

(2)分散共有メモリシステム

- 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
 - <http://www.ics.uci.edu/~javid/dsm.html>
 - <http://www.cs.umd.edu/~keleher/dsm.html> (TreadMarks開発者のP.Keleherが管理していたが2004/01現在上記URLに変わっている)

分散共有メモリ環境の実現

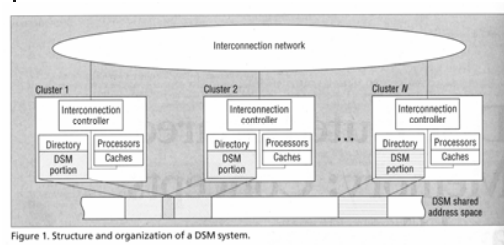
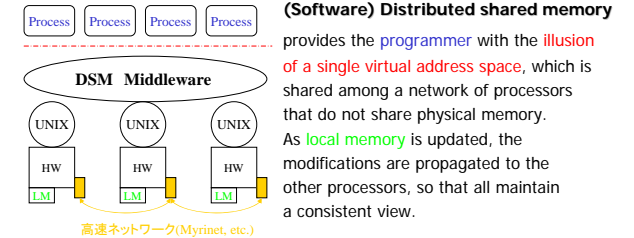


Figure 1. Structure and organization of a DSM system.

分散共有メモリ環境の実現



(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.

[極端な例]
Read Only Data の replication だけを提供

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
⇒機能限定DSM
- Multiple Reader/Single Writer algorithms
⇒多くのHW DSM, 無効化型の一貫性制御
- Multiple Reader/Multiple Writer algorithms
⇒Page-based SW DSMにおけるFalse Sharing対策
更新型の一貫性制御

Avenues for performance improvement
Directory 構成法と一貫性制御プロトコル
SWオーバーヘッドの軽減

DSM mechanismの実装レベル

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

Software DSM

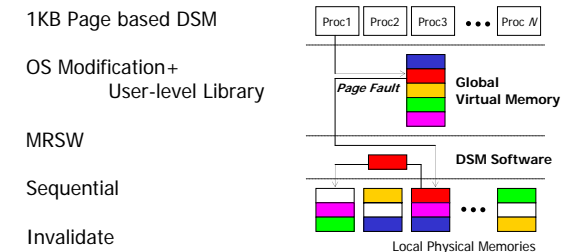
IMPLEMENTATION	TYPE OF IMPLEMENTATION	TYPE OF ALGORITHM	CONSISTENCY MODEL	GRANULARITY UNIT	COHERENCE POLICY
IVY	User-level library + OS modification	MRSW	Sequential	1 Kbyte	Invalidate
Mermaid	User-level library + OS modifications	MRSW	Sequential	1 Kbyte, 8 Kbytes	Invalidate
Munin	Runtime system + linker + library + preprocessor + OS modifications	Type-specific (SRW, MRSW, MRMW)	Release	Variable size objects	Type-specific (delayed updates, invalidate)
Midway	Runtime system + compiler	MRSW	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大学 "IVY"



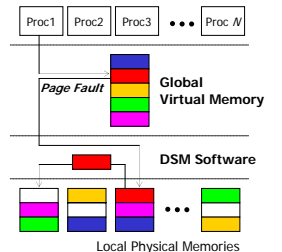
1KB Page based DSM

OS Modification +
User-level Library

MRSW

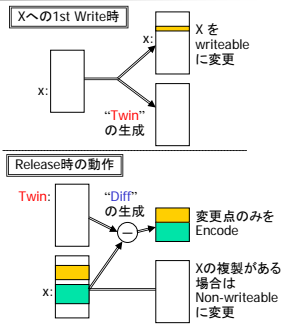
Sequential

Invalidate



Rice大学 “TreadMarks”

4KB Page based DSM
 User-level Library
 TLBでWriteを検出 Twin作成
 Release時に Diff を作成
 Acquire時に Patch をあてる
 MRMW
 Lazy Release
 Update/Invalidate

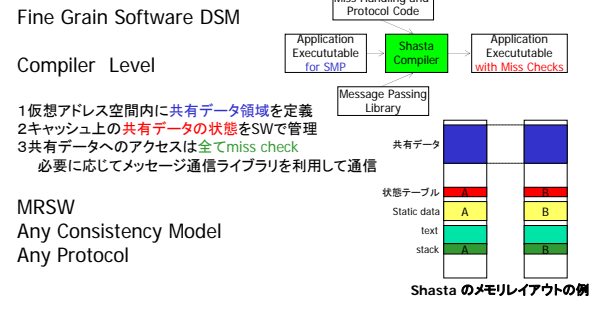


False Sharing対策
 無駄な一貫性制御の排除

Midway

- Entry Consistency**
 In Midway, there is an explicit binding of locks to the data 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

- 1 仮想アドレス空間内に共有データ領域を定義
- 2 キャッシュ上の共有データの状態をSWで管理
- 3 共有データへのアクセスは全てmiss check
 必要に応じてメッセージ通信ライブラリを利用して通信

MRSW
 Any Consistency Model
 Any Protocol

Hybrid DSM implementation

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

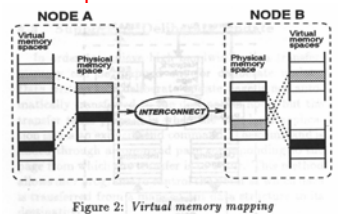


Figure 2: Virtual memory mapping

Hybrid DSM implementation

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

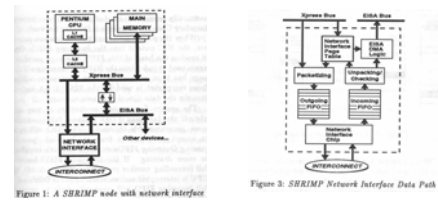


Figure 1: A SHRIMP node with network interface
 Figure 3: SHRIMP Network Interface Data Path

Hardware DSM 実装

- CC-NUMA
 - Directory-based
 → JUMP-1, Cenju-4, Origin2000, Asama, その他多数
 - Broadcast-based Reflective Memory
 → Memory Channel
- COMA Family

Hardware DSM 実装

- COMA Family

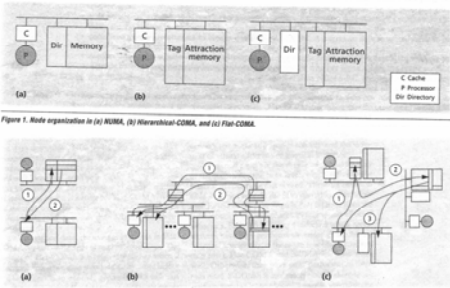
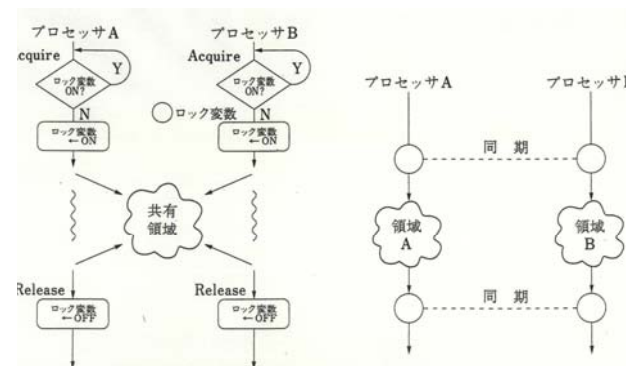


Figure 1. Node organization in (a) NUMA, (b) Hierarchical-COMA, and (c) Flat-COMA.

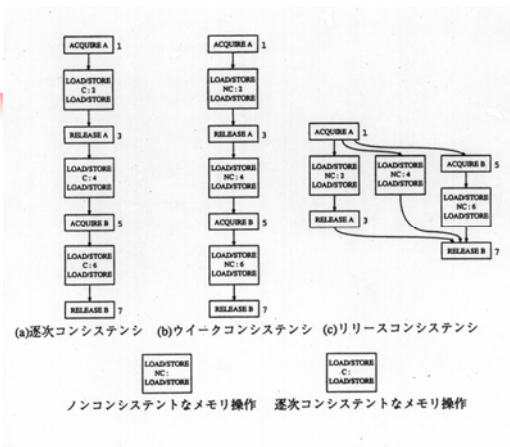
Memory Consistency Models

- 仮定1**
 データを共有する場合には、生産者と消費者が定義でき、プログラムの不確実性をなくすためには生産者と消費者の間には必ず何らかの同期が存在する
- 仮定2**
 同期のためのメモリ操作 と それ以外のメモリ操作が 区別可能である。



(a) 臨界領域へのアクセス (b) 異なる領域へのアクセス

図 4.31 共有データへのアクセス

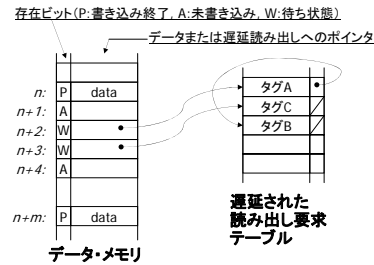


Memory Consistency Models

「特殊ケース」同期と通信の融合

I-structure

- プログラミングモデルで単一代入則を保证
- 専用HWによる同期機構 (マッチングメモリ)



Memory Consistency Models

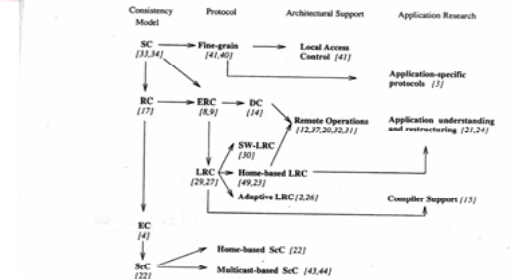


Fig. 1. Research in SVM. The figure treats lazy release consistency (LRC) and eager release consistency (ERC) as different protocols implementing the RC consistency model, though they are in fact slightly different consistency models. SW-LRC is single-writer LRC protocol. SC, RC, DC, EC, and SC are the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

Important design choices in building DSM systems

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

本資料のPDF版 と 講義に関する参考文献リスト を以下からアクセス(学内のみ)できるようにしています。

<http://www.lab3.kuis.kyoto-u.ac.jp/members/moris/lecture/PDS/>

(最終版は全講義終了後に掲載します)

レポート課題(森担当分)

以下について、A4 2~3ページ程度にまとめて報告。

- 予算規模1億円程度のクラスタ計算機を概念設計せよ。この際、設計した計算機の特徴、用途、実現可能性について議論せよ。
- 以下のいずれか一方を調査せよ。
 - Wave Pipeline技術について調査せよ。
 - Software Transactional Memory (STM) について調査せよ。この際、STMを効率的に実現するにはどのようなハードウェア支援が必要かを明記すること。

提出方法: レポートには表紙をつけ、タイトルとして「並列分散システム論 レポート(森担当分)」と書いた上で、1. 氏名 2. 学籍番号 3. 入進学年(学年) 4. 所属(研究科, 専攻, 研究室)を明記して、ホッチキス(stapler)等で固定し提出。
 締切: 12/18(火)
 提出先: 情報学研究科事務室教務窓口(工学部10号館)前の「大学院」と書かれたレポートボックス