I/O
  - Automatic Update Release Consistency(AURC)



Figure 2: *Virtual memory mapping*

# Hybrid DSM implementation

- ## SHRIMP@Princeton Univ.
  - ### Virtual Memory Mapped I/O
  - ### Automatic Update Release Consistency(AURC)
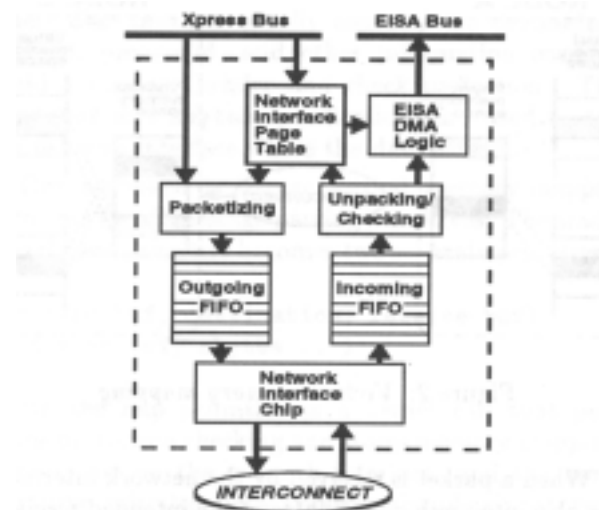


Figure 1: *A SHRIMP node with network interface*



Figure 3: *SHRIMP Network Interface Data Path*

# Hardware DSM

- CC-NUMA
  - Directory-based
    - → JUMP-1, Cenju-4, Origin2000, AsamA,
  - Broadcast-based ……. Reflective Memory
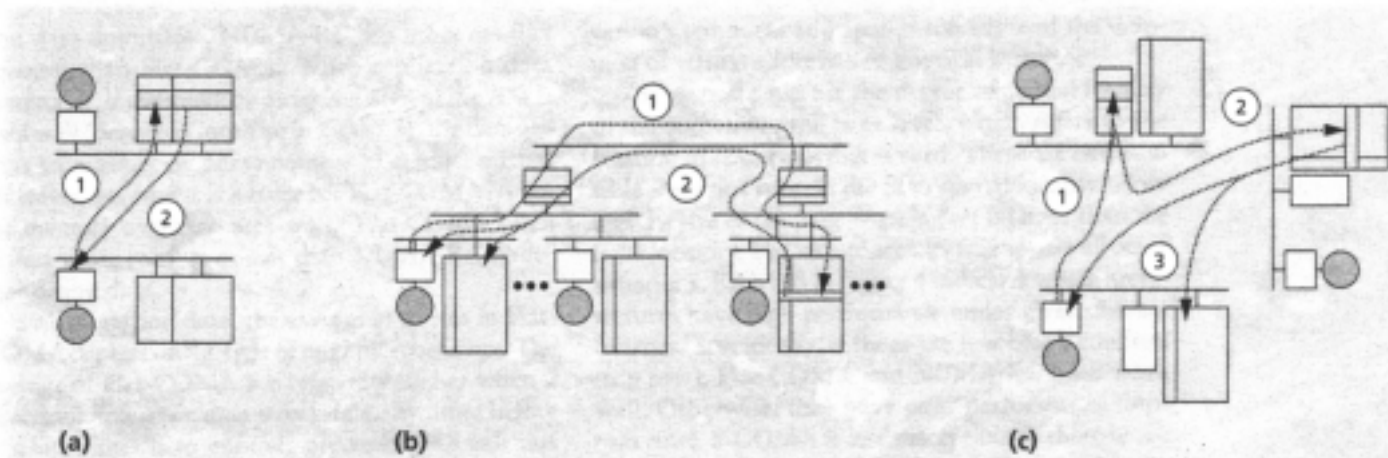    - → Memory Channel
- COMA  Family

# Hardware DSM

- COMA  Family



Figure 1. Node organization in (a) NUMA, (b) Hierarchical-COMA, and (c) Flat-COMA.

C Cache
P Processor
Dir Directory

# Memory Consistency Models

# Memory Consistency Models

## I-structure

- 

-     HW

P: _____ , A: _____ , W: _____



```
      n:    P │   data
    n+1:    A │
    n+2:    W │          ●
    n+3:    W │          ●
    n+4:    A │

    n+m:    P │   data
```
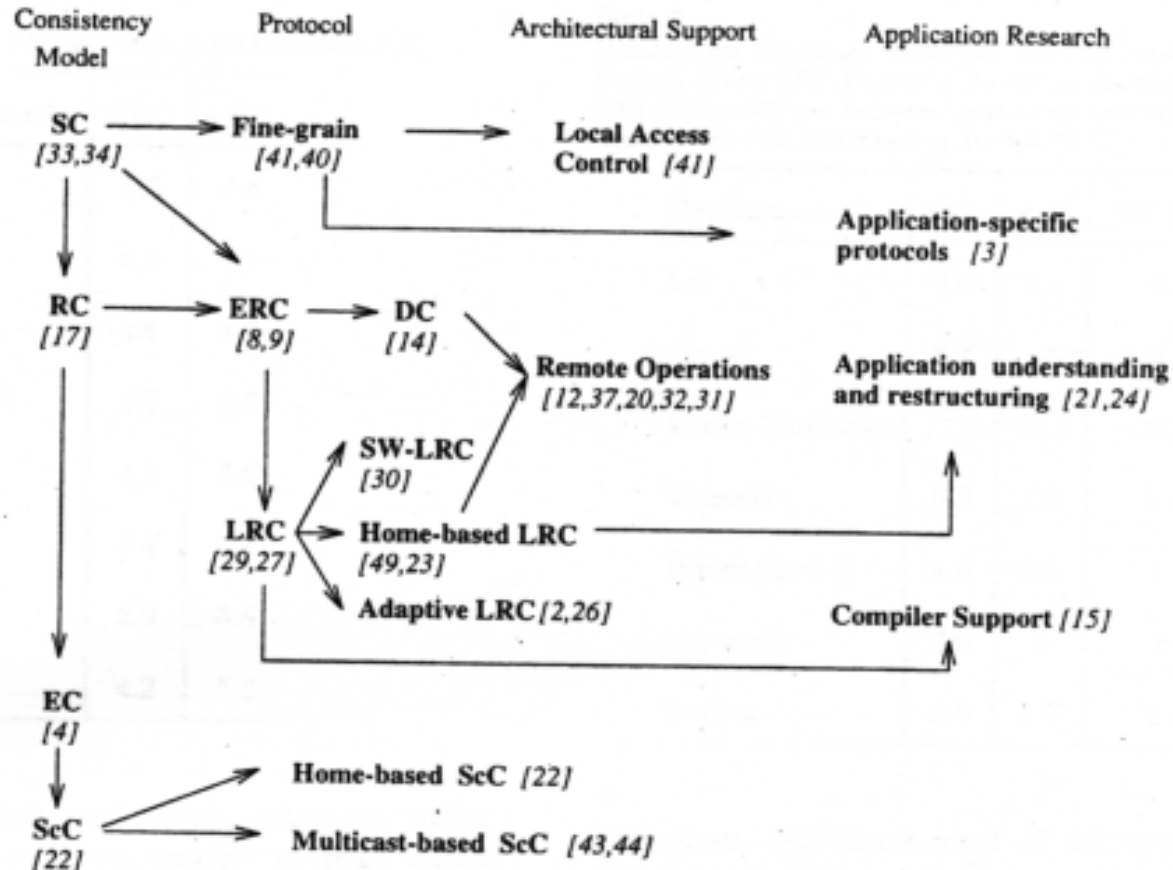
A
C
B

# Memory Consistency Models



Fig. 1. Research in SVM. The figure treats lazy release consistency (LRC) and eager release consistency (ERC) as different protocols implementing the RC consistency model, though they are in fact slightly different consistency models. SW-LRC is single-writer LRC protocol. SC, RC, DC,

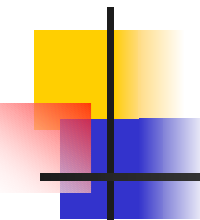# Important design choices in building DSM systems

- Cluster configuration
- Interconnection networks
- Shared data structure
- Coherence unit granularity
- DSM management responsibility
- Coherence policy

A4 1~2

8B10B

Wave Pipeline

SMT(or HT by Intel)

Memory Consistency Model          LRC(Lazy Release Consistency)

29

(                              1/22                              )

PDF