

第 5 章 マルチプロッサ

5 . 1 概論

5 . 1 . 1 分類

「超」の規模：1 TFLOPSを実現できる

プロセッサ台数

地球シミュレータ：4 0 TFLOPS

ASCI プロジェクト

(1) プロセッサアーキテクチャ

SPP (Scalar Parallel Processor)

1TFLOPSの実現：1千台程度

マイクロプロセッサ：1 GFLOPS

スーパスカラ方式

VLIW方式

スーパーパイプライン方式

マルチスレッド方式

VPP (Vector Parallel Processor)

メモリ系単純化

ベクトルプロセッサの性能：10 GFLOPS程度

「超」並列：100台程度

台数の少ない構成方式：実行性能優

制御単純、理解しやすい、ユーザからの連続性

VPPが優位

SPP：最新鋭のマイクロプロセッサを安価に利用

通信やメモリのレイテンシ（遅延）に耐え得

るアーキテクチャ

（２）メモリアーキテクチャ

共有か非共有か

- ・ 共有メモリ方式
- ・ 非共有メモリ方式

集中か分散か

集中共有、分散共有、

分散非共有、集中非共有の4つの場合

(3) ネットワークアーキテクチャ

トポロジ

交換方式 (ストアド・フォワード、ワームホール)

並列計算機：特定の並列アーキテクチャと

特定応用との結び付き

数種の並列アーキテクチャに収斂する必要

5.1.2 基本方式

(1) 粒度

プログラム分割

プロセッサ割付け

負荷の分散と粒度の関係

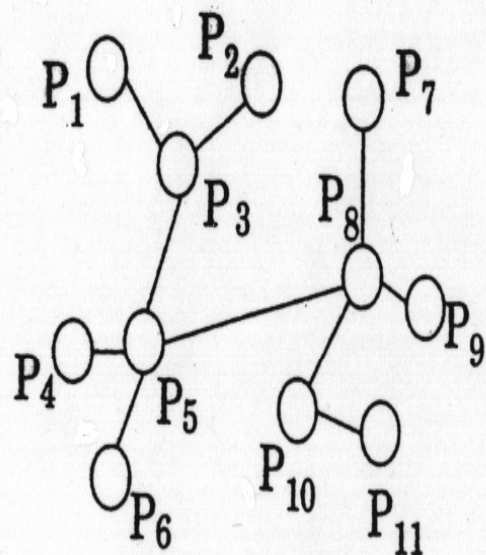
(2) プログラム分割

メモリ共有方式

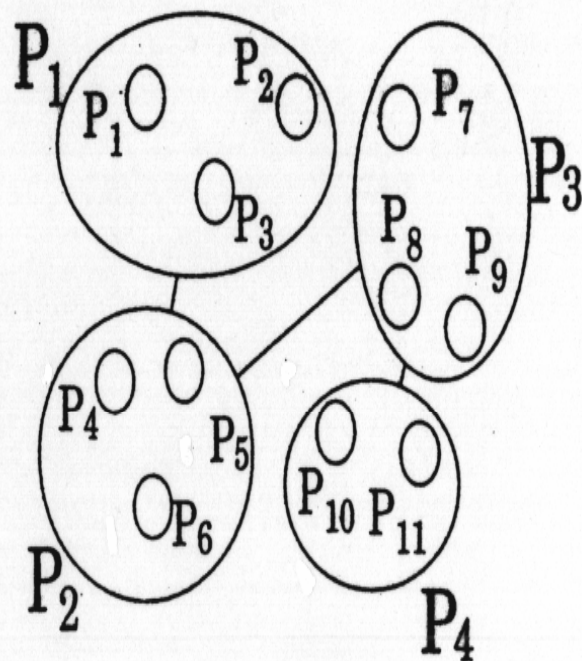
分割：命令系列のみ

データ：共有

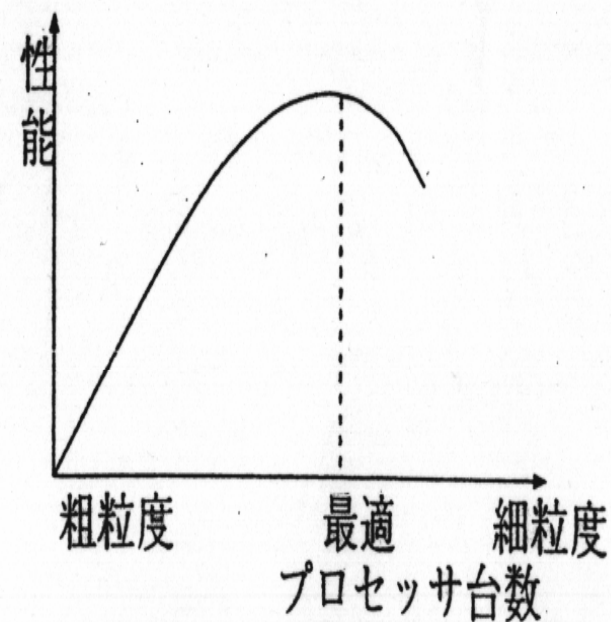
処理内容・データ構造：



(a) 細粒度



(b) 粗粒度



(c) 粒度の性能に及ぼす影響

非定型、コンパイル時にはデータ分割不可
のような非定型応用

共有メモリ：

物理的に集中 / 分散共有方式で実現

細粒度並列処理向き

小規模のマルチプロセッサ：

集中共有メモリ方式

大規模なシステム：

分散共有メモリ方式

メモリ共有を分散非共有方式で実現

ワークステーションを用いたネットワ
ークコンピューティング
仮想（論理）アドレス空間
ページフォールト

メッセージ交換方式

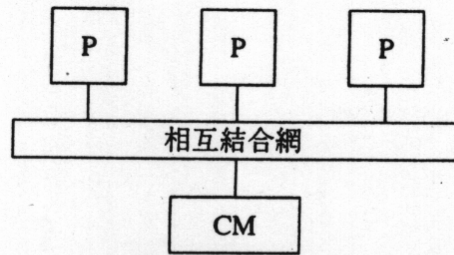
分割：命令系列とデータ

処理内容：定型的応用

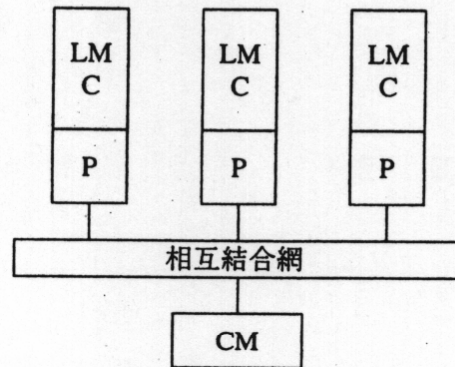
オブジェクト指向モデル

離散系シミュレーション

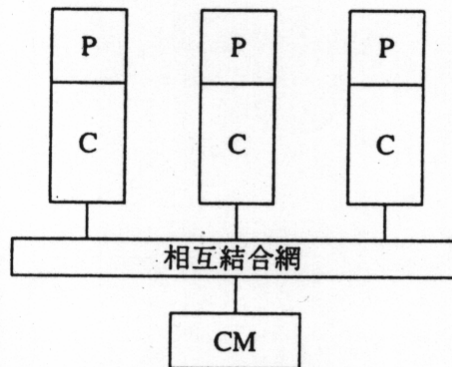
粗粒度並列処理



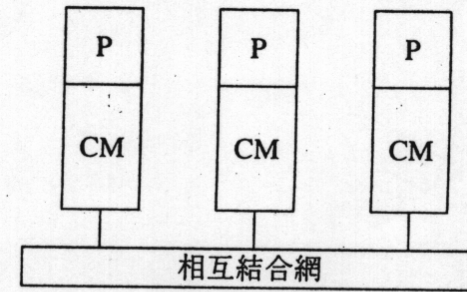
(a) 集中共有メモリ



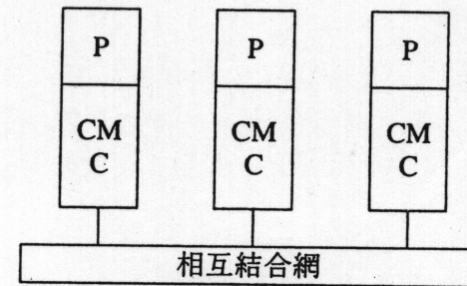
(b) 集中共有メモリ (ローカルメモリ付)



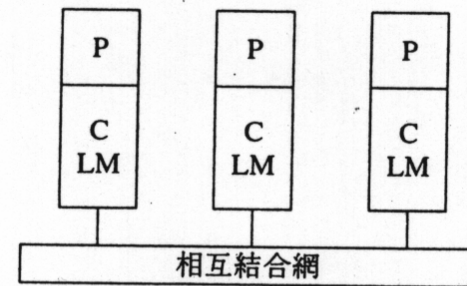
(c) 集中共有メモリ (キャッシュ付)



(d) 分散共有メモリ



(e) 分散共有メモリ (キャッシュ付)



(f) メッセージ交換 (分散非共有メモリ)

P: プロセッサ
 CM: 共有メモリ
 LM: ローカルメモリ (非共有メモリ)
 C: キャッシュメモリ

表 5.1 マルチプロセッサ例

機種	結合網	台数(プロセッサ)	性能	メモリ
CM-5	トリー	16000ノード(1ノード: 1SPARC+ 4ベクトル装置)	2T	非共有
GC	3次元トーラス	1024クラスタ(1クラスタ: 16台のT-9000)	400G	非共有
Paragon	2次元メッシュ	4096 (iPSC860)	300G	非共有
nCUBE-3	ハイバキューブ	65000 (独自プロセッサ)	6.5T	非共有
TERA-3D	3次元トーラス	2048(DEC α)	300G	共有
SP2	多段網	512ノード(1ノード: 16台の RS/6000)	136G	非共有
KSR	階層リング	1088 (独自プロセッサ)	43.5G	共有キャッシュ (キャッシュオンリ)
Exemplar	リング	16クラスタ(1クラスタ: 8台 のPA, クロスバ)	198M	共有キャッシュ (SCI プロトコル)
AP-1000	2次元トーラス	1024 (SPARC)	8.53G	非共有
VPP-500	クロスバ	222 (1.6Gベクトル, 300M スカラ (VLIW))	355G	共有
SX-4	クロスバ(ノード内) 光接続(ノード間)	16ノード(1ノード: 32台の PE, 1PE: 2G ベクトル)	1T	共有(ノード内) 非共有(ノード間)
Cenju-3	多段網	256 (V_R 4400)	12.8G	非共有
ADENART	ハイバクロス	256 (独自)	2.56G	非共有
Cyberflow	2次元トーラス	64 (データフロー)	640M	非共有
CP-PACS	ハイバクロスバ	1024 (PA+機能追加)	300G	非共有
SR2201	ハイバクロスバ	1024 (PA+機能追加)	300G	非共有
JUMP-1	RDT	512 (SuperSPARC+)	30G	共有キャッシュ

T: TFLOPS, G: GFLOPS, M: MFLOPS, 共有: 共有メモリ方式,
 共有キャッシュ: キャッシュコヒーレンス制御付き共有メモリ方式,
 非共有: メッセージパッシング方式, PA: HP社製マイクロプロセッサ,
 SPARC: サンマイクロシステムズ社製マイクロプロセッサ,
 V_R : MIPS社製マイクロプロセッサ

Institution	Name	Number of nodes	Basic topology	Data bits/link	Network clock rate (MHz)	Peak BW/link (MB/sec)	Bisection (MB/sec)	Year
Thinking Machines	CM-2	1024 to 4096	12-cube	1	7	1	1,024	1987
Intel	Delta	540	2D grid	16	40	40	640	1991
Thinking Machines	CM-5	32 to 2048	multistage fat tree	4	40	20	10,240	1991
Intel	Paragon	4 to 2048	2D grid	16	100	175	6,400	1992
IBM	SP-2	2 to 512	multistage fat tree	8	40	40	20,480	1993
Cray Research	T3E	16 to 2048	3D torus	16	300?	600	122,000	1997
Intel	ASCI Red	4536 ($\times 2$ CPUs)	2D grid			800	51,600	1996
IBM	ASCI Blue Pacific	1336 ($\times 4$ CPUs)				150		
SGI	ASCI Blue Mountain	1464 ($\times 2$ CPUs)	fat hypercube			800	200 \times nodes	1998
IBM	ASCI Blue Horizon	144 ($\times 8$ CPUs)	multistage Omega			115		1999
IBM	SP	1 to 512 ($\times 2$ to 16 CPUs)	multistage Omega			500		2000
IBM	ASCI White	484 ($\times 16$ CPUs)	multistage Omega			500		2001

Figure 8.18 Characteristics of interconnections of some commercial supercomputers. The bisection bandwidth is for the largest machine. The 2D grid of the Intel Delta is 16 rows by 35 columns and the ASCI Red is 38 rows by 32 columns. The fat-tree topology of the CM-5 is restricted in the lower two levels, hence the lower bandwidth in the bisection. Note that the Cray T3E has two processors per node, and the Intel Paragon has from two to four processors per node.

5 . 2 メモリ共有型

キャッシュの装備を前提

データ参照の局所性を生かすこと

無駄な通信を極力削減すること

5 . 2 . 1 構成方式

集中共有メモリ

集中共有メモリ（ローカルメモリ付き）

集中共有メモリ（キャッシュ付き）

分散共有メモリ

分散共有メモリ（キャッシュ付き）

共有メモリ

UMA (Uniform Memory Access) モデル

NUMA (Non-Uniform Memory Access) モデル

4.5 メモリコンシステンシモデル

書込みの順序づけ (ordering)

プロセッサA,Bが書込み

プロセッサC: Aの後Bが到着

プロセッサD: Bの後Aが到着

これでよいのか？

4.5.1 プロセッサコンシステンシモデル

プロセッサA

プロセッサB

```
data = new;          while (flag != set) {}
```

```
flag = set;          data copy = data;
```

同一のプロセッサからのライト:

その順で反映

2つ以上のプロセッサで発せられた書込み:

順序については何も制約なし

プロセッサ A

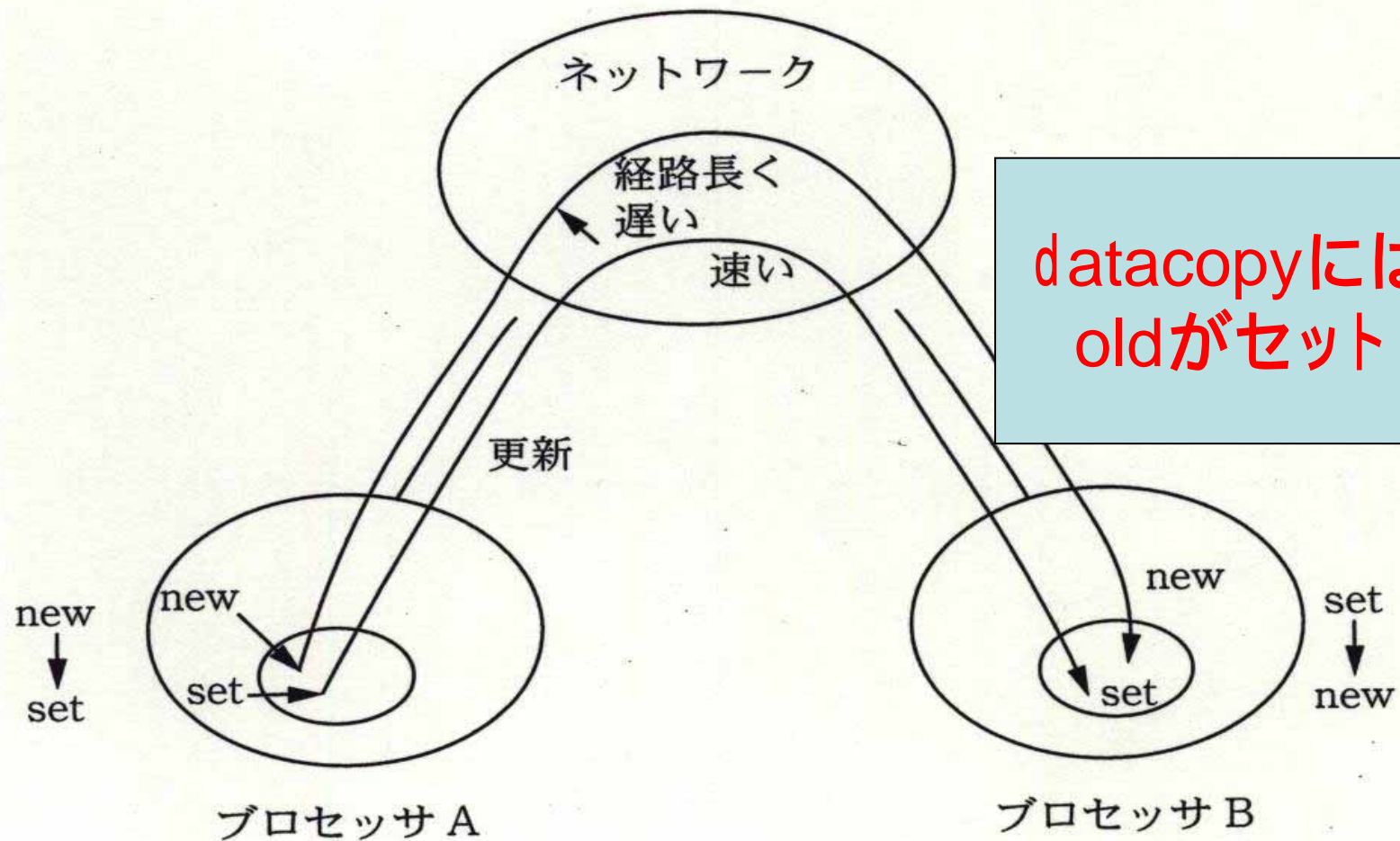
① data=new ;

② flag=set ;

プロセッサ B

③ while (flag != set) {}

④ datacopy=data ;



(a) プロセッサコンシステンシモデル違反

4.5.2 逐次 (sequential) コンシステンシ

$X = Y = 0$

プロセッサ A

プロセッサ B

$X = 1$

$Y = 1$

IF $Y = 0$ KILL B

IF $X = 0$ KILL A

A,Bともに相打ちは
起こるか？

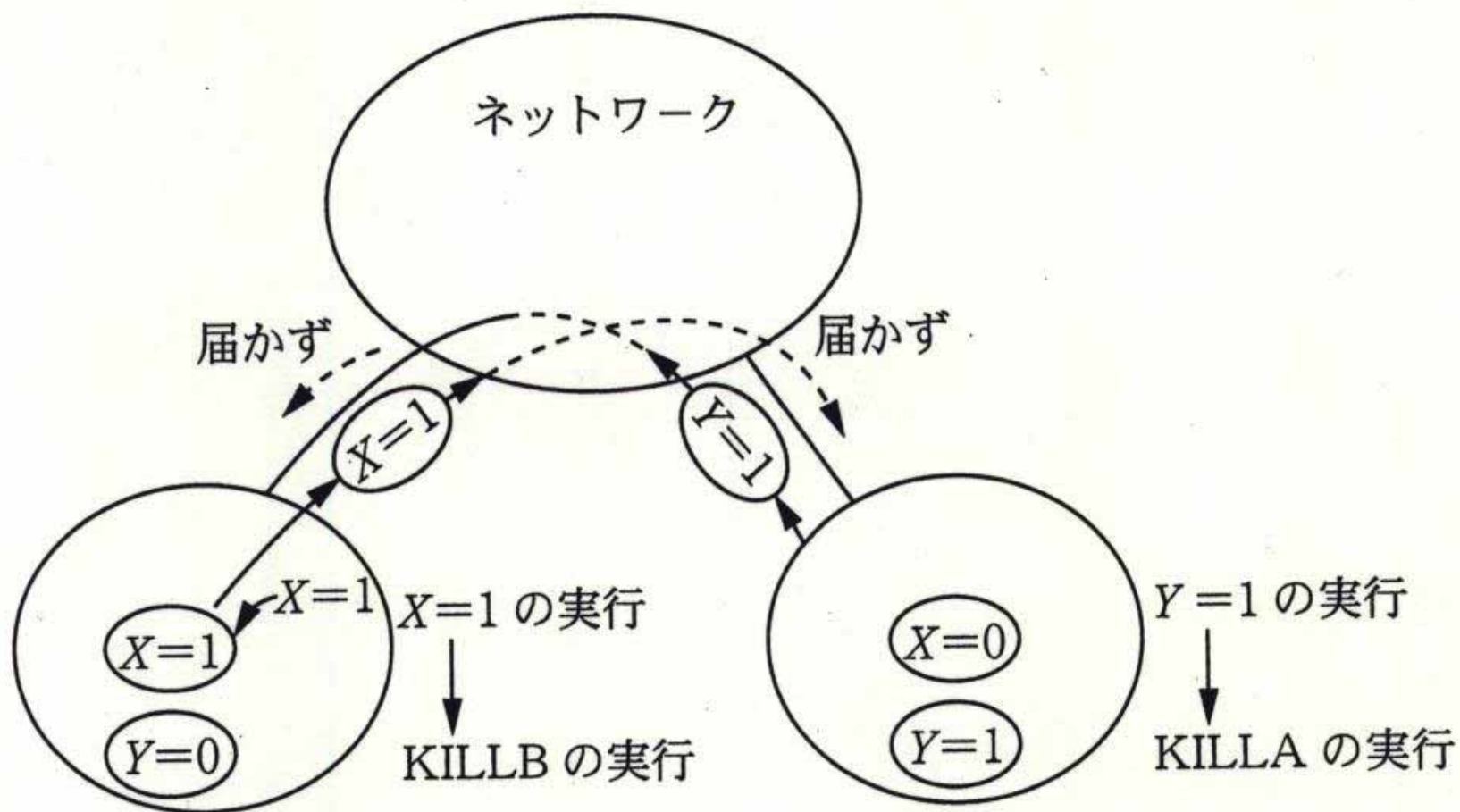
(1) Lamportの逐次コンシステンシ定義

マルチプロセッサ上での

並列プログラムの実行結果 =

並列プログラムを単一のプロセッサで

時分割で実行したときと同一



(b) 逐次コンシステンシモデル違反

図 4.29 メモリコンシステンシ問題

(2)十分条件

並べられた命令順に実行

メモリ操作の大域的完了後、後続メモリ操作
発行

大域的書込み完了:

書込み操作がすべてのプロセッサに反映

大域的読込み完了:

リードデータが大域的に書込み完了

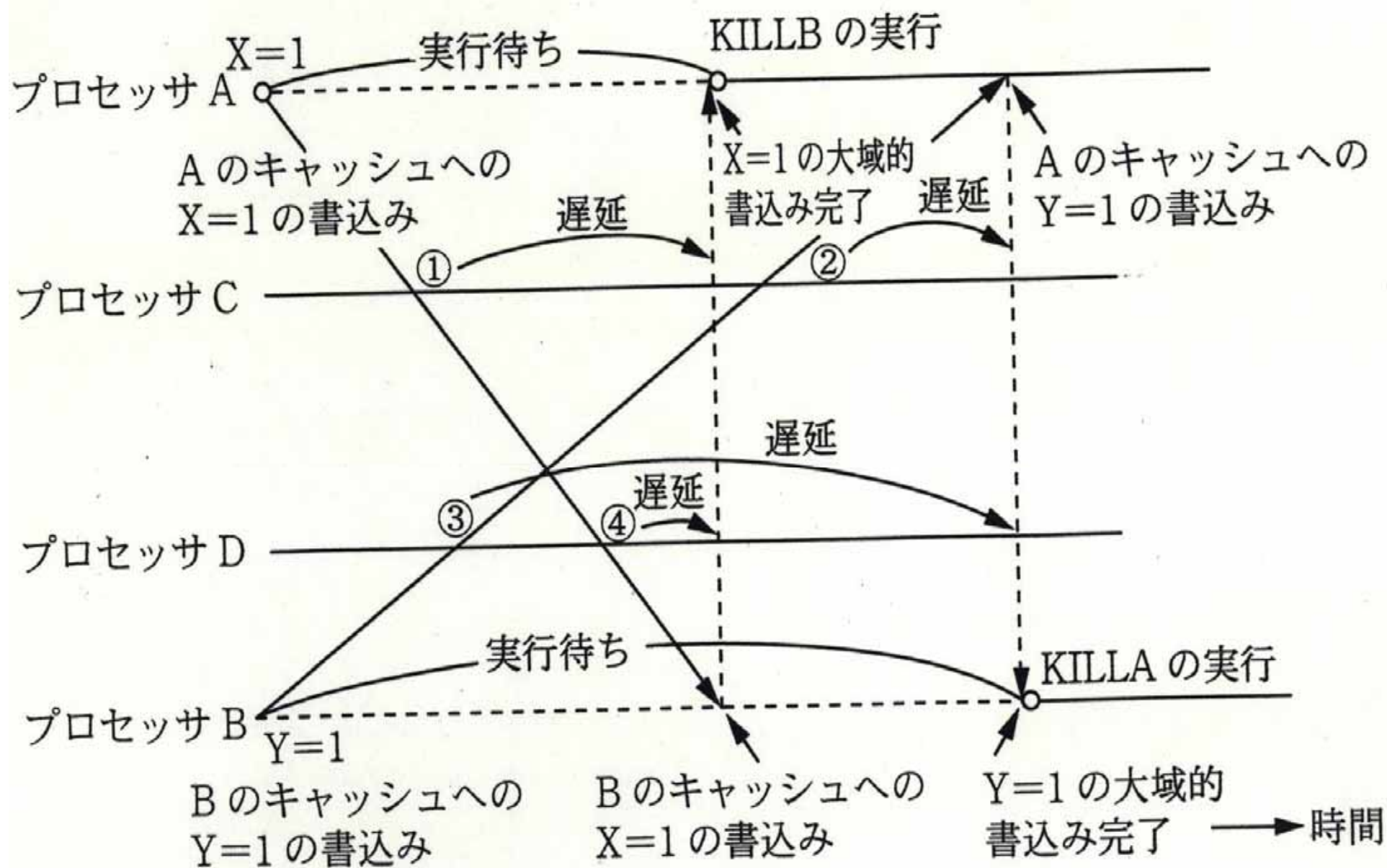


図 4.30 逐次コンシステンシの実現

(3) 逐次コンシステンシの緩和

臨界領域へのアクセス

異なる領域へのアクセス

Weakコンシステンシ

同期命令

先行命令の大域的完了後発行

後続命令は同期命令の大域的完了後発行

Releaseコンシステンシ

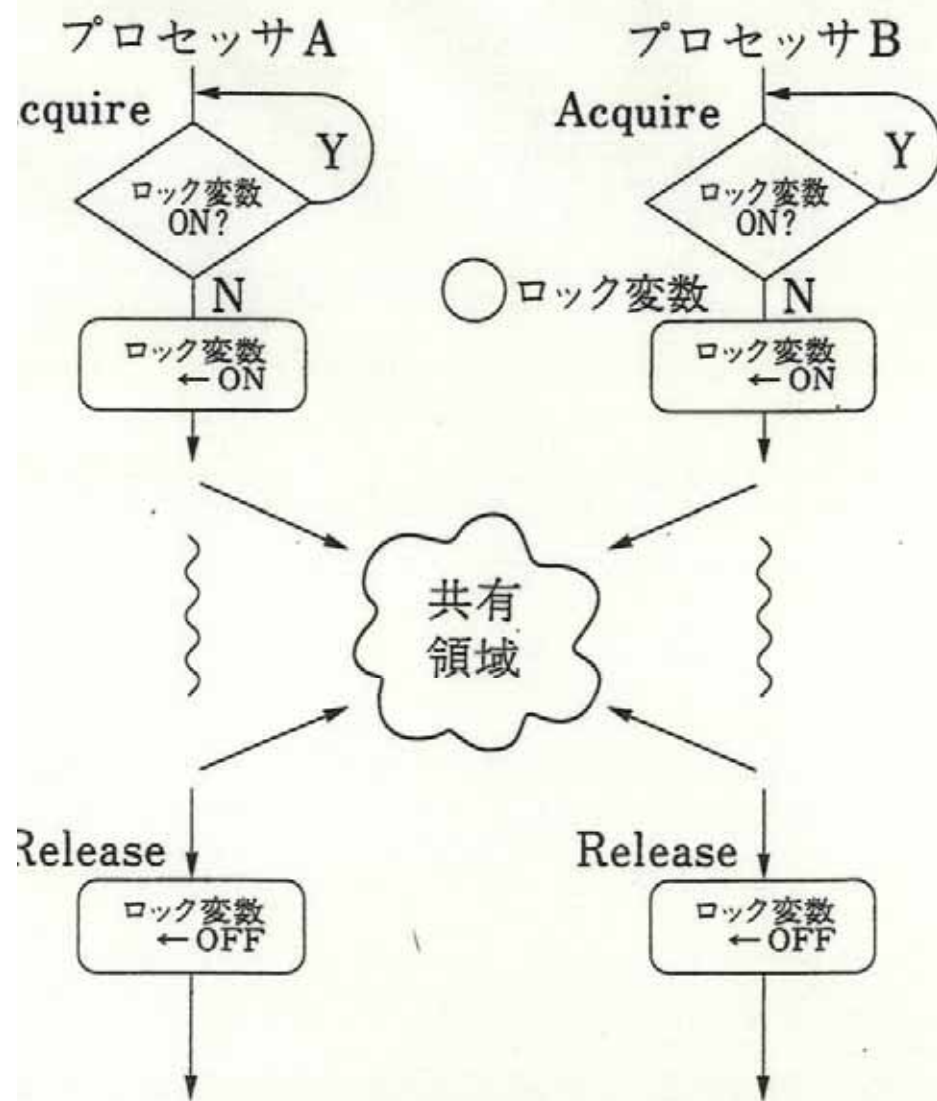
Acquire同期命令

後続命令はAcquireの大域的完了後発行

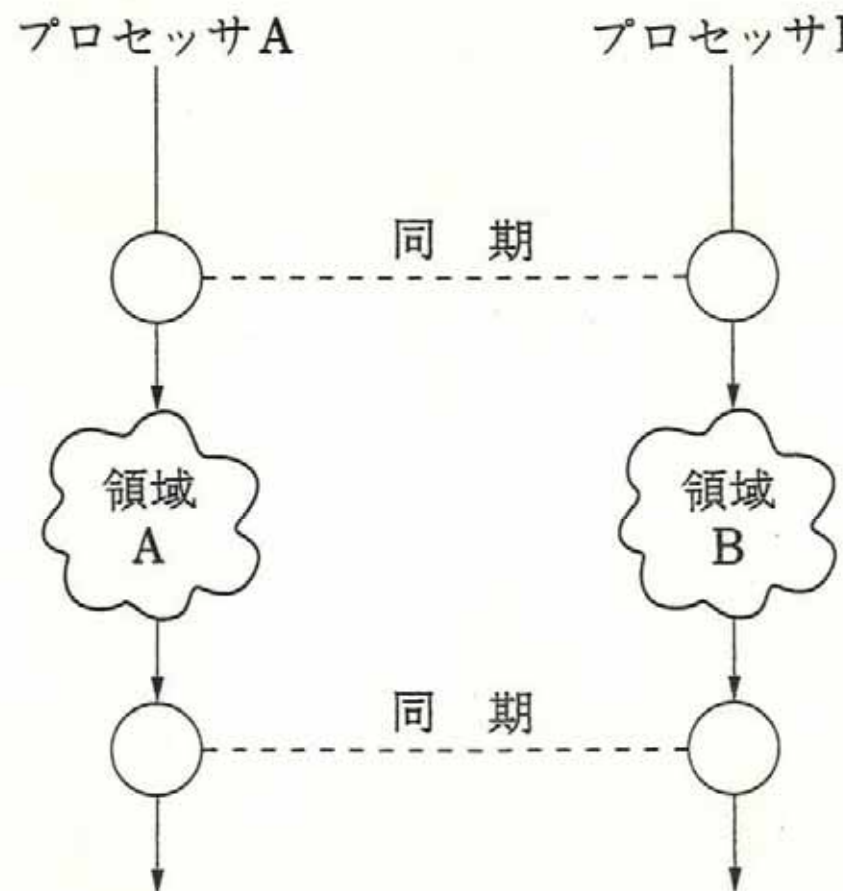
Release同期命令

先行命令の大域的完了後発行

後続命令は発行可

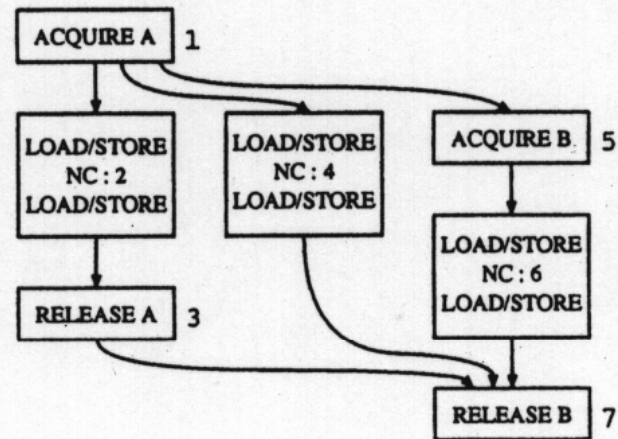
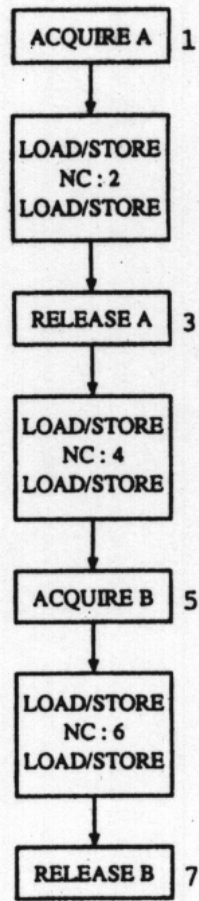
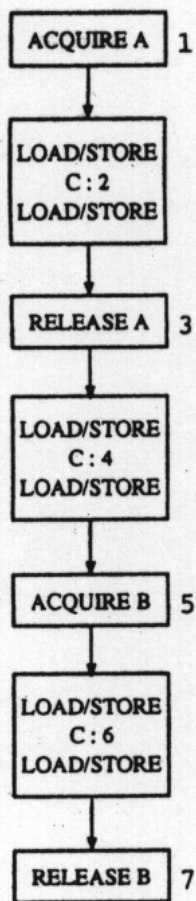


(a) 臨界領域へのアクセス

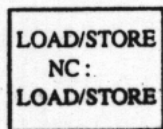


(b) 異なる領域へのアクセス

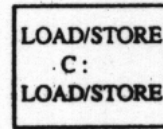
図 4.31 共有データへのアクセス



(a)逐次コンシステンシ (b)ウイークコンシステンシ (c)リリースコンシステンシ



ノンコンシスレントなメモリ操作



逐次コンシスレントなメモリ操作

4.4 キャッシュコヒーレンス

(1) ハードウェアによる方式

スヌープキャッシュ方式

ディレクトリ方式

(2) ソフトウェアによる方式

4.4.1 スヌープキャッシュ法

(1) 単純な方式の場合

スヌープコントローラ

(2) 各ブロックにタグを持たす方式

ライト時に無効化または更新

キャッシュ間転送時の主記憶に書き戻し

(3) バスの特徴

放送能力: 分散制御

排他制御: 逐次コンシステンシ

スヌープ方式

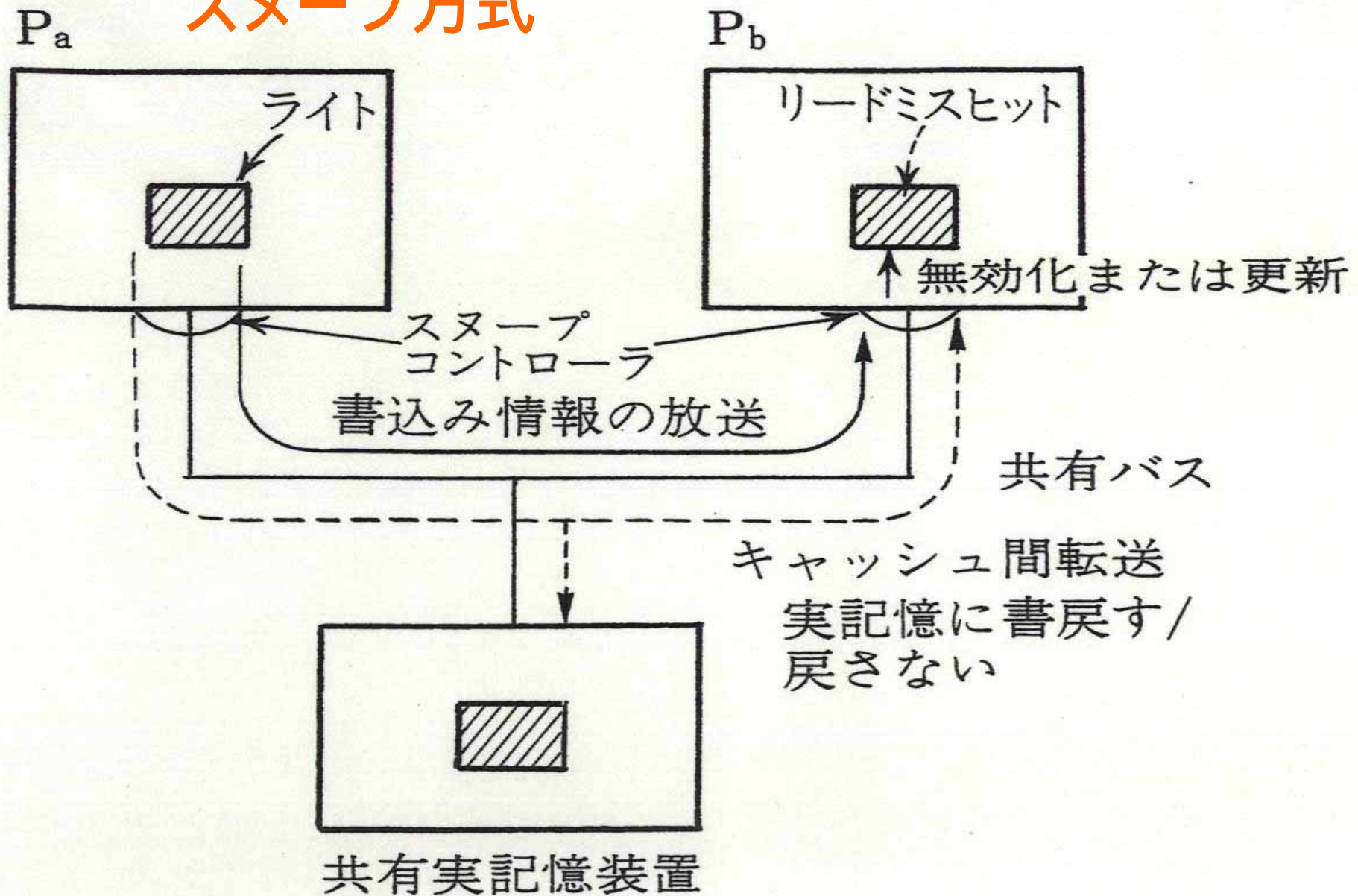
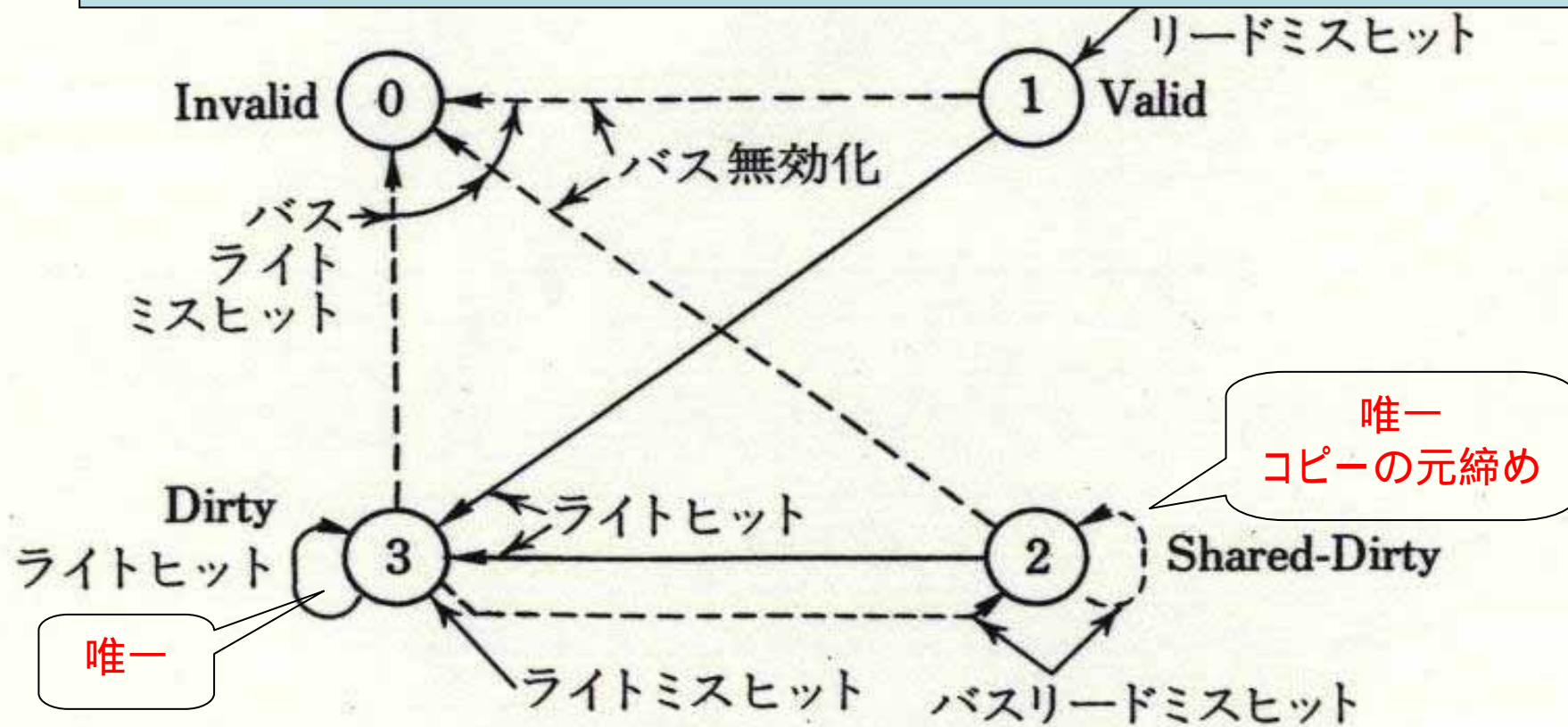


図 4.24 スヌープキャッシュの基本方式

表 4.3 スターフキャッシュ方式

共有ブロックへの書込み時処理 キャッシュ間転送時の主記憶更新	ブロードキャスト無効化	ブロードキャスト更新
変更のあるブロックのキャッシュ間転送時に主記憶へは書き戻さない	Berkeley State 0: Invalid State 1: Valid (clean, potentially shared, unowned) State 2: Shared-Dirty (modified, potentially shared, owned) State 3: Dirty (modified, only copy, owned)	Dragon State 0: Valid-Exclusive (clean, only copy) State 1: Shared-Clean (clean, one or more copy) State 2: Shared-Dirty (modified, one or more copy) State 3: Dirty (modified, only copy)
変更のあるブロックのキャッシュ間転送時に主記憶へも書き戻す	Illinois State 0: Invalid State 1: Valid-Exclusive (clean, only copy) State 2: Shared (clean, possibly other copies) State 3: Dirty (modified, only copy)	Firefly State 0: Valid-Exclusive (clean, only copy) State 1: Shared (clean) State 2: Dirty (dirty, only copy)

プロセッサA	状態	プロセッサB	状態	プロセッサC	状態
LOAD	Valid				
STORE	Dirty				
	Shared	Dirty			
	InValid				
		LOAD	Valid	LOAD	Valid
		STORE	Dirty		InValid



(a) バークレイプロトコル

MOESI ADM64で採用

- ・Modified Dirty 唯一、メモリ内容と不一致
- ・Owned 唯一、メモリ内容と不一致、コピーを持つのはShared
- ・Exclusive、Clean、唯一
- ・Shared 他にOwnedがあればメモリ内容と不一致、なければメモリ内容と一致
- ・Invalid 無効

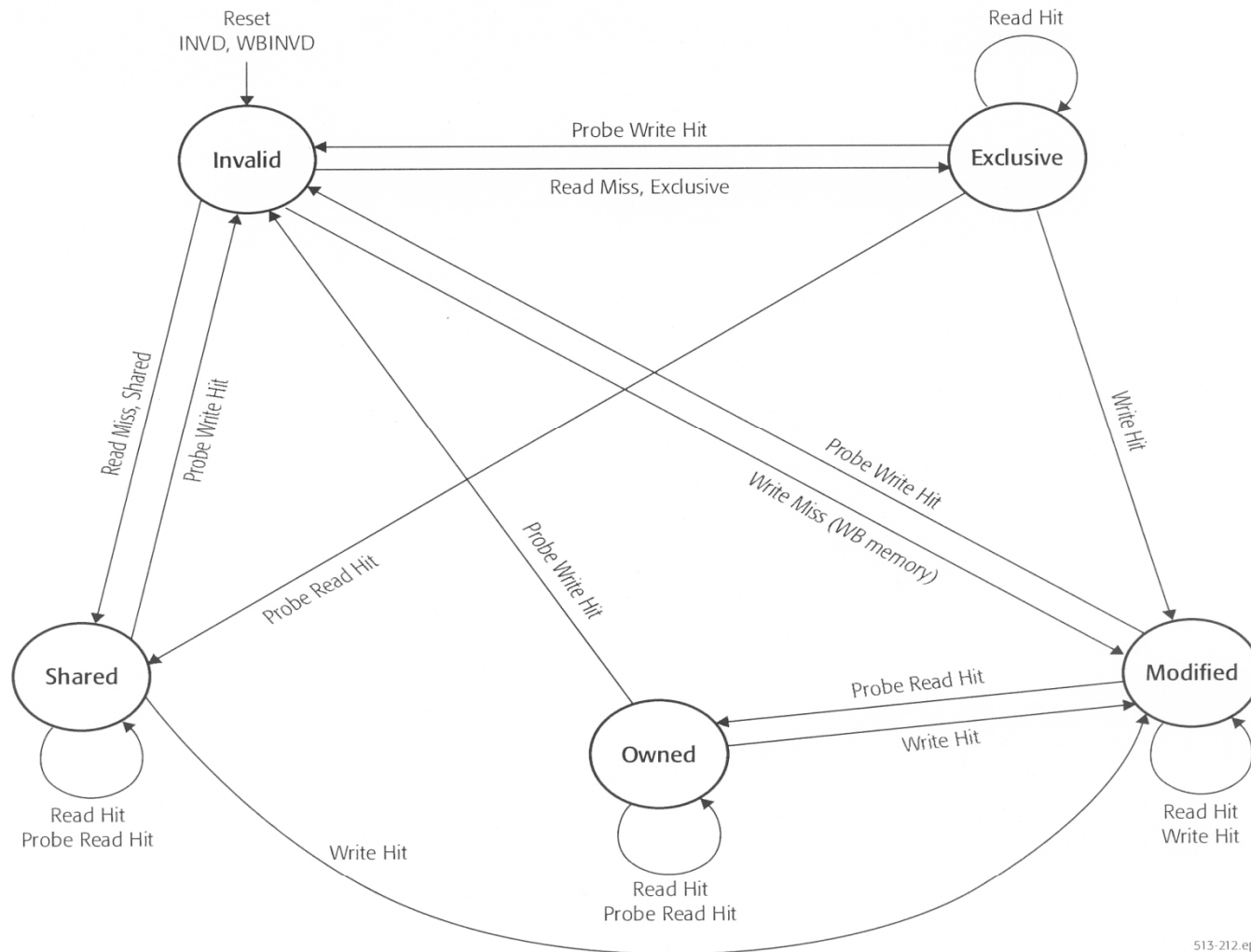


Figure 7-2. MOESI State Transitions

To maintain memory coherency, external bus masters (typically other processors with their own internal caches) need to acquire the most recent copy of data before caching it internally. That copy can be in main memory or in the internal caches of other bus-mastering devices. When an external master has a cache read-miss or write-miss, it *probes* the other mastering devices to determine whether the most recent copy of data is held in any of their caches. If one of the other mastering devices holds the most recent copy, it provides it to the requesting device. Otherwise, the most recent copy is provided by main memory.

表1:キャッシュラインの状態遷移の例

	事象	P0	P1	P2
0	初期状態	Invalid	Invalid	Invalid
1	P0 Read	Exclusive	Invalid	Invalid
2	P1 Read;P0応答データ供給	Shared	Shared	Invalid
3	P2 Read;P0応答データ供給	Shared	Shared	Shared
4	P0 Write;Invalidate応答を待つWrite	Modified	Invalid	Invalid
5	P2 Read;P0応答データ供給	Owned	Invalid	Shared
6	P1 Write;Invalidate、P0(Owner)応答データ供給、P2 Invalidate応答	Invalid	Modified	Invalid
7	P2 Read;P1応答データ供給	Invalid	Owned	Shared
8	P1 Writeback	Invalid	---	Shared

4.4.2具体例

4.4.3ディレクトリ方式

一般のネットワーク利用

メモリ側で集中管理

フルマップ

リミテッド

チェーン

一般のネットワーク

Snoop方式はX

放送機能が弱い

排他制御が困難

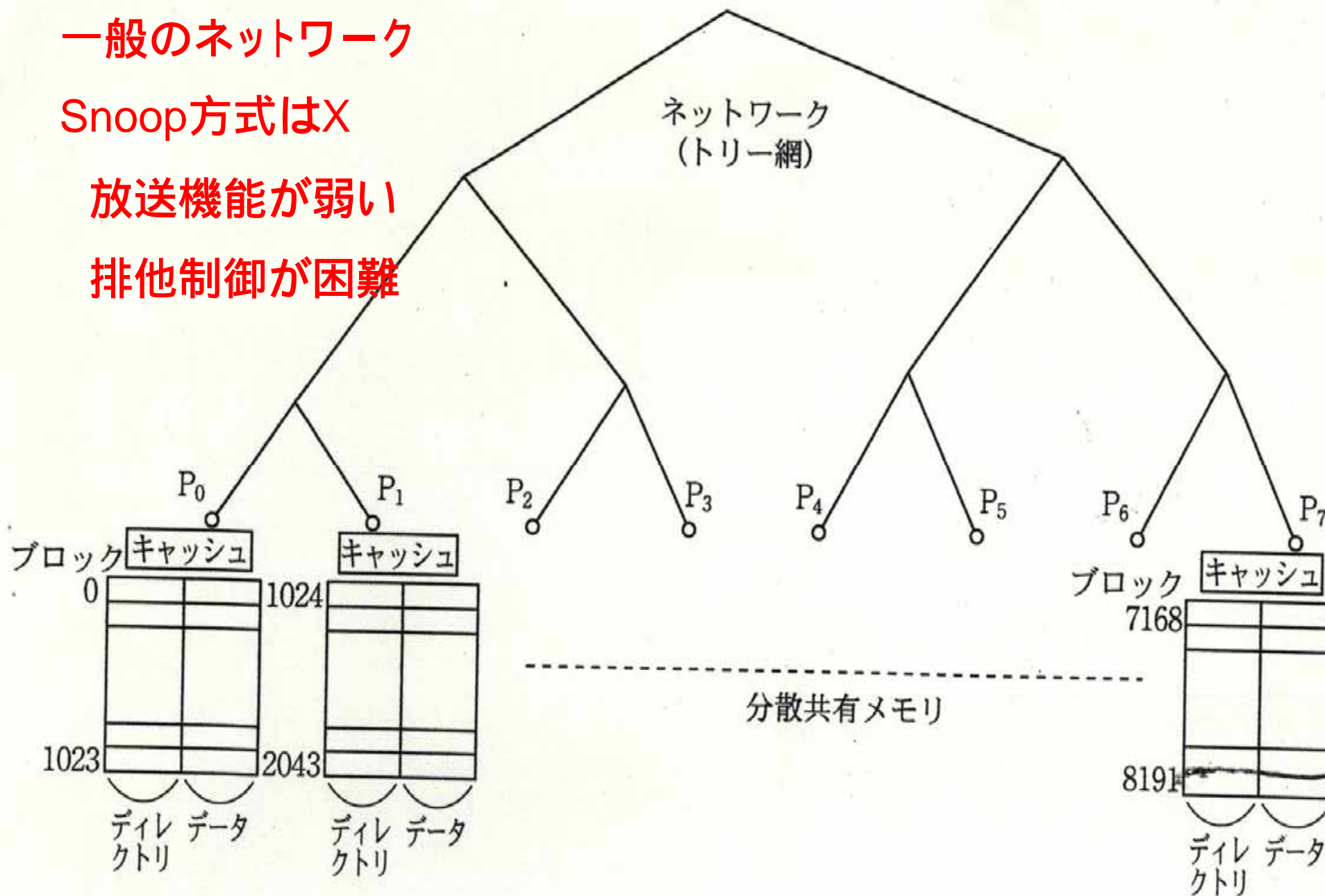


図 4.26 ディレクトリ方式

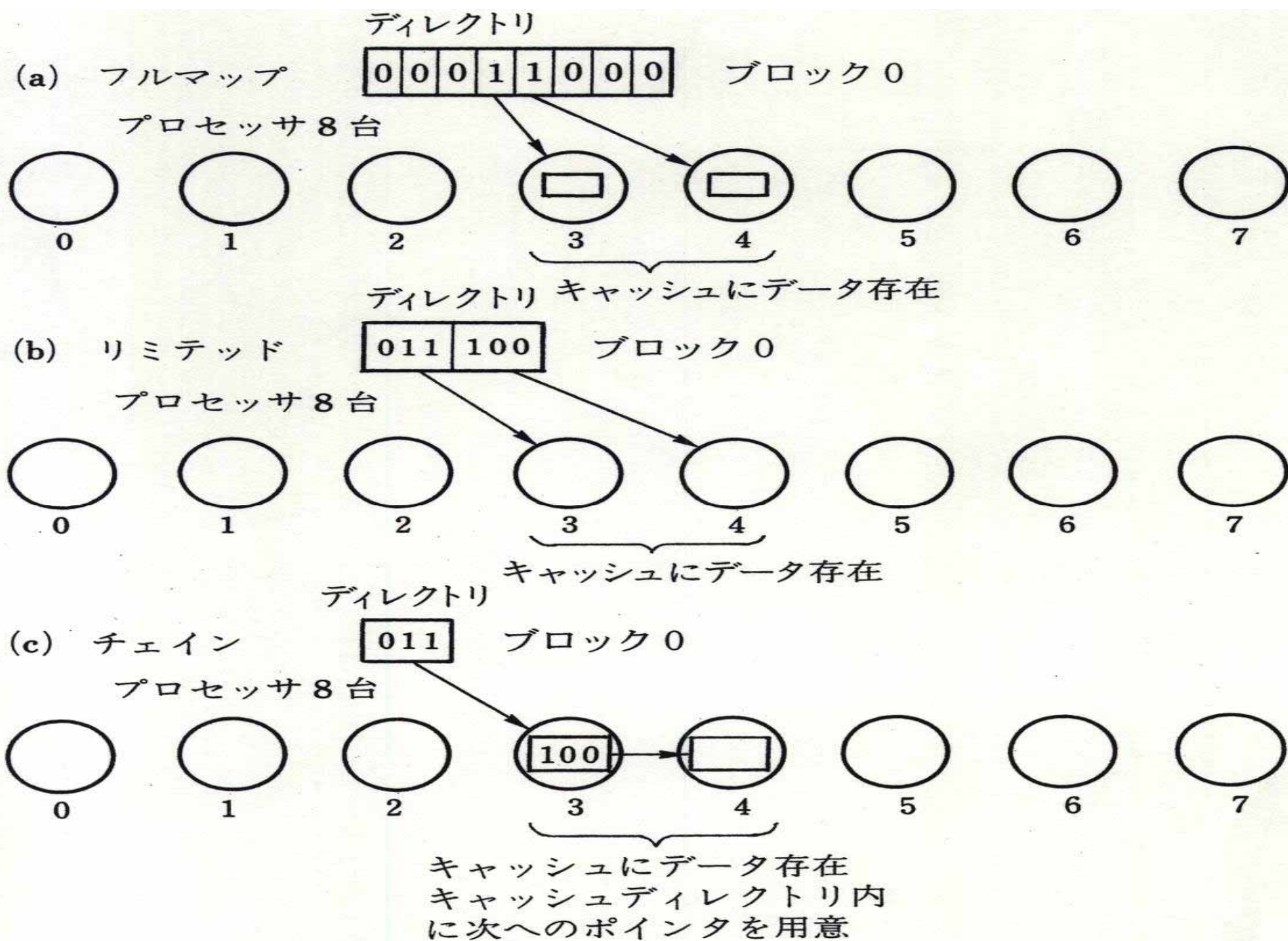
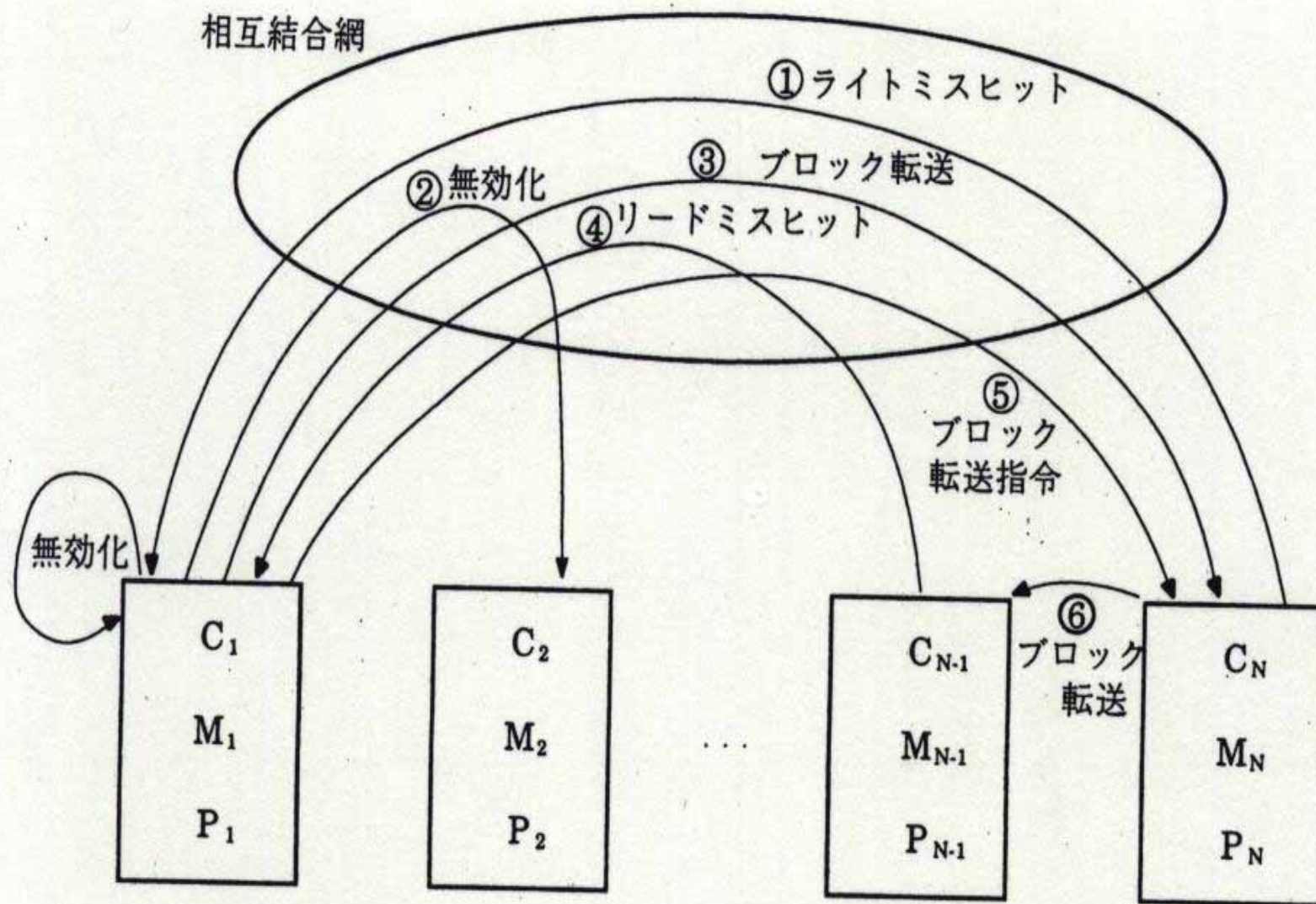
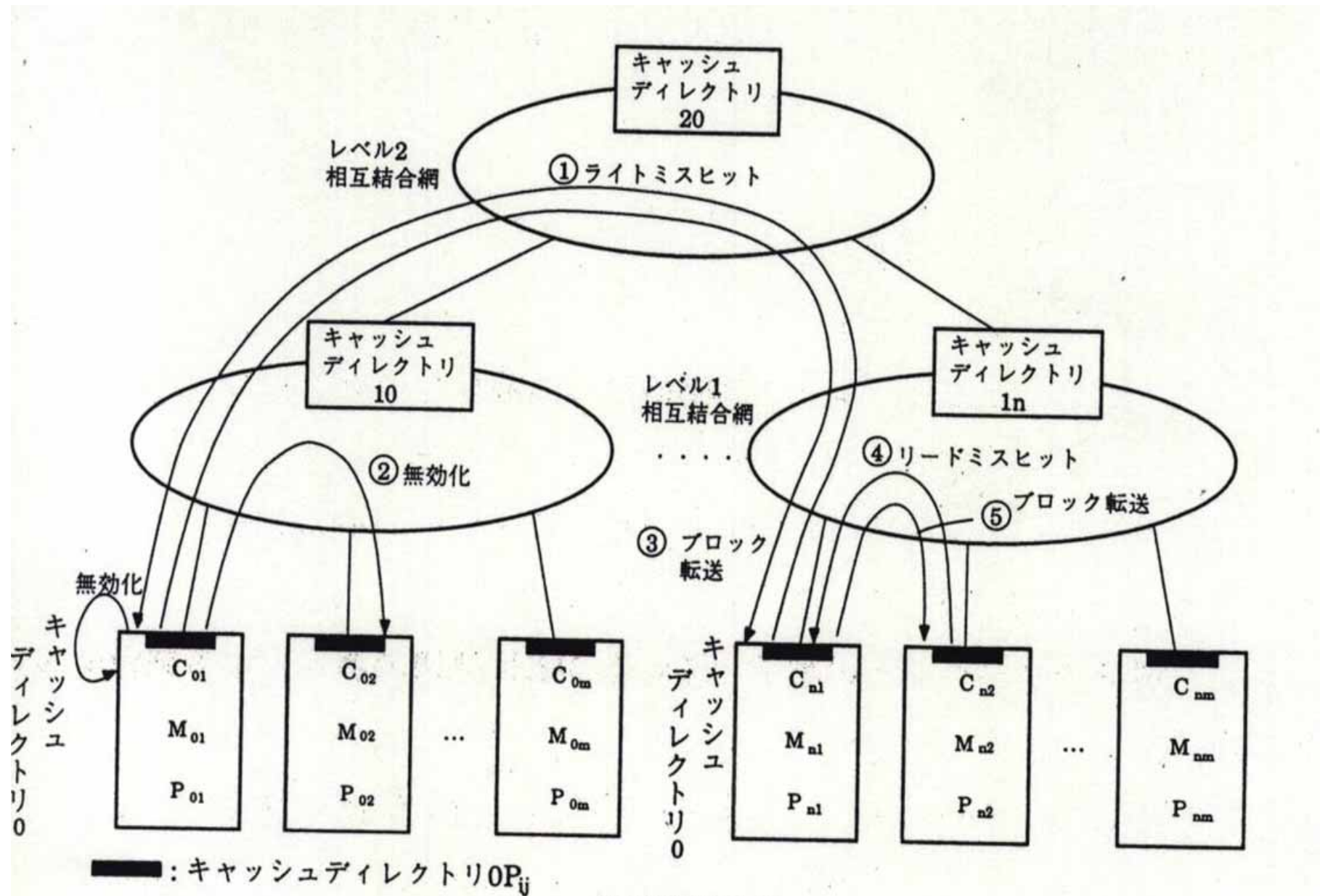


図 4.27 ディレクトリ方式



(a) 通常のディレクトリ方式



(b) COMA方式

図 5.11 キャッシュオンリメモリ方式

COMA方式

コヒーレンスの動作を応用例で見ると
線形1次方程式の反復解法

$$A X = a \quad X = b + B X$$

$$X = D^{-1} ((D - A)X + a) : \text{ヤコビ法}$$

直接法: ガウス消去法、LU分解法

スヌープキャッシュの2つのプロセッサで並列実行

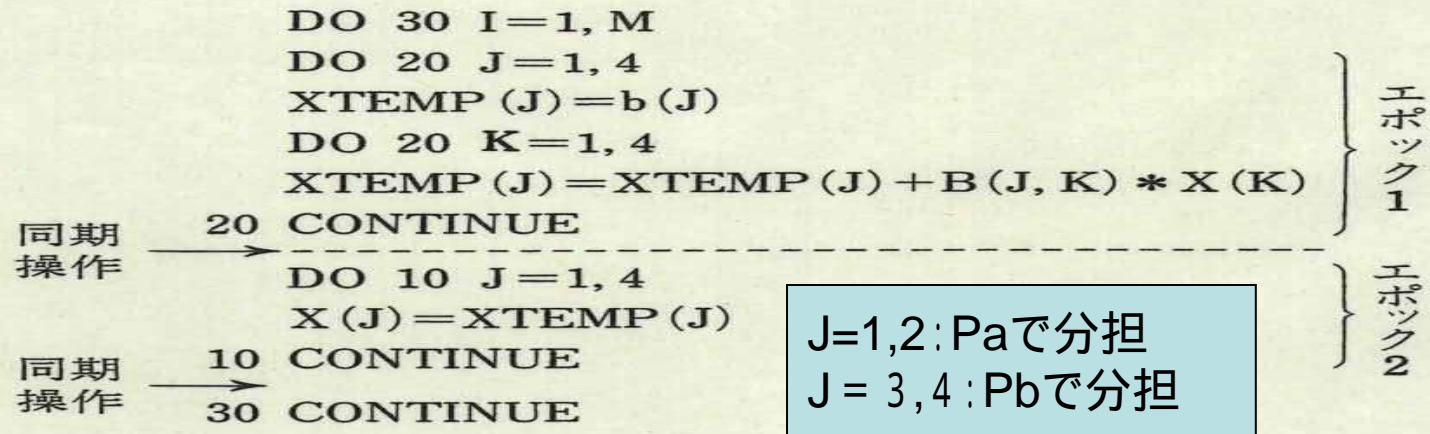
同期命令が必要: 先走り禁止

新しい値 / 古い値を使ってしまう

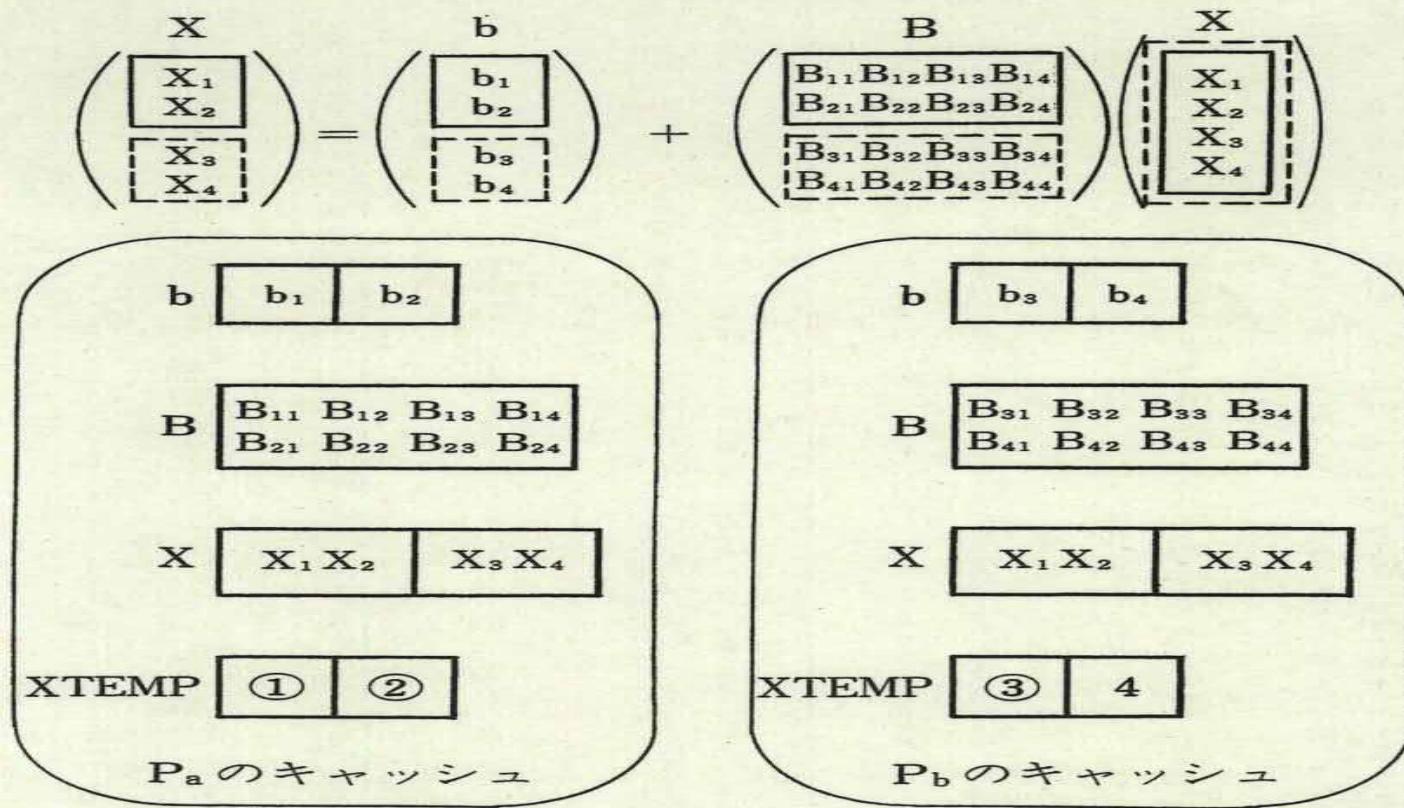
無効化 / 更新どちらがよいか

この場合は更新

コヒーレンス制御が必要のないものがある



(a) 線形方程式の反復解法



(b) キャッシュの状態

図 4.28 ソフトウェアによるキャッシュコヒーレンス制御方式

4.4.5ソフトウェアによる方式

コンパイラで共有データを検出

リードオンリ、排他的専有利用データ:対象外
ストアスルーを前提:最新データはメモリに存在、ハード簡単
ソフト的な選択的無効化

マーキング法:共有データにマークを付け、キャッシュしない

キャッシュ制御命令:共有データブロックをひとつ

ひとつ命令で無効化

ハードウェアによるキャッシュ全面的無効化

同期命令でキャッシュ全面無効化

リードオンリデータはCache Read命令で復活、有効化

共有データはMemory Readでミスヒット、

メモリからキャッシュへ転送、その後ヒット

ハードウェアによる選択的無効化機構の設置

キャッシュブロックごとに無効化ビットを持たせ、一斉に選択
的に無効化

5.3.5 通信量の削減法

(1) データ属性に基づいたプロトコル切替え

読出し専用 (リードオンリ) データ

専有データ

無効化が有利なデータ

更新が有利なデータ

最終書込み時ブロードキャストが

有利なデータ

バリア同期変数

ロック変数

不可分操作

T S (Test and Set) 命令

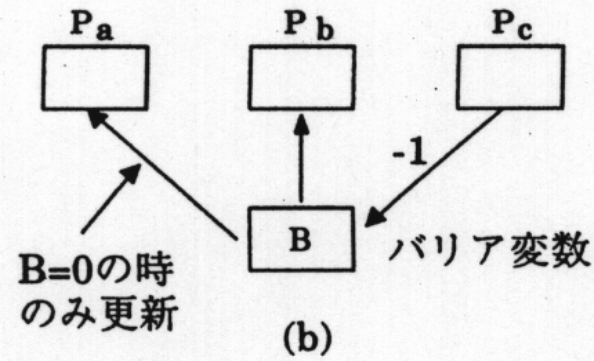
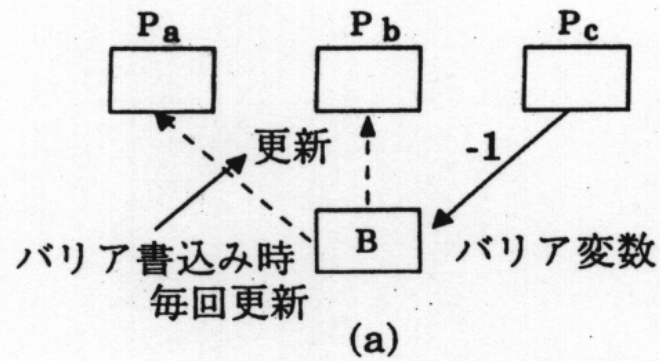
原始的な方式

T S 命令で行う方式

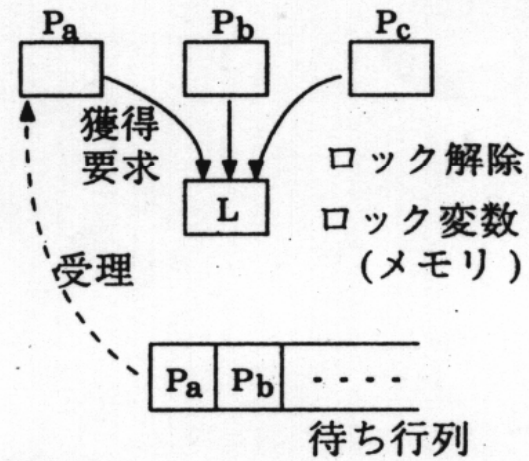
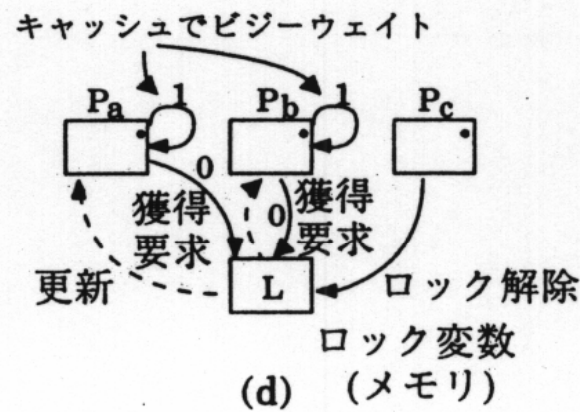
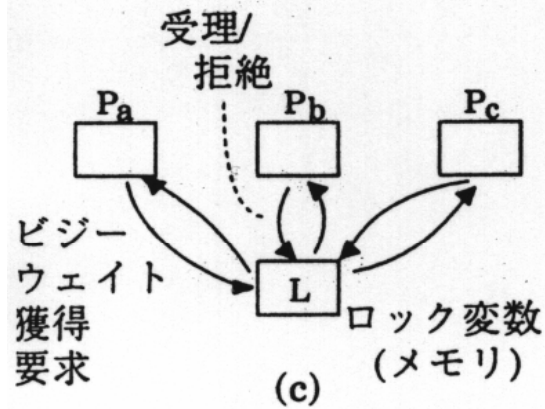
ビジーウエイト

少し効率のよい方式

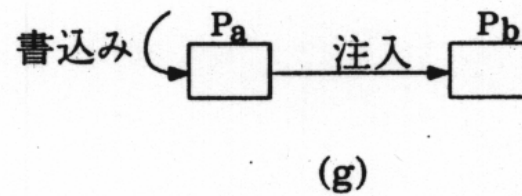
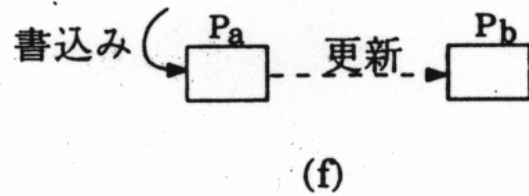
更新方式



(1) バリア同期変数



(2) ロック変数



(3) 通信データ

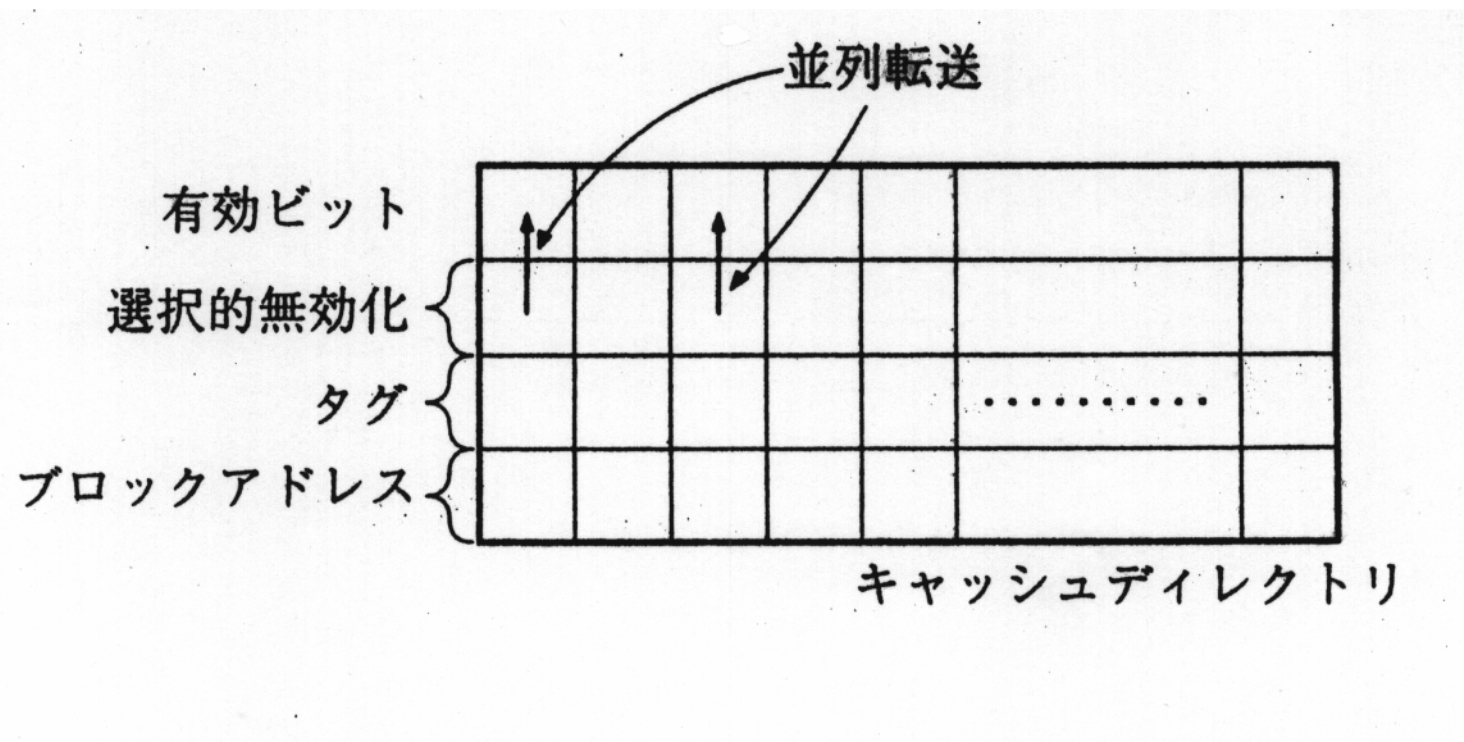
通常のロード命令でキャッシュからロック変数読み出し

0 であれば、TS 命令によるロック獲得
要求の待ち行列

通信データ

注入 (injection)

(2) キャッシュ管理の積極的導入



5.3.6 キャッシュ

オンリメモリ

プログラムやデータ：

システム内のどこかのキャッシュに分散格納

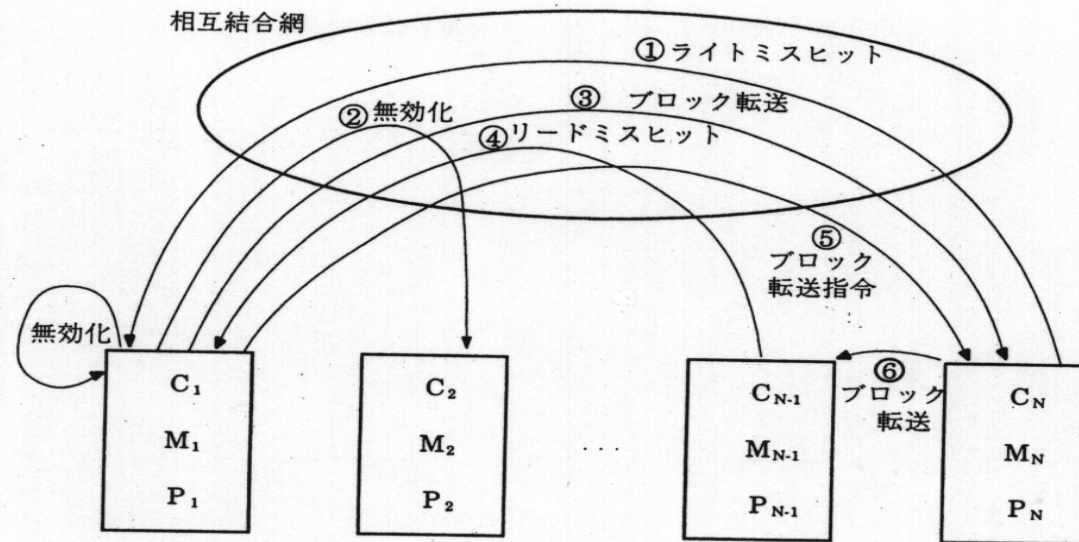
階層構造

キャッシュディレクトリのハードウェア構成：

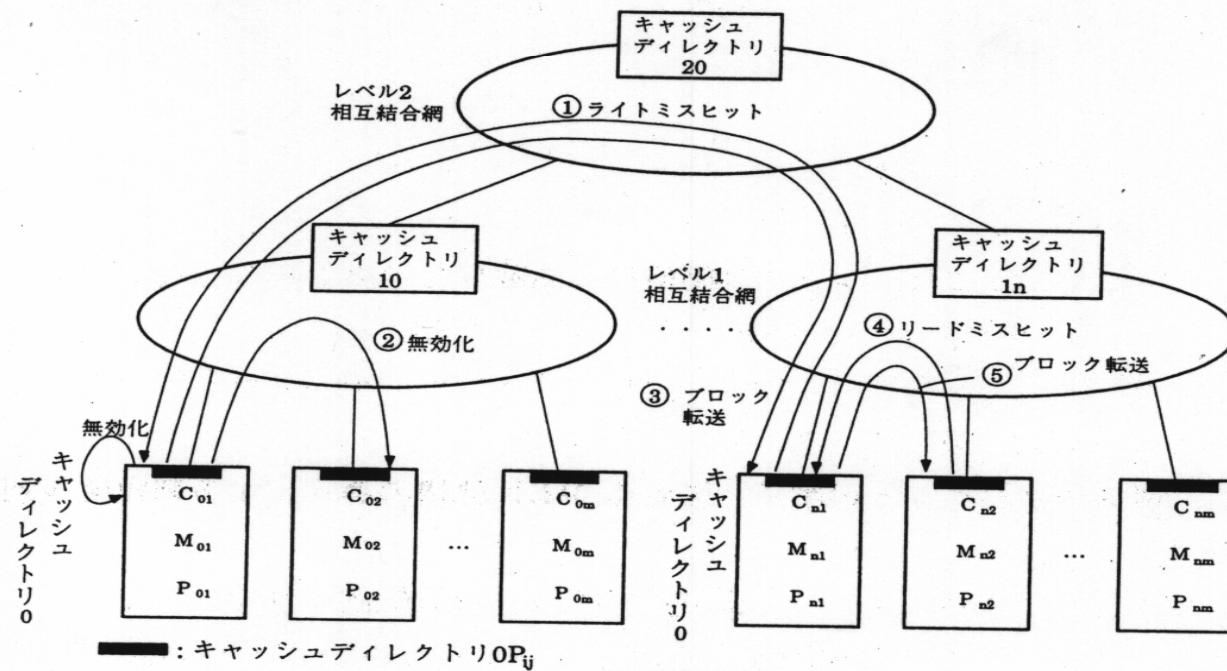
連想メモリで構成

上位包含性 (multi-level inclusion)

データ参照の局所性



(a) 通常のディレクトリ方式



(b) COMA方式

図 5.11 キャッシュオンリメモリ方式

5 . 4 J U M P - 1

JUMP-1、スタンフォード大学のDASH、
MITのAlewifeなど

JUMP-1

クラスタ構成：1 2 8 クラスタ

クラスタ：

4 台のSuper SPARC Plus（動作周波数は6 0 MHz）

2 台のメモリ（各6 4 MB）、

2 台の知的メモリ機構（MBP）

数本の高速シリアルリンク（STAFF-Link）

プロセッサ：各 1 MBの 2 次キャッシュ付き
階層トーラス網（以後、RDT網とよぶ）

5.4.1 特長

- （ 1 ） キャッシュコヒーレンス制御機構を
内蔵したメモリ共有方式
- （ 2 ） 通信オーバーヘッドの削減を図れる
MBP機構の導入による細粒度並列処理方式
レイテンシ削減
無駄な通信の削減
- （ 3 ） スケーラブルで直径の小さい階層トーラス網

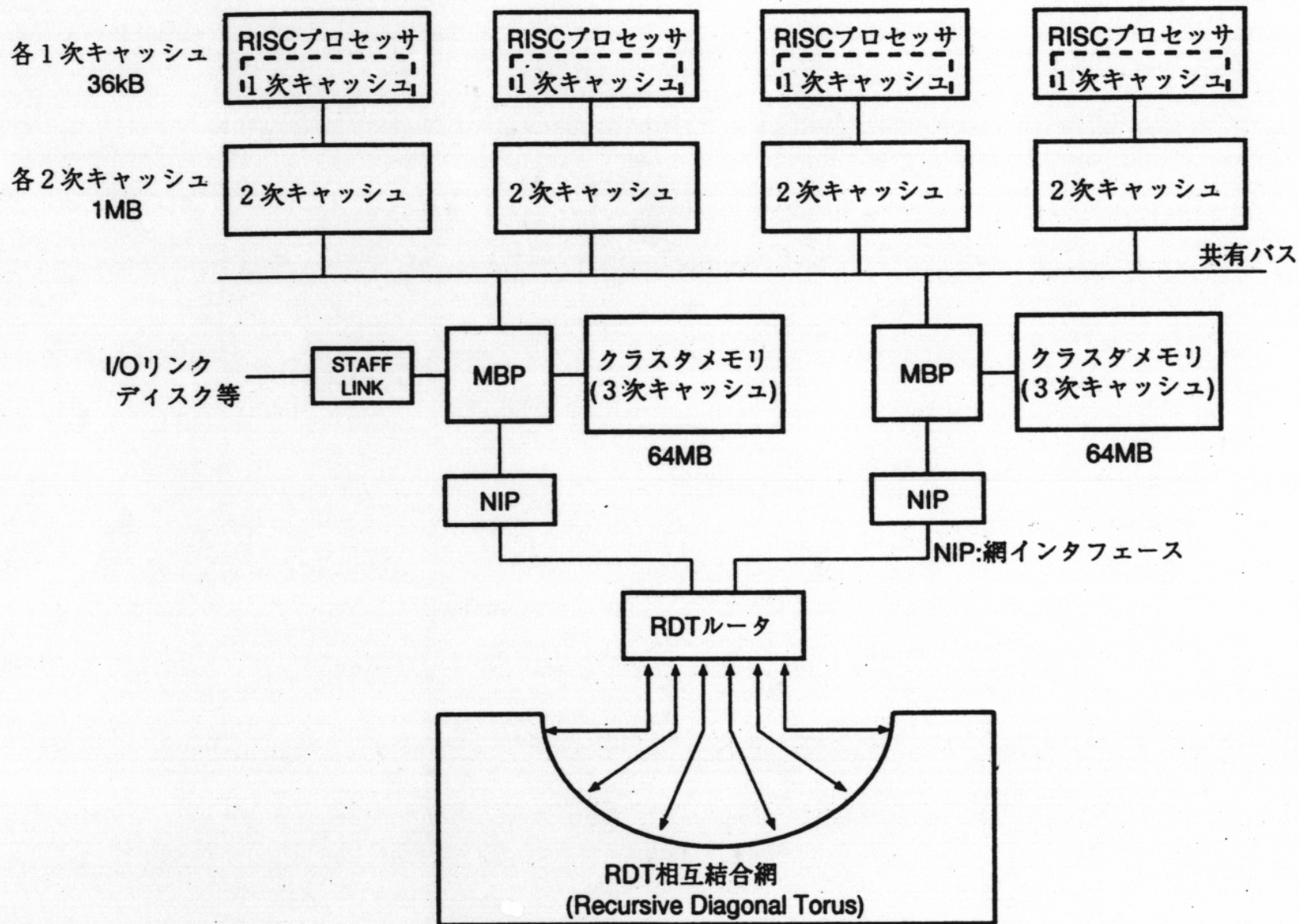


図 5.12 JUMP-1 のクラスタ内構成

(4) 高速な入出力機構

5.4.2 メモリ

アーキテクチャ

(1) 仮想共有メモリ

クラスタメモリ :

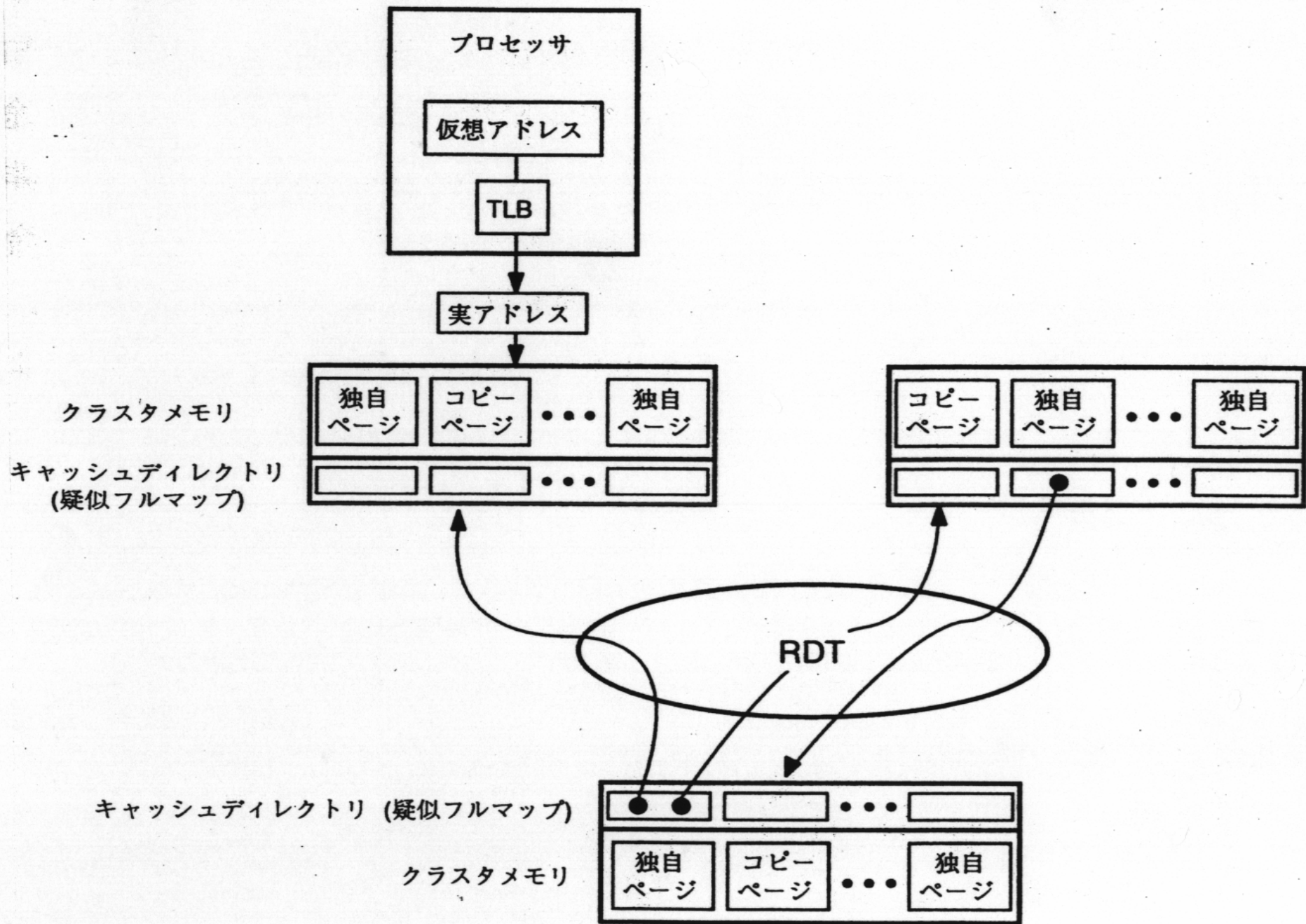
独自 (original) 空間とコピー (copy) 空間

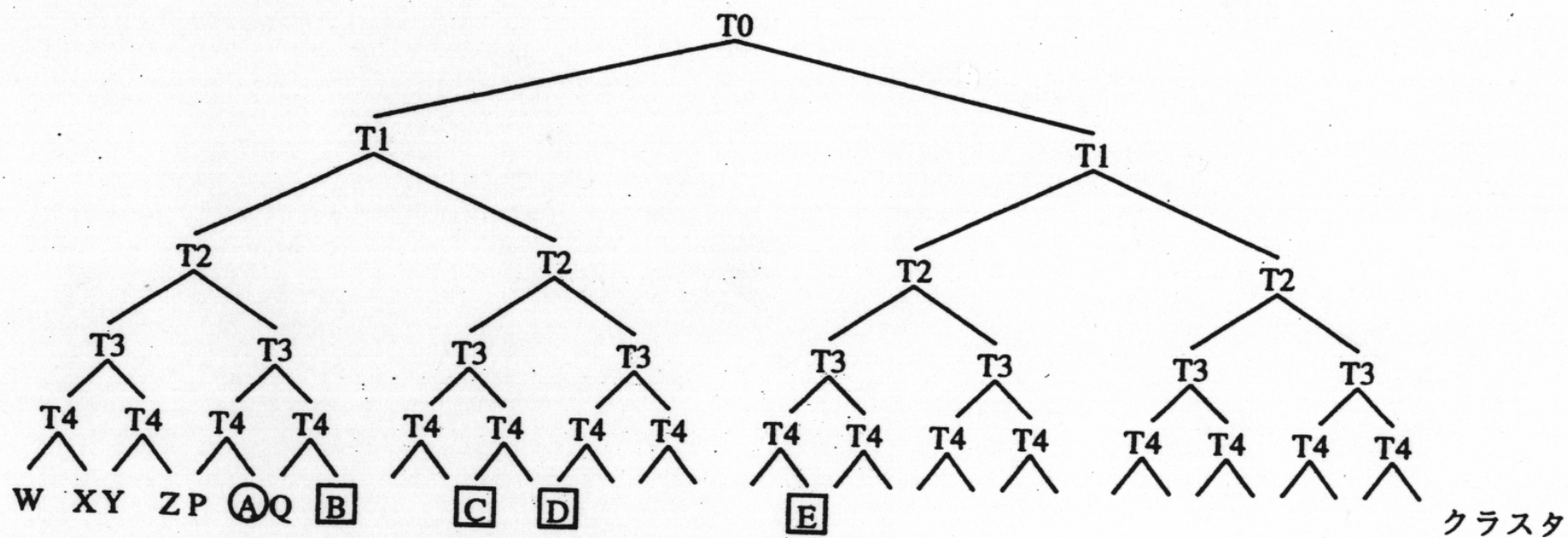
3 次キャッシュ

疑似フルマップ方式

(2) クラスタ間キャッシュコヒーレンス

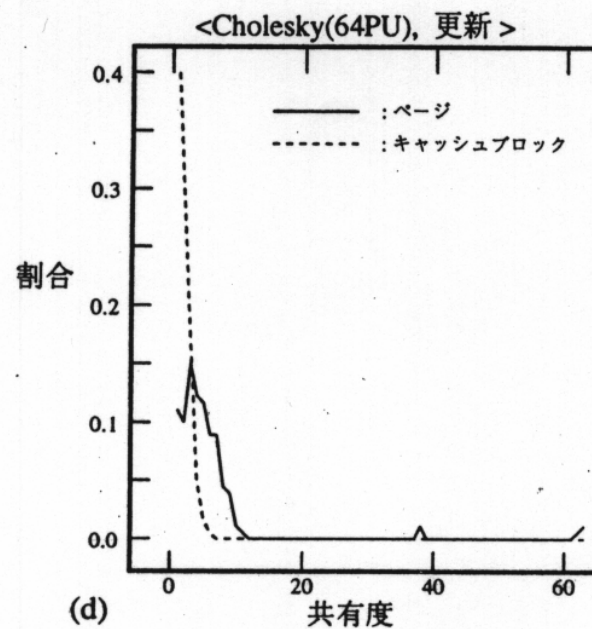
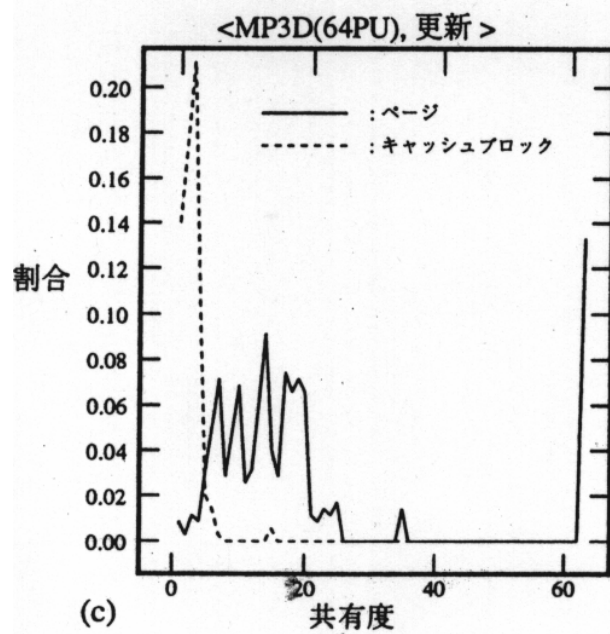
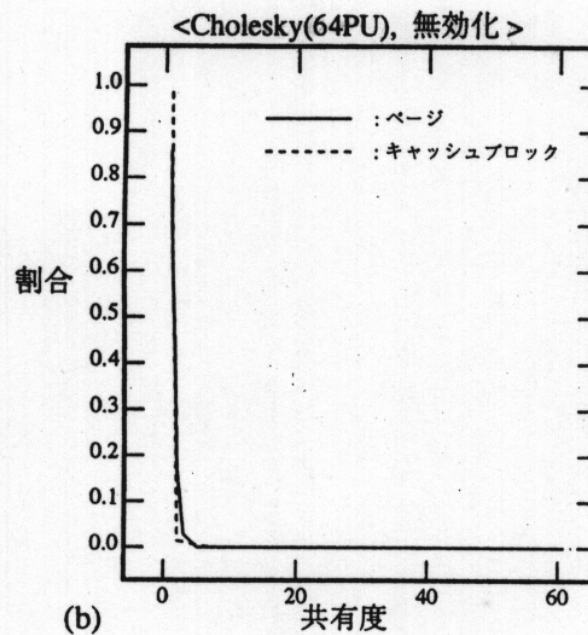
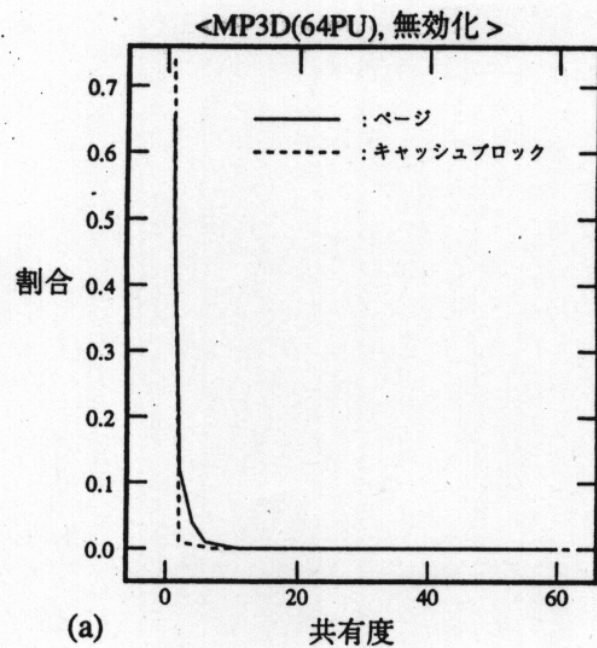
(a) 疑似フルマップ方式

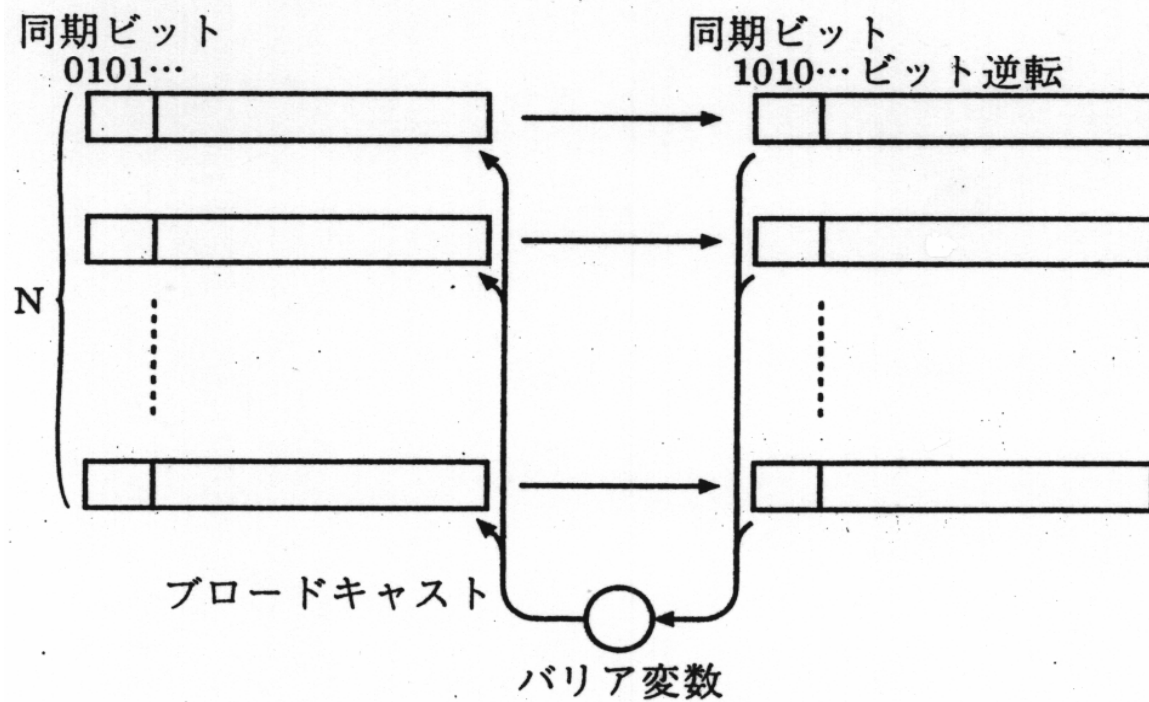
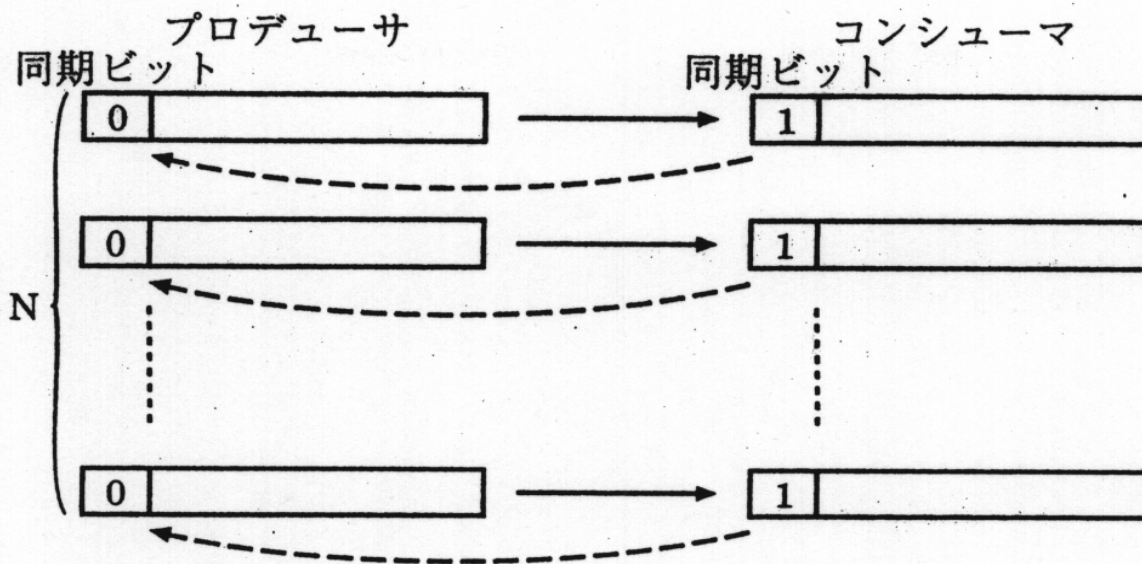




	T4	T3	T2	T1	T0	MD
A->B	0:0	0:1	0:0	0:0	0:0	0
A->B,C,D	0:0	0:1	0:0	0:1	0:0	0
A->B,E	0:0	0:1	0:0	0:0	0:1	0
A->ブロードキャスト	1:0	0:1	1:0	0:1	0:1	0
A->E	0:1	1:0	1:0	1:0	0:1	1

図 5.14 疑似フルマップ [188]





放送機構

結合操作 (Combining)

(b) 疑似フルマップ方式の評価

- ・ 共有ブロックの大きさ
- ・ 書込み時のポリシ

リミテッド、チェーン方式と比較

- ・ ディレクトリ容量
- ・ ネットワークを流れるトラフィック量

5.4.3 高速同期 ・ 通信機構

(1) キャッシュ

タグ：INV(Invalid),EX-Dty(Exclusive Dirty),
LS-CIn(Locally Clean),
LS-Dty(Locally Dirty),
GS-CIn(Globally Clean)の5状態

Iストラクチャのfull/emptyビット
(8バイトごと)

2次キャッシュの特長

- ・ ページ属性として無効化と更新
データ属性に応じたプロトコル選択
- ・ 実アドレスの上位5ビットをコマンドとして

使用

プリフェッチと注入

注入の例：FIFOキューの先頭が消費された
場合

リードオール

更新データのマージ

早期共有解除

自浄機構

(2) キャッシュコヒーレンスプロトコルの動的切り
替え

(3) I - ストラクチャとFIFO機構

1 対 1 通信

I - ストラクチャ

異なるプロセッサ間 (N 個) で 1 対 1 通信

N 個のコンシューマへのデータ到着

バリア同期機構で高速に検知

各プロデューサに放送

ビット逆転方式 (最初のデータ転送では

0、次の転送では 1、その次では 0、

でデータの利用可能状態を示す)

1 対多通信

多対 1 通信

FIFOキュー

多対多通信

(4) バリア機構

5.4.3 ネットワーク

アーキテクチャ

ブロードキャスト

ランク 0 グループ

グループでランク 1 グループ

8 進木構成

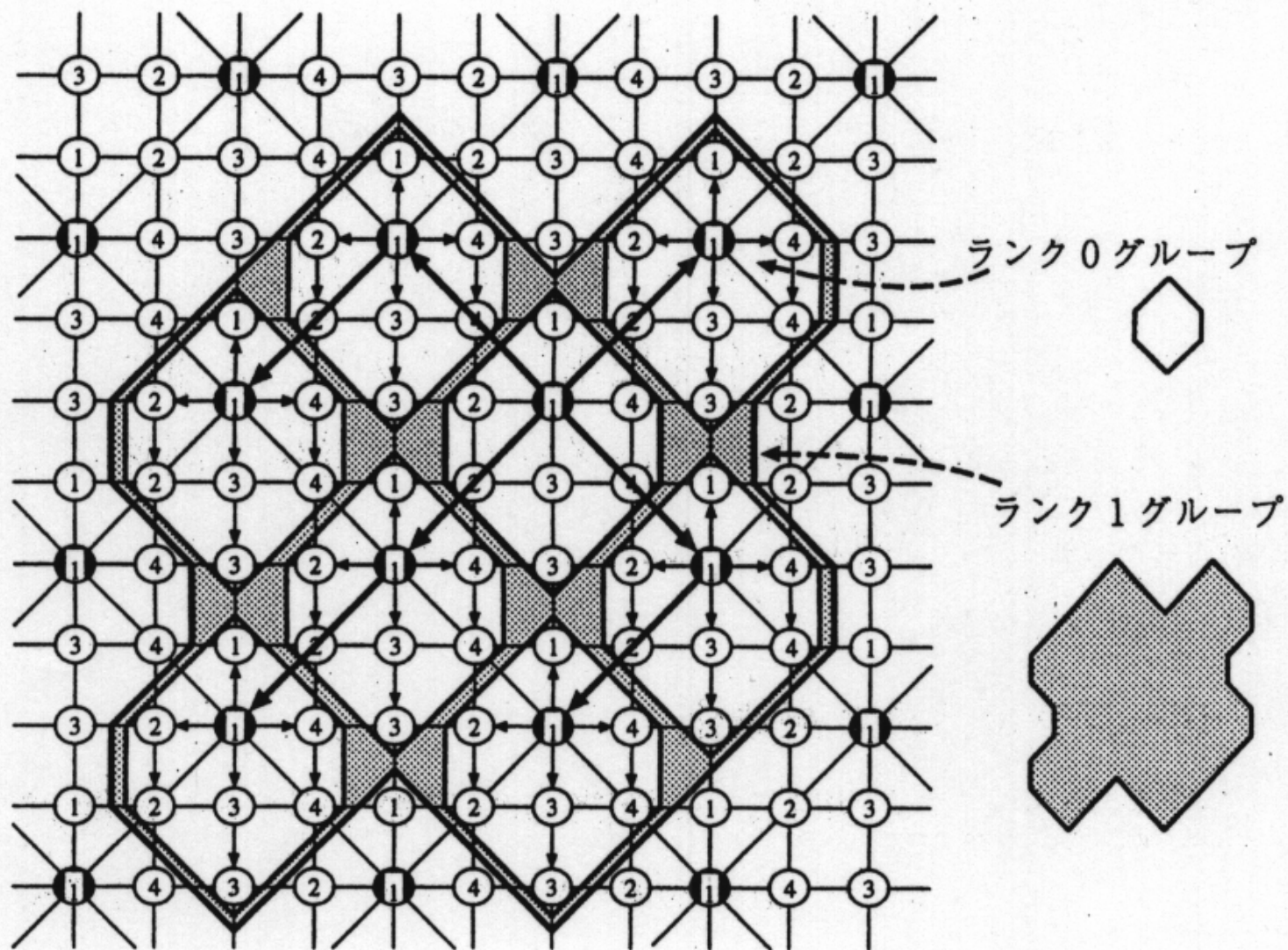
ブロードキャスト、バリア同期や疑似フルマップの制御

ルーティング：ベクトルルーティング

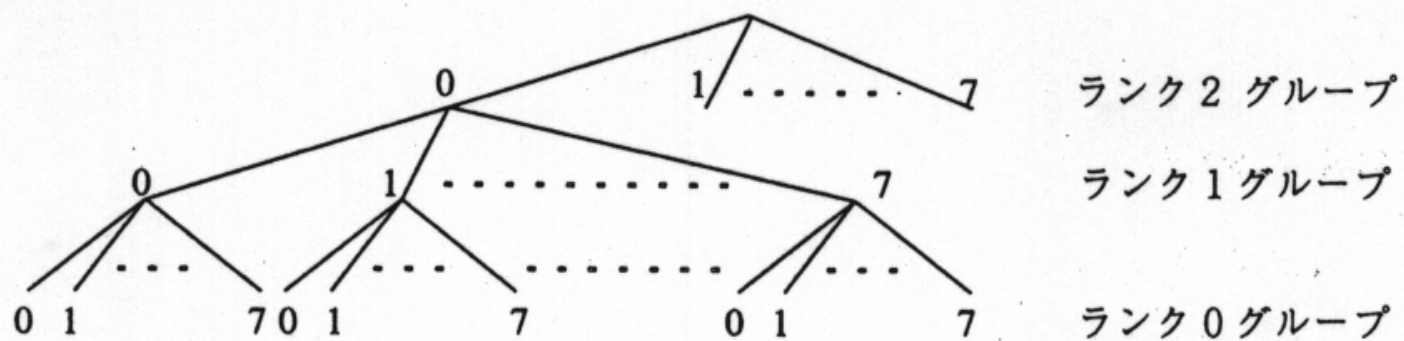
RDTの性能：

3次元トーラスとハイパキューブの中間的な特性

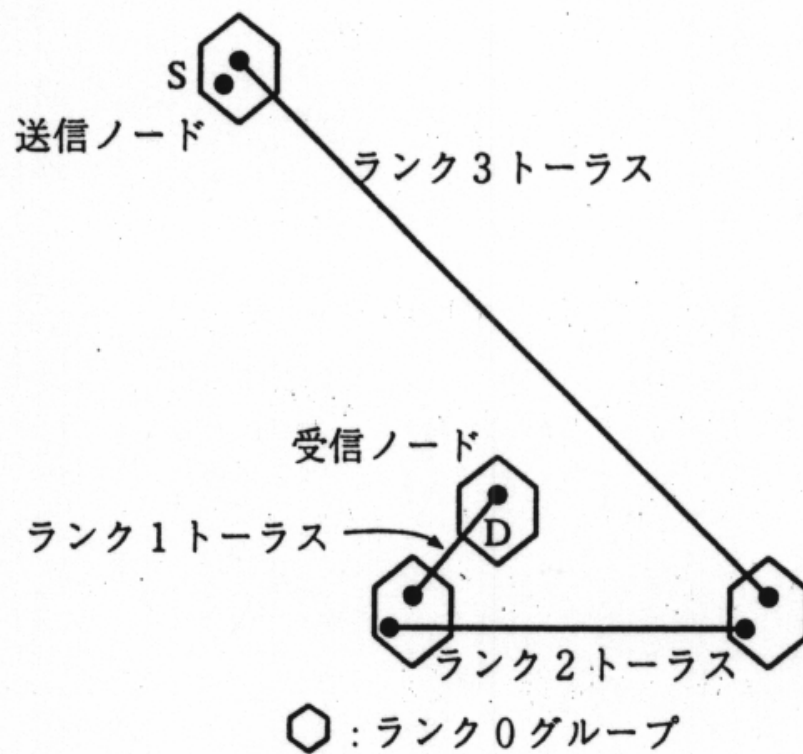
ハイパキューブ網やシャフル網の高速にシミュレート



(a) RDTの構造



(b) 8進木構成



(c) ベクトルルーティング

5.5 メッセージ交換型マルチプロセッサ

OS

1 OSの必要性

複雑な機能の提供: 命令セットの拡張

大容量記憶, ファイル操作,

ネットワーク, 入出力

リソースの共有利用: 効率的な管理,

セキュリティ

多数のユーザプロセスの実行

多数のサーバプロセスの実行

2 OSの機能

記憶管理

多重仮想記憶

ファイル管理

ディレクトリ管理, ファイル更新・編集

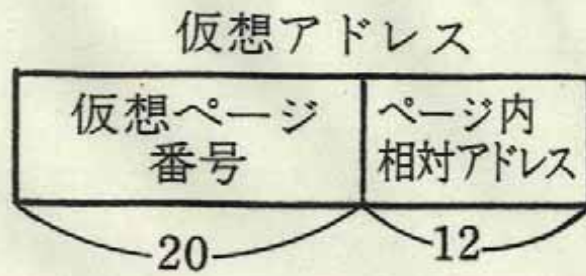
効率のよい記憶方式

入出力管理

デバイスドライバ

スケジューリング

記憶管理



1次元アドレス

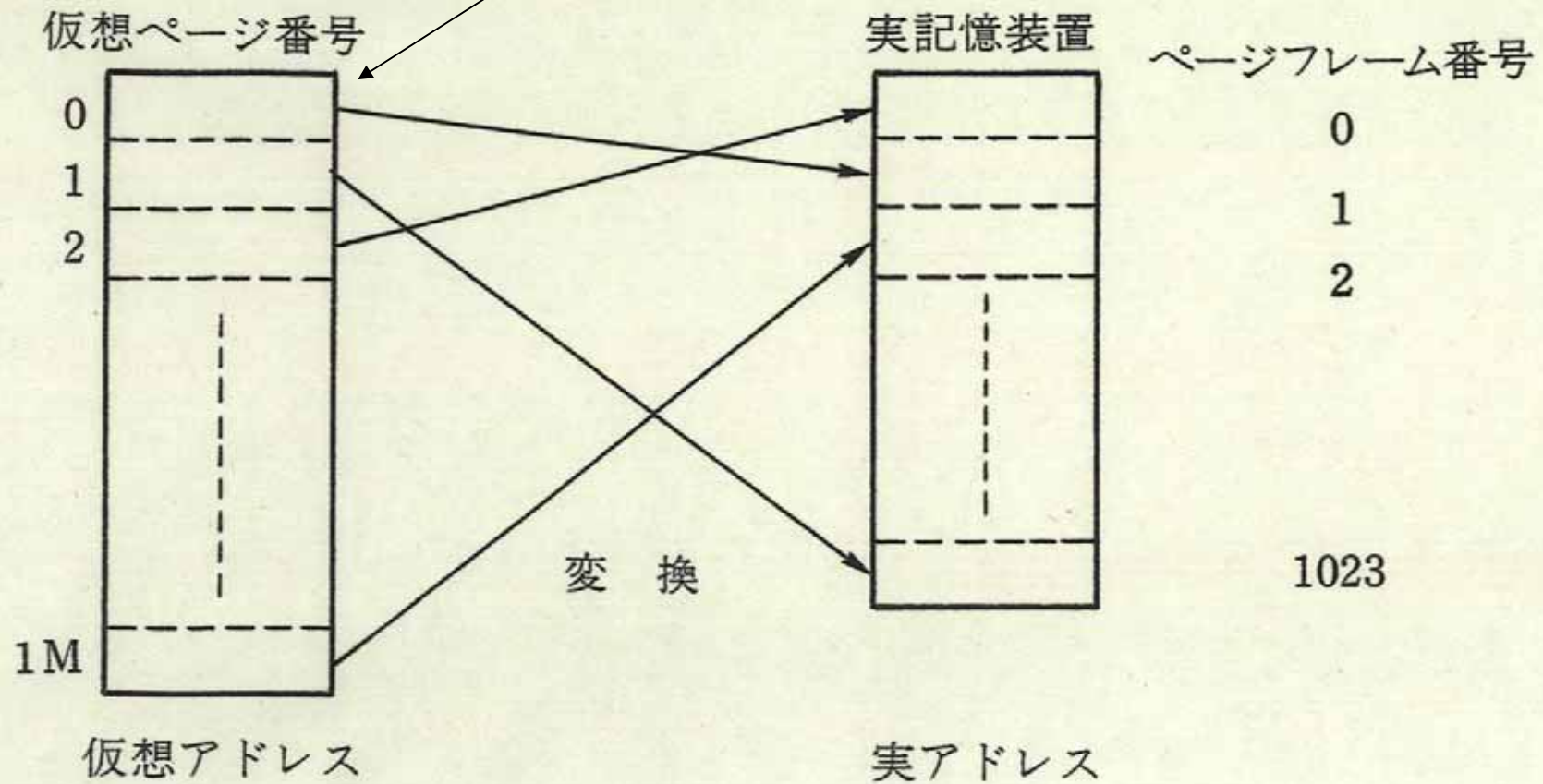
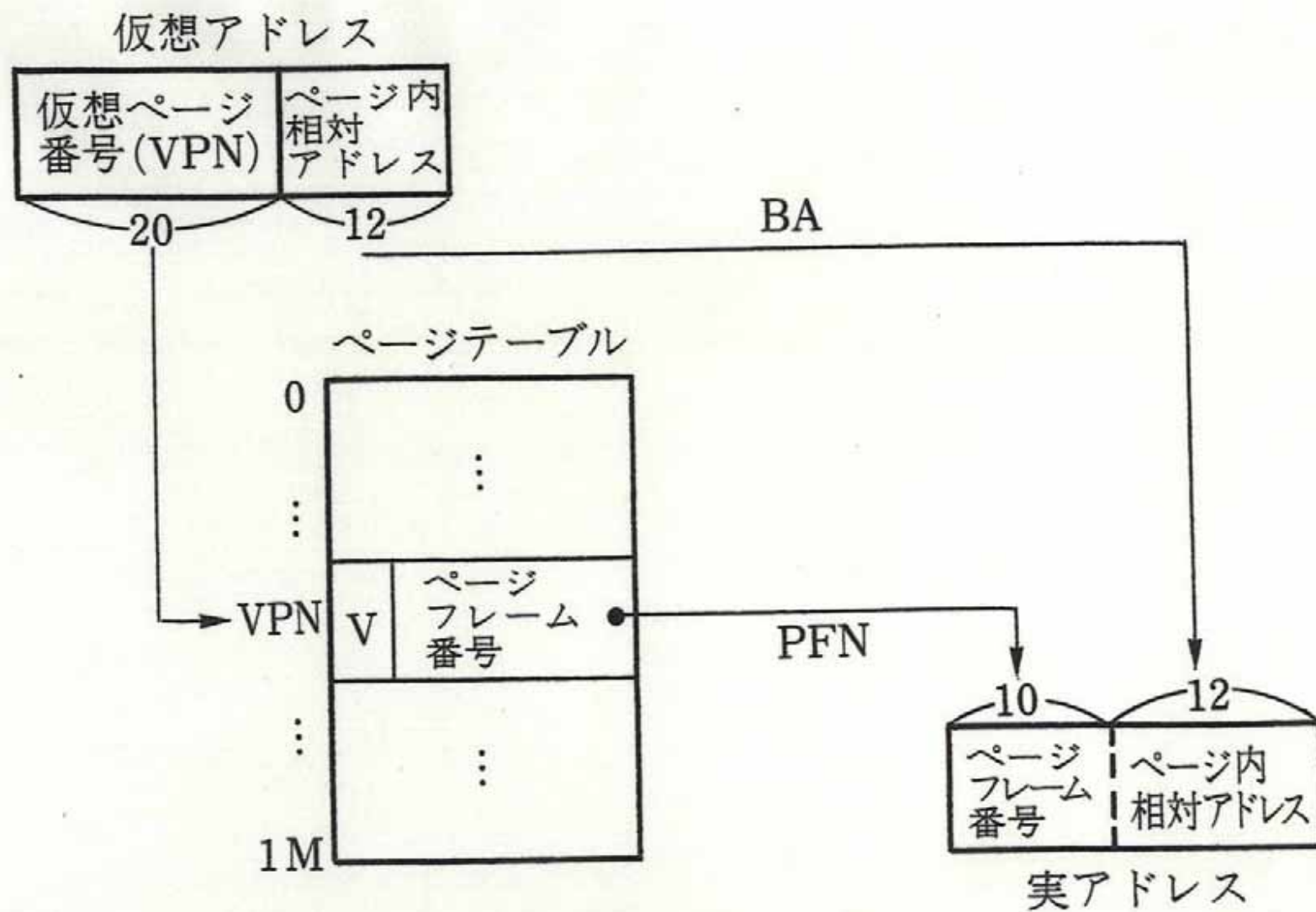


図 4.6 ページング方式



(a) ページテーブル

図 4.8 直接写像方式

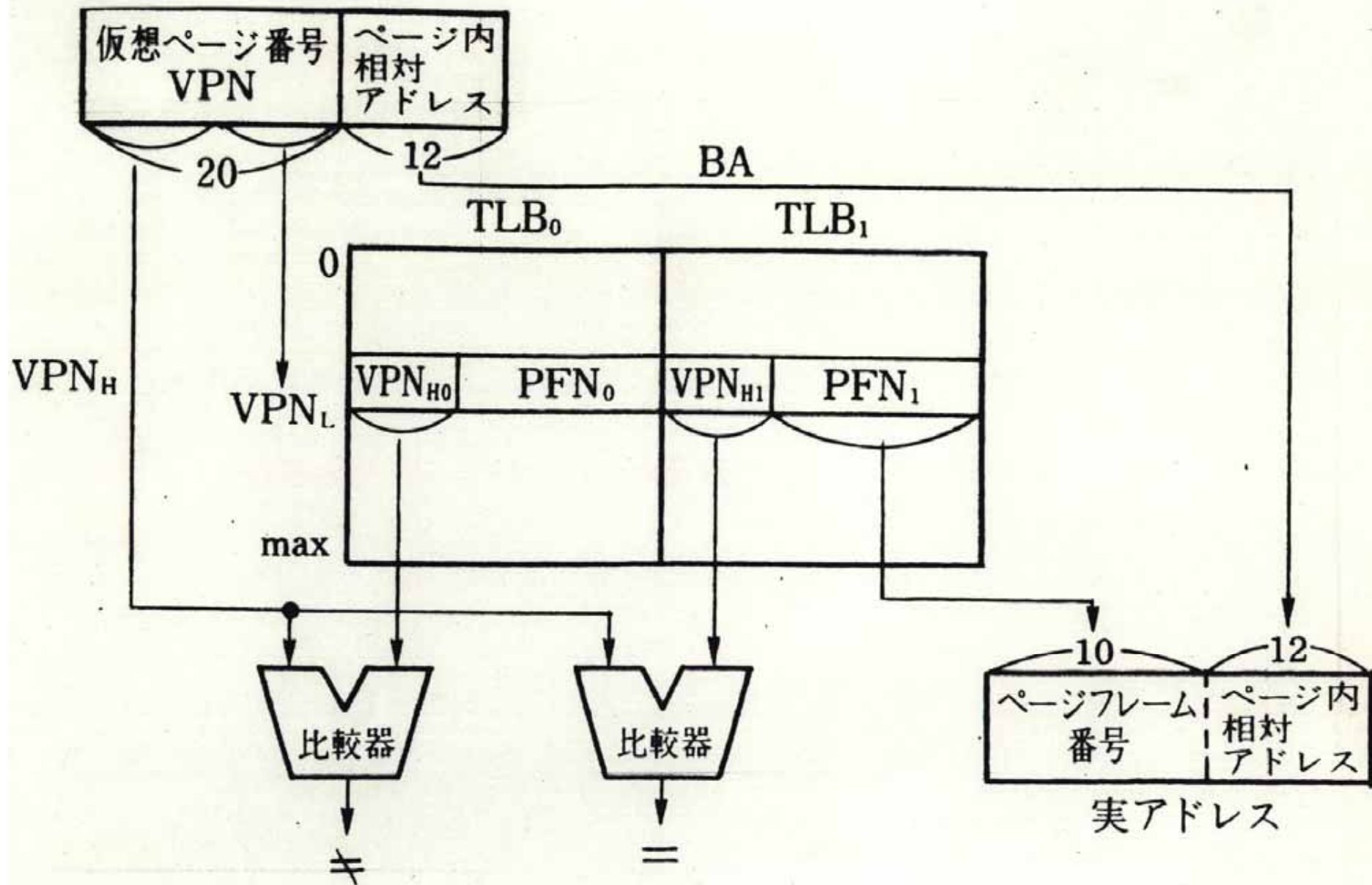


図 4.11 TLB の構成

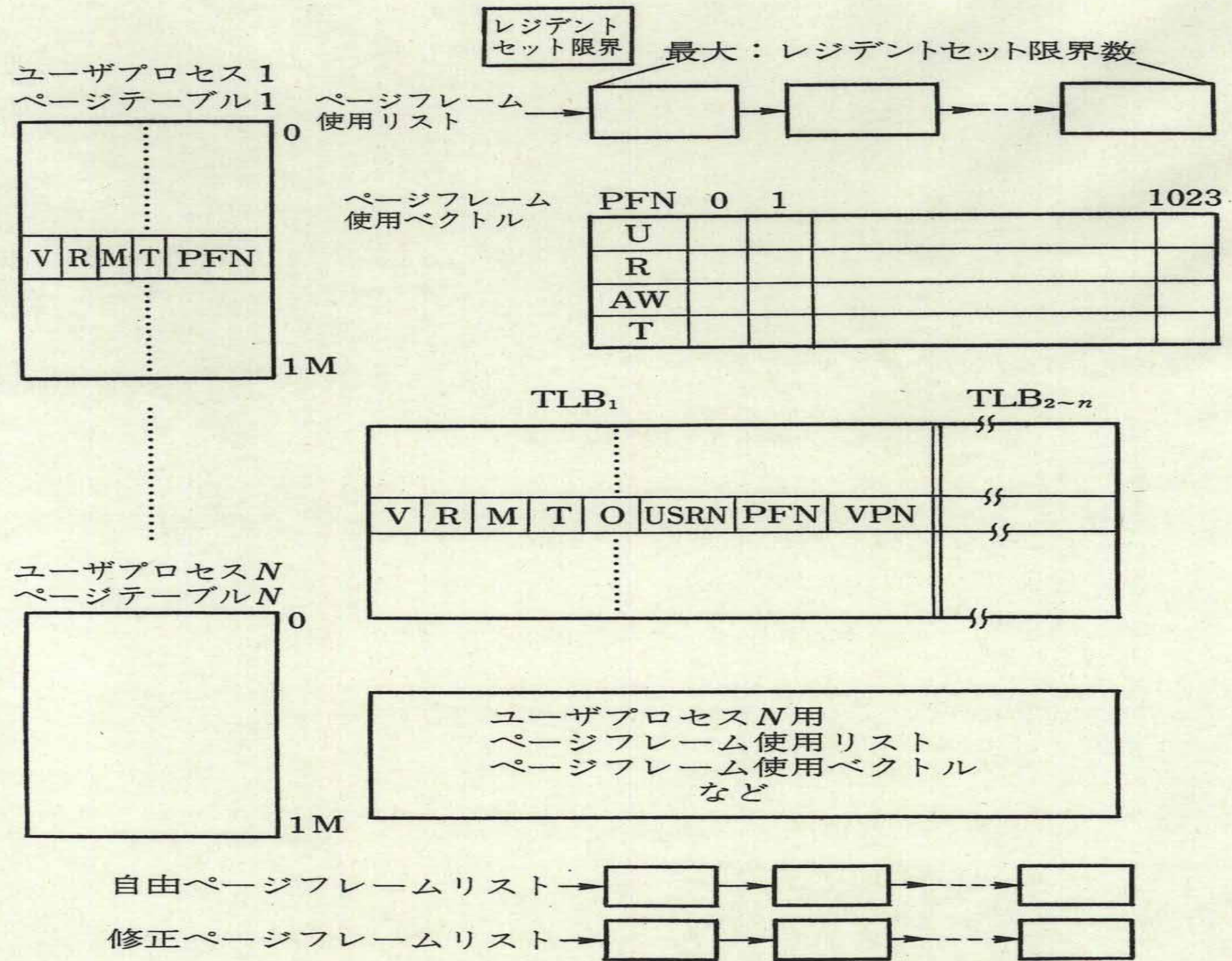


図 4.13 ページフレームの管理

マイコンピュータ

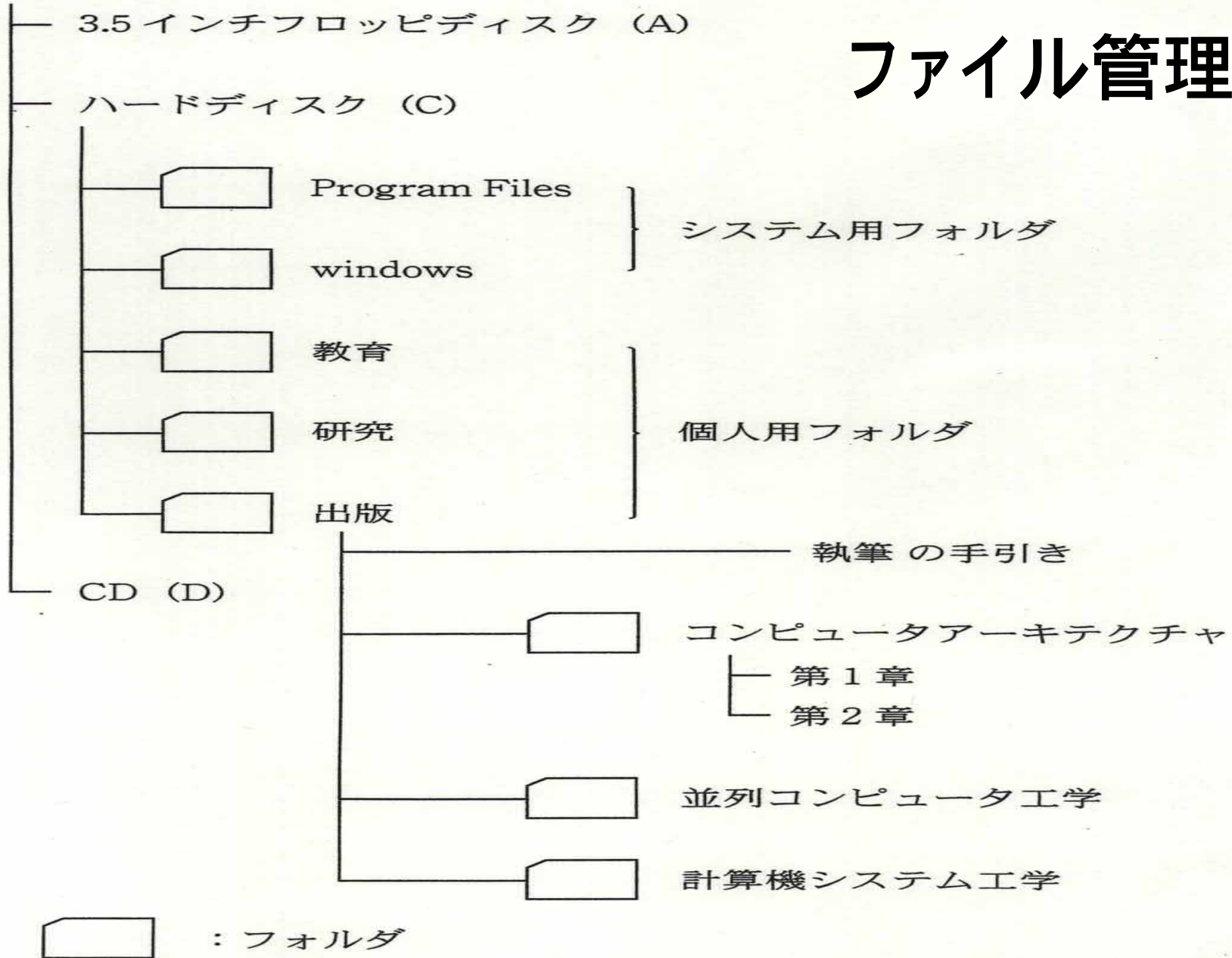


図 1.11 ファイルのディレクトリ構造

通信管理

アプリケーション層

FTP, Telnet, WWW, SMTP, ...

TCP / UDP層

TCP: コネクション型

バーチャルチャネル

UDP: コネクションレス型

IP層

中継ノード間での通信制御

IPアドレス

インタフェース層

マルチキャスト機能

物理層

イーサネット, ATM

プロセス管理

OSの中のOS

たくさんのプロセスの擬似的な並行実行

プロセス: OSにより管理実行されているプログラム

並列プロセスの例

長時間かかるCプログラムを実行させながら

WORDで文書作成し, プリントアウト

プリンタが動作し始め

EXCELで表計算処理を開始していると

時計の表示

電子メールの受信通知の表示

プロセスの状態

実行中: プロセスがOSからプロセッサを割り当てられて実行中の状態

レディ: プロセスは実行可能状態であるが、OSのプロセッサ割当てがなされていない状態

待ち: プロセスに必要なデータが揃っていないので、待たねばならない状態

プロセススイッチ

契機：割込み

外部割込み：入出力，タイマ
マシンチェック

内部割込み：演算例外，
命令例外，ページ
フォルト，トレース
スーパーバイザコール

多重プログラミング，

TSS (Time Sharing System)

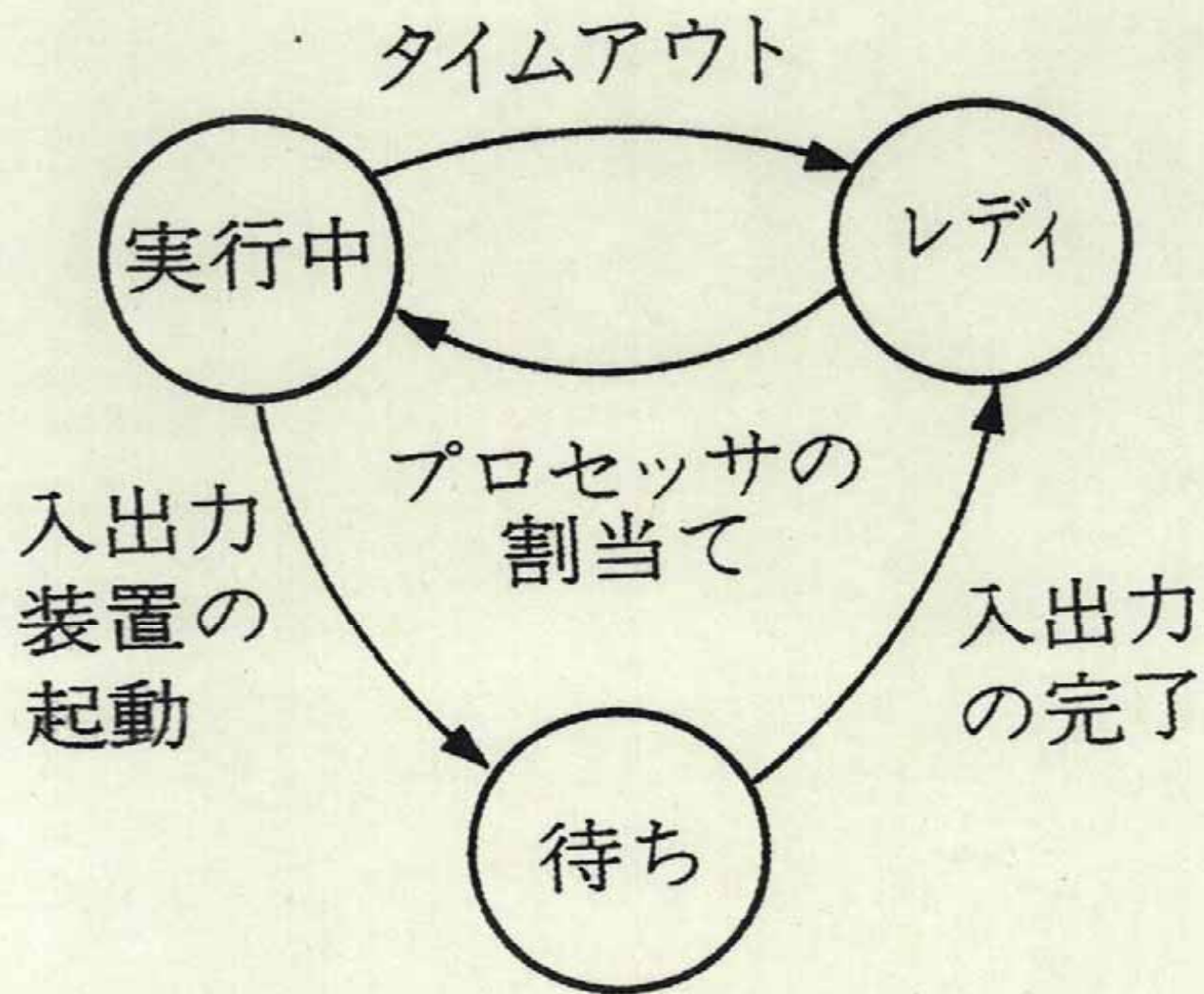


図 1.12 プロセス状態遷移

プロセスA, B, C: レディ

OS: プロセスAを選択し, 実行させる

プロセスA: 実行中

プロセスB, C: レディ

プロセスA: ディスクアクセス

I/O命令実行でOSに割込み

OS: プロセスAのI/O処理し, Aを待ちへ
プロセスBを選択し, 実行へ

プロセスB: 実行中

プロセスC: レディ

プロセスA: 待ち

プロセス B : 一定時間実行 (タイムクアラム)

10 msec

タイマ割込み

OS が B をレディへ, C を選択

プロセス C : 実行

プロセス B : レディ

プロセス A : 待ち

ディスクからI/O割込み

OSがチェックし、プロセスAのディスク

アクセス終了を知る

OSはプロセスAをレディへ

プロセスC: 実行中

プロセスB: レディ
プロセスA: レディ

プロセスC: ディスクアクセス

OS: プロセスCを待ちへ、

プロセスAを実行へ

プロセスA: 実行

プロセスB: レディ

プロセスC: 待ち

- ・プロセススケジューリング

- ・プロセス間通信

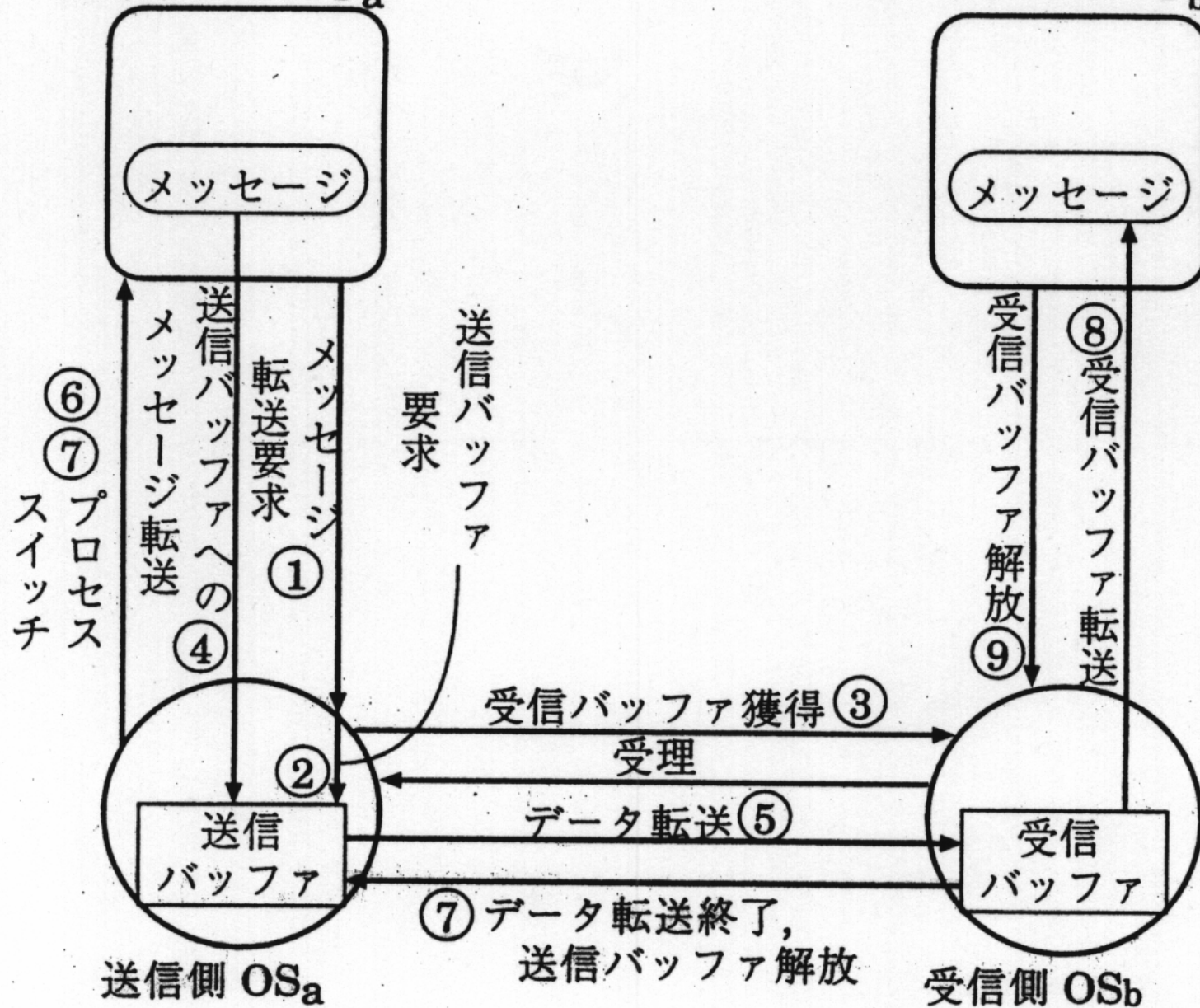
 - 同期：排他制御，事象待ち

 - メッセージ交換

 - メモリ空間共有，非共有

ユーザプロセス P_a

ユーザプロセス P_b



ユーザプロセスでは

Paよりスーパーバイザコール（OS呼出し）で
メッセージ転送要求

送信側OSaでは

送信バッファ獲得

受信側OSbに受信バッファ獲得依頼

Paのメッセージを送信バッファに転送

送信バッファのメッセージのパケット化。

OSbの受信バッファに相互結合網を通して転
送

プロセススイッチでPaは待機状態

データ転送終了。OSbより割込みを受理。

送信バッファ解放。Pa実行可状態

受信側OSbでは

受信バッファにデータ転送後

外部割込みによって起動

受信バッファよりPbの所定の領域にデータ
転送

受信バッファを解放

Pbを実行可状態

(4) メッセージ処理の高速化

メッセージのコピーによるオーバヘッド

割込み時でのレジスタ退避などのオーバヘッド

相互結合網での各種受理信号（ACK信号など）の
ための通信回数の増加

高速化の基本

ユーザプロセス自身によって通信を制御

OSのオーバヘッドを削減

通信用プロセッサの設置

通信処理：汎用マイクロプロセッサでよいのか？

- ・ 通信自体がイベント起動で、処理の切換え
頻繁
- ・ データが一過性。キャッシュメモリの参照
の局所性に不適合

通信処理とデータ転送のオーバーラップ

資源の専有

無駄なメッセージコピーの削減

種々の転送モードの用意

キャッシュ無効化の高速化

割込み回数の削減

(5) 放送機能の強化

1 対多

多対多

(6) 同期操作の高速化

バリア同期は

バリア到達の通知

バリア解除までに待機

ユーザレベル通信

OSの介在の少ない方式

Zero - Copy

仮想記憶によるセキュリティ確保

(1)DMAを用いた方式

固定DMA領域にコピー : AM-II, Hamlyn

軽いアドレス変換カーネルを毎回起動、BIP, LFC

TLBキャッシュ : VMMC-2, U-NET

(2)プログラムモードバス方式

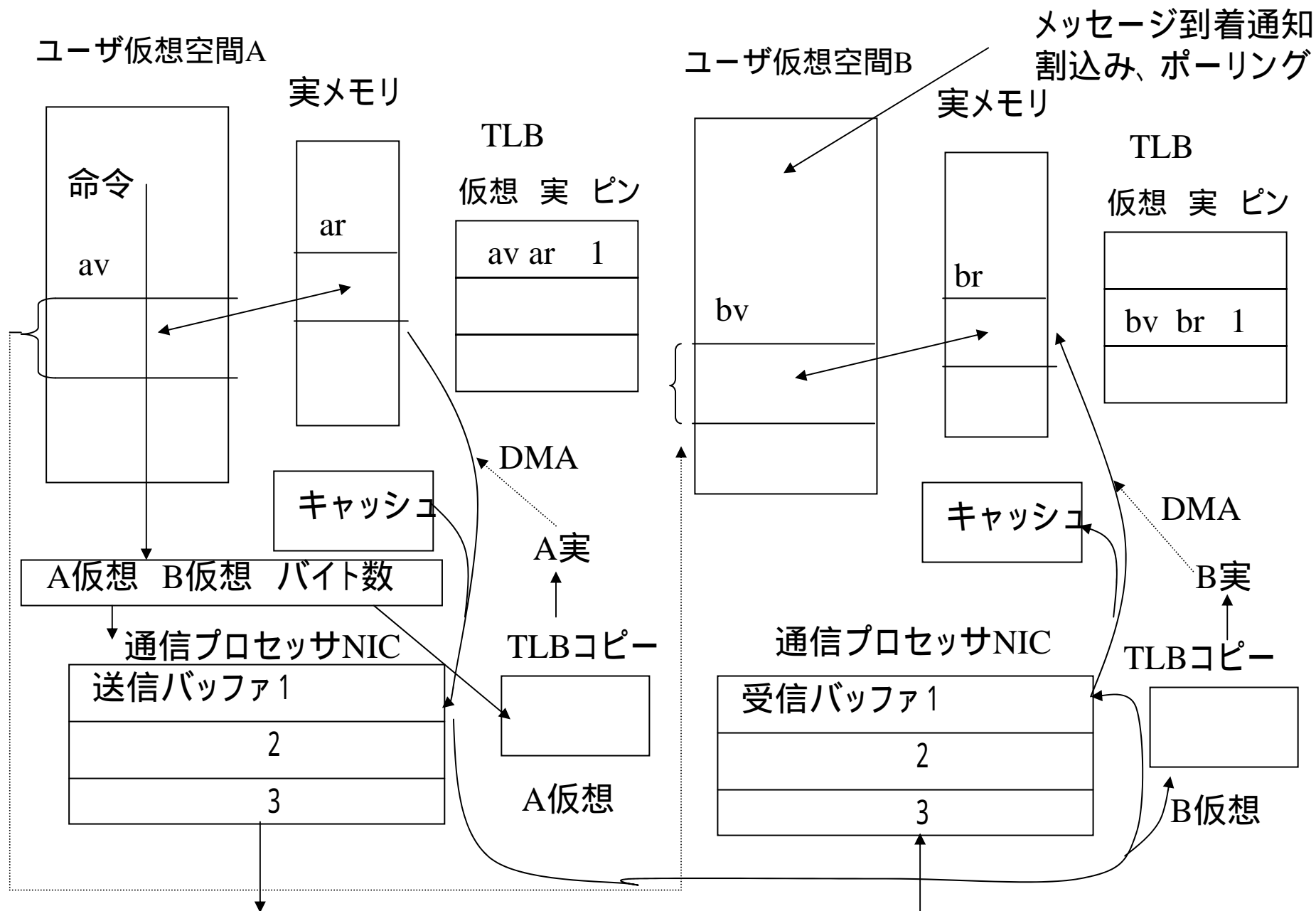
Write Combining あり : FM, LFC, AM-II, Hamlyn, BIP

同上なし

Table 1. Characteristics of 11 communication systems built for Myrinet.

System	Data transfer (host-M)	Translation	Protection	Control transfer	Reliability	Multicast support
AM-II ¹	PIO & DMA*	DMA areas	Yes	Polling + interrupts	Reliable, network interface: alternating bit, host: sliding window	No
FM ²	PIO	DMA area (recv)	No	Polling	Reliable, host-level credits	No
FM/MC ³	PIO	DMA area (recv)	No	Polling + interrupts	Reliable, unicast: host-level credits, multicast: network- interface-level credits	Yes (on network interface)
PM ⁴	DMA	Software TLB* on network interface	Yes (gang scheduling)	Polling	Reliable, ACK/NACK protocol on network interface	Yes (multiple sends)
VMMC ⁵	DMA	Software TLB on network interface	Yes	Polling + interrupts	Reliable, exploits hardware backpressure	No
VMMC-2 ⁶	DMA	UTLB* in kernel, cached on network interface	Yes	Polling + interrupts	Reliable	No
LFC ⁷	PIO	User translates	No	Polling + interrupts + watchdog	Reliable, unicast: network-interface- level credits, multicast: network-interface-level credits	Yes (on network interface)
Hamlyn ⁸	PIO & DMA	DMA areas	Yes	Polling + interrupts	Reliable, exploits hardware backpressure	No
Trapeze ⁹	DMA	DMA to page frames	No	Polling + interrupts	Unreliable	No
BIP ¹⁰	PIO & DMA	User translates	No	Polling	Reliable, rendezvous and backpressure	No
U-Net ¹¹	DMA	TLB on network interface (U-Net/MM)	Yes	Polling + interrupts	Unreliable	No

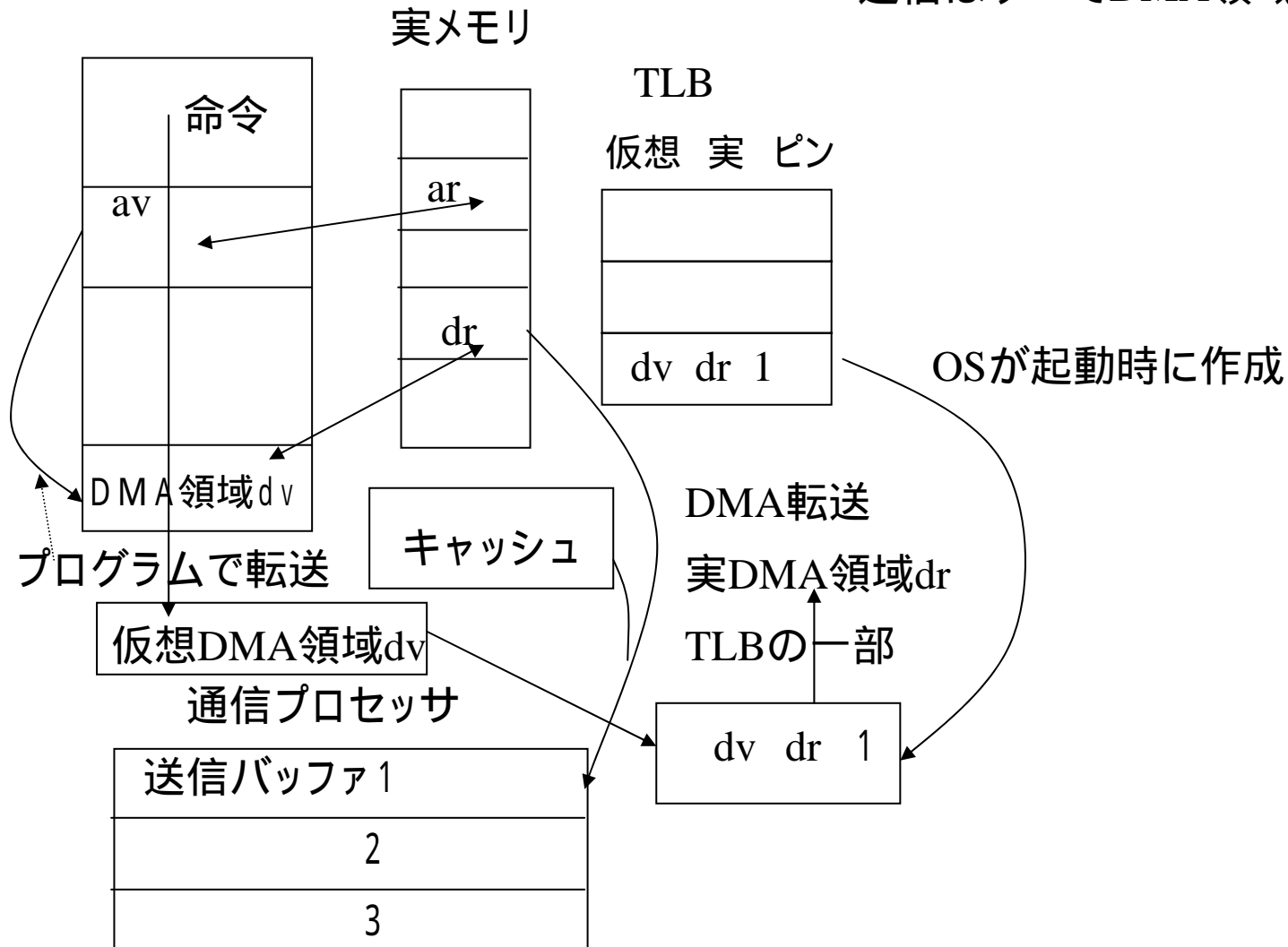
DMA基本方式 転送実領域:貼付け(ピン)、通信プロセッサ TLBアクセスの必要



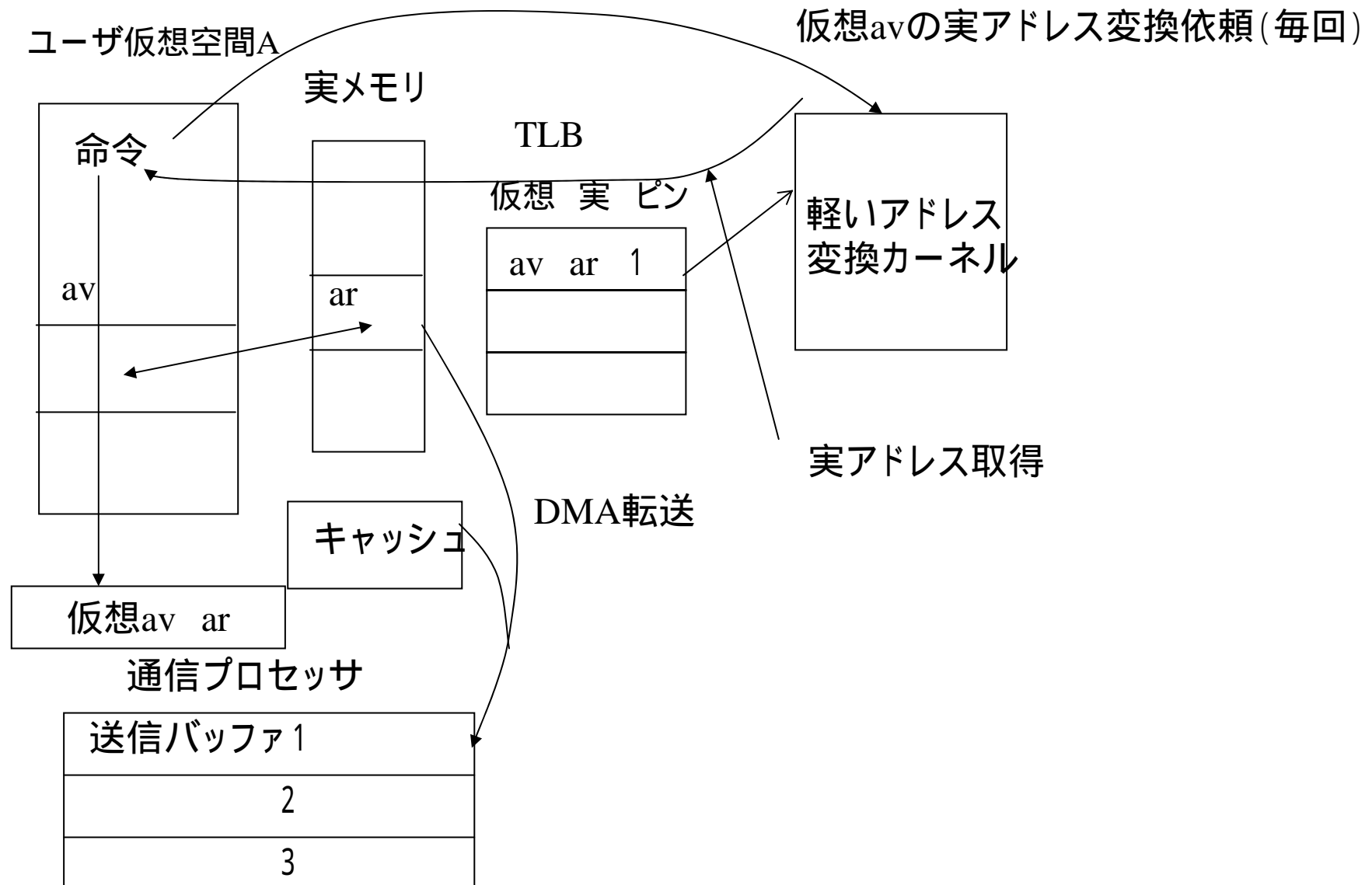
方式 : 仮想空間AのDMA領域をOSに依頼して貼付け(起動時一度だけ)

ユーザ仮想空間A

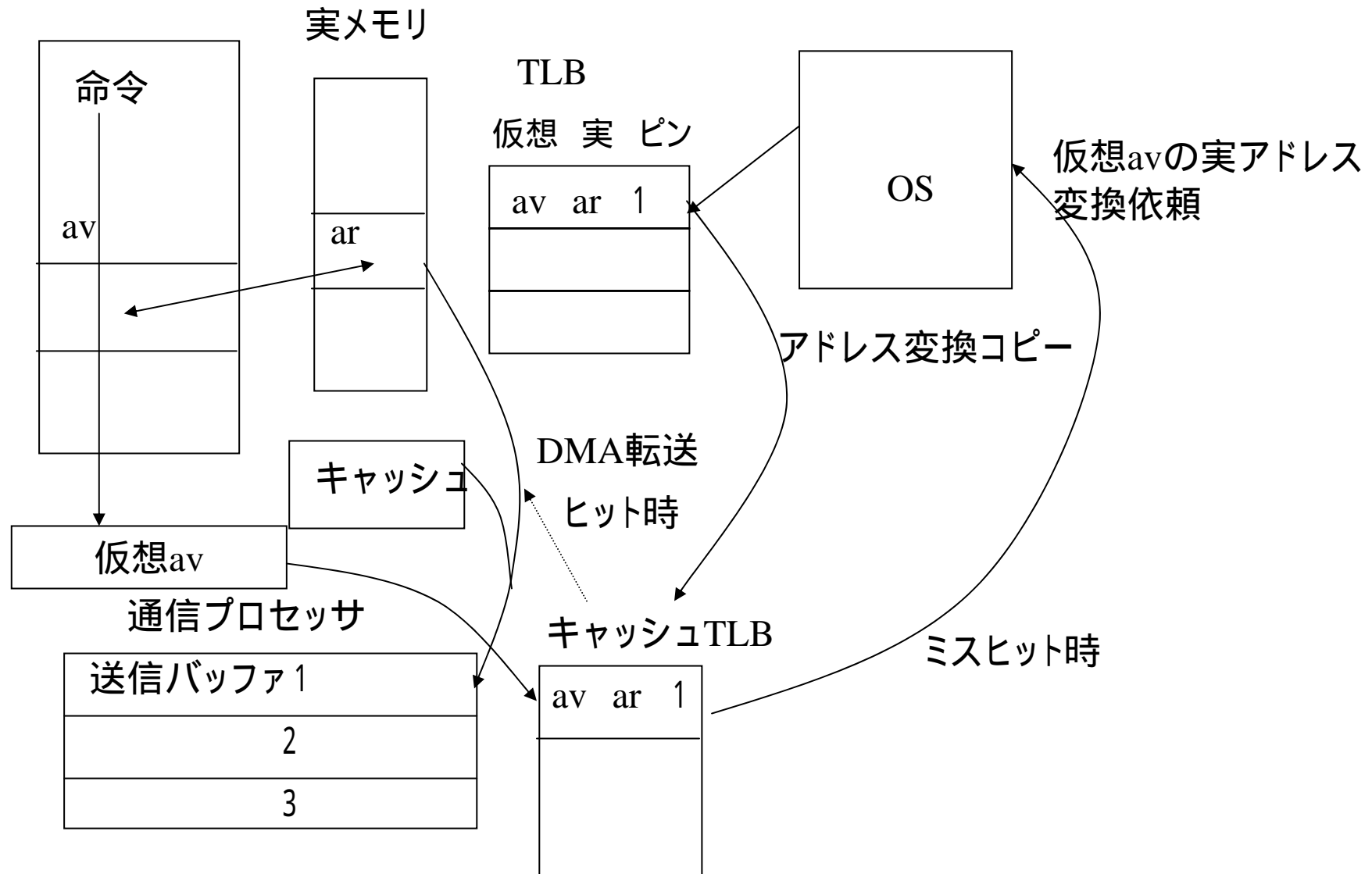
通信はすべてDMA領域を通して行う



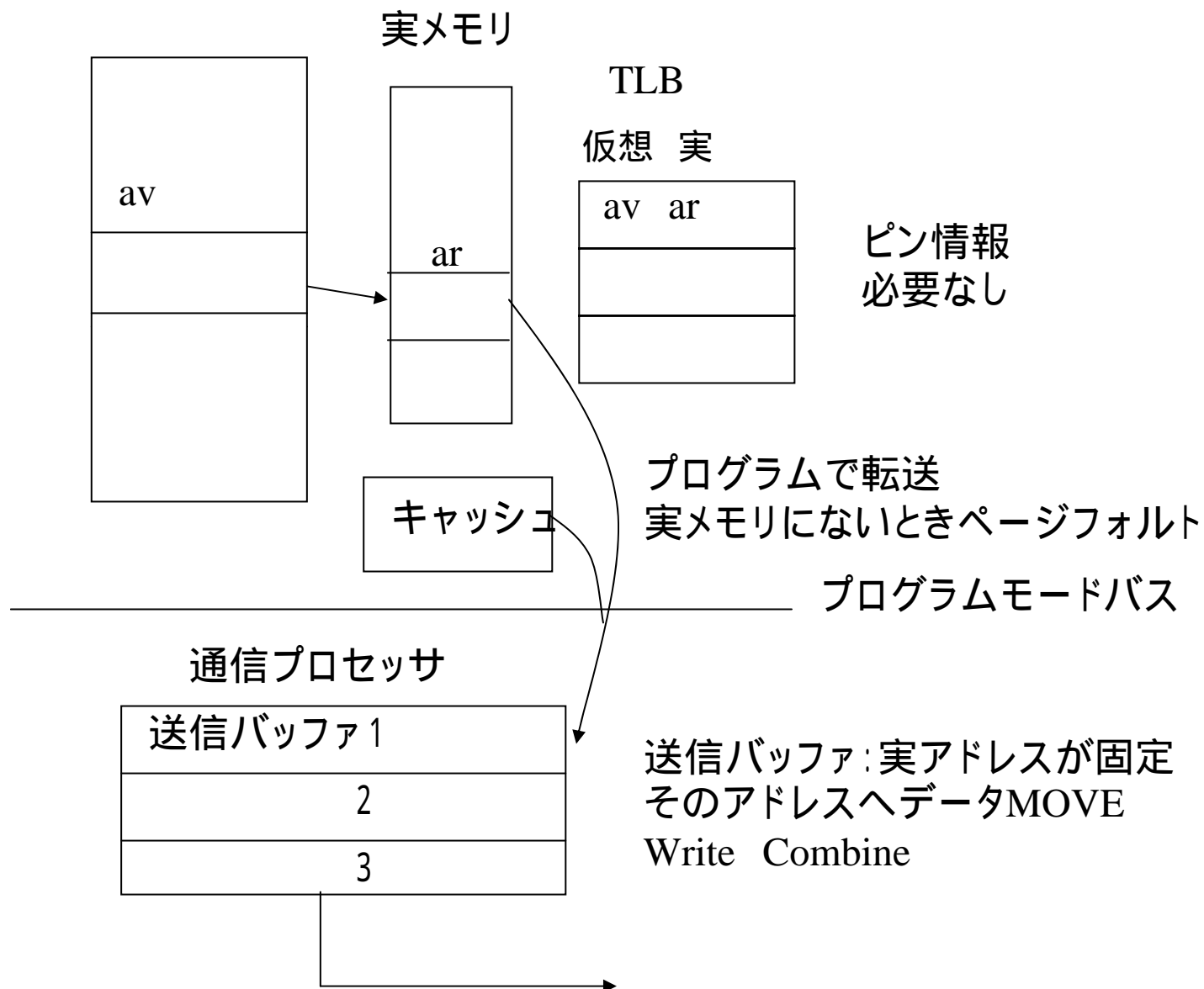
方式：軽いアドレス変換カーネル

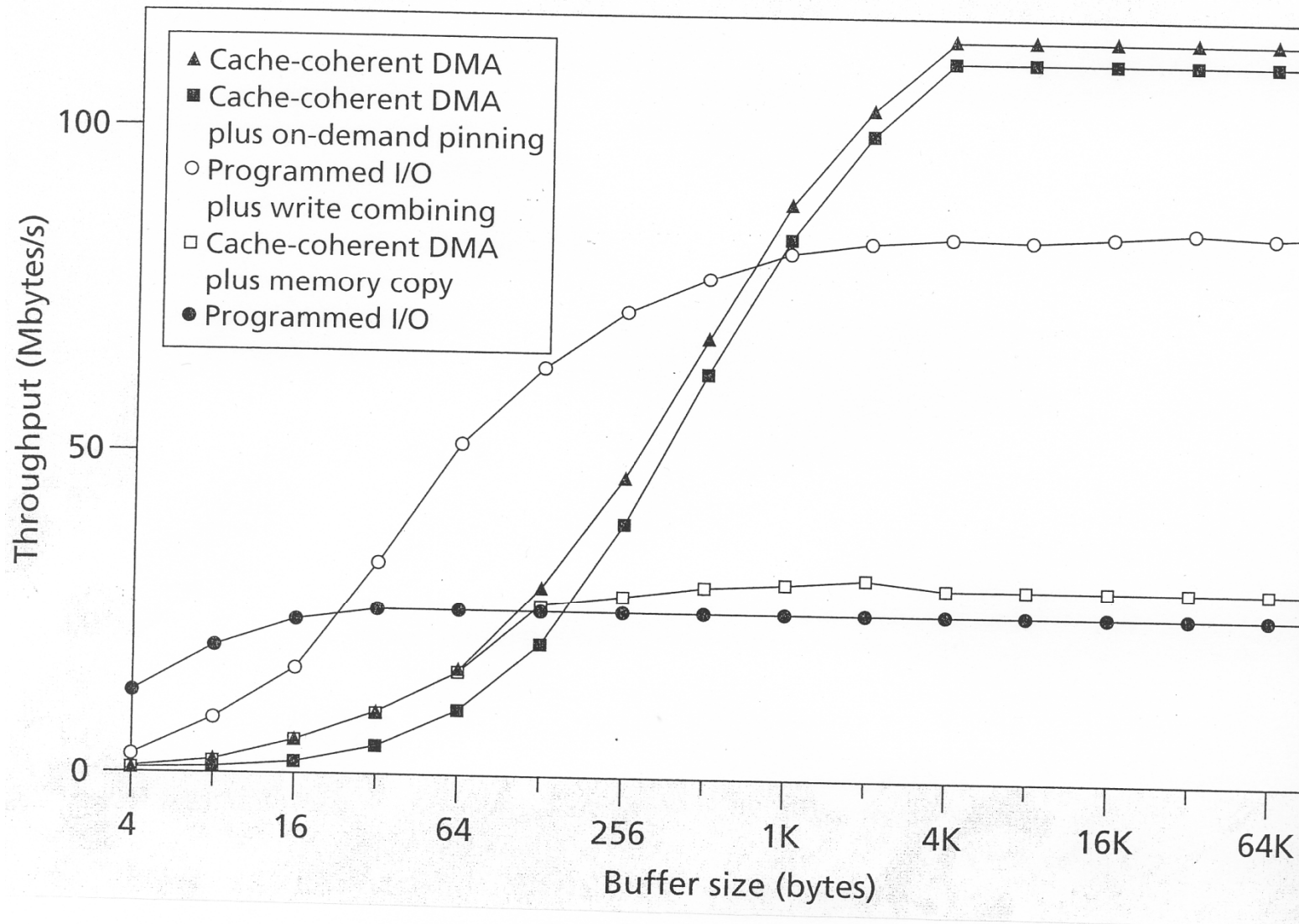


ユーザ仮想空間A



方式 、 プログラムモードバス
DMAを使わない方式
ユーザ仮想空間A





MPI Benchmark	MX/Myrinet Myricom 10G Myrinet switch	MX/Ethernet Fulcrum 10G Ethernet switch	MX/Ethernet Fujitsu 10G Ethernet switch	OpenIB with Intel MPI Mellanox InfiniBand
PingPong latency	2.4 μ s	2.4 μ s	2.8 μ s	4.0 μ s
One-way data rate (PingPong)	1204 MByte/s	1201 MByte/s	1002 MByte/s	964 MByte/s
Two-way data rate (SendRecv)	2397 MByte/s	2162 MByte/s	1762 MByte/s	1902 MByte/s

MX: Myrinet Express: メッセージパッシングソフト

Myri-10G: 10 Gigabit/s, dual protocol NIC

5.5.2 A P 1 0 0 0 の構成

メッセージ交換型マルチプロセッサの代表例

富士通のAP1000、NECのCenju-3、MITのJ-Machine

(1) プロセッサ構成

ノードプロセッサ：SPARC、通信プロセッサ（MSC）

トーラス網

ラインセンド / バッファレシーブ機構

1 0 B , 1 0 K B のメッセージ転送の場合

通常的方式：112 μ sec/1590 μ sec

本機構：32.8 μ sec/491 μ secに短縮

ストライド転送

(2) 相互結合網

T-net : 各プロセッサでのメッセージ交換を行う

2次元トーラス網

リンク25MB/secの転送能力

ワームホール方式

ルーティング : 次元順

S-netはバリア同期用のトリー網

B-netはホスト計算機からの放送やホスト計算機
へのデータ収集

(3) 後継機

AP1000+

PUT/GET機構とその専用ハードウェア (MSC+)

PUT

GET

キャッシュ：ストアスルー方式

通信手順

送信プロセスMSC+にコマンド

コマンド：受信プロセッサ番号、転送先の論理アドレス、送信元の論理アドレス

MSC+ : 送信データの論理アドレスを物理アドレスに変換

物理アドレスよりDMAでデータ転送

転送終了すると、指定されたメモリ領域にあるフラグを更新

受信側MSC+ : データ転送

受信データの論理アドレスを物理アドレスに変換

宛先領域にDMA転送

フラグの変更

受信プロセッサのメモリに書込み時：

キャッシュ無効化

データ転送の終了通知：

割込みオーバヘッド大

フラグを陽にチェック

通信の保護：論理（仮想）アドレス

AP1000でPUT命令実現：68 μ sec

AP1000+で実現すると5.1 μ sec

通信はすべてユーザプロセスで実現

最新のシステム例

FUJITSU PRIMEPOWER HPC2500

地球シミュレータ(NEC SX-8)

日立SR11000

CRAY X1

BlueGene/L

(1) ベクトルパラレル VS スカラパラレル

(2) メモリ共有 VS メッセージパッシング

(3) ネットワーク
クロスバ VS トーラスなど他網

、 、

、

TOP500LIST- June2000

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>Sandia National Laboratories</u> United States	<u>ASCI Red</u> Intel	9632	1999	2379	3207
2	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI Blue-Pacific SST,</u> <u>IBM SP 604e</u> IBM	5808	1999	2144	3856.5
3	<u>Los Alamos National Laboratory</u> United States	<u>ASCI Blue Mountain</u> SGI	6144	1998	1608	3072
4	<u>IBM/Naval Oceanographic Office</u> <u>(NAVOCEANO)</u> United States	<u>SP Power3 375 MHz</u> IBM	1336	2000	1417	2004
5	<u>Leibniz Rechenzentrum</u> Germany	<u>SR8000-F1/112</u> Hitachi	112	2000	1035	1344
6	<u>High Energy Accelerator Research</u> <u>Organization /KEK</u> Japan	<u>SR8000-F1/100</u> Hitachi	100	2000	917	1200
7	<u>Government</u> United States	<u>T3E1200</u> Cray Inc.	1084	1998	891	1300.8
8	<u>US Army HPC Research Center at NCS</u> United States	<u>T3E1200</u> Cray Inc.	1084	2000	891	1300.8
9	<u>University of Tokyo</u> Japan	<u>SR8000/128</u> Hitachi	128	1999	873	1024
10	<u>Government</u> United States	<u>T3E900</u> Cray Inc.	1324	1997	815	1191.6

TOP500LIST-June2001

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI White, SP Power3</u> <u>375 MHz</u> IBM	8192	2000	7226	12288
2	<u>NERSC/LBNL</u> United States	<u>SP Power3 375 MHz 16</u> <u>way</u> IBM	2528	2001	2526	3792
3	<u>Sandia National Laboratories</u> United States	<u>ASCI Red</u> Intel	9632	1999	2379	3207
4	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI Blue-Pacific SST,</u> <u>IBM SP 604e</u> IBM	5808	1999	2144	3856.5
5	<u>University of Tokyo</u> Japan	<u>SR8000/MPP</u> Hitachi	1152	2001	1709.1	2074
6	<u>Los Alamos National Laboratory</u> United States	<u>ASCI Blue Mountain</u> SGI	6144	1998	1608	3072
7	<u>Naval Oceanographic Office</u> (NAVOCEANO) United States	<u>SP Power3 375 MHz</u> IBM	1336	2000	1417	2004
8	<u>Osaka University</u> Japan	<u>SX-5/128M8 3.2ns</u> NEC	128	2001	1192	1280
9	<u>National Centers for Environmental</u> <u>Prediction</u> United States	<u>SP Power3 375 MHz</u> IBM	1104	2000	1179	1656
10	<u>National Centers for Environmental</u> <u>Prediction</u> United States	<u>SP Power3 375 MHz</u> IBM	1104	2001	1179	1656

TOP500LIST- June2002

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>The Earth Simulator Center</u> Japan	<u>Earth-Simulator</u> NEC	5120	2002	35860	40960
2	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI White, SP Power3 375 MHz</u> IBM	8192	2000	7226	12288
3	<u>Pittsburgh Supercomputing Center</u> United States	<u>AlphaServer SC45, 1 GHz</u> Hewlett-Packard	3016	2001	4463	6032
4	<u>Commissariat a l'Energie Atomique (CEA)</u> France	<u>AlphaServer SC45, 1 GHz</u> Hewlett-Packard	2560	2001	3980	5120
5	<u>NERSC/LBNL</u> United States	<u>SP Power3 375 MHz 16 way</u> IBM	3328	2001	3052	4992
6	<u>Los Alamos National Laboratory</u> United States	<u>AlphaServer SC45, 1 GHz</u> Hewlett-Packard	2048	2002	2916	4096
7	<u>Sandia National Laboratories</u> United States	<u>ASCI Red</u> Intel	9632	1999	2379	3207
8	<u>Oak Ridge National Laboratory</u> United States	<u>pSeries 690 Turbo 1.3GHz</u> IBM	864	2002	2310	4492.8
9	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI Blue-Pacific SST,</u> <u>IBM SP 604e</u> IBM	5808	1999	2144	3856.5
10	<u>IBM/US Army Research Laboratory (ARL)</u> United States	<u>pSeries 690 Turbo 1.3GHz</u> IBM	768	2002	2050	3993.6

TOP500LIST-June2003

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>The Earth Simulator Center</u> Japan	<u>Earth-Simulator</u> NEC	5120	2002	35860	40960
2	<u>Los Alamos National Laboratory</u> United States	<u>ASCI Q - AlphaServer SC45, 1.25 GHz</u> Hewlett-Packard	8192	2002	13880	20480
3	<u>Lawrence Livermore National Laboratory</u> United States	<u>MCR Linux Cluster Xeon 2.4 GHz - Quadrics</u> Linux Networx/Quadrics	2304	2002	7634	11060
4	<u>Lawrence Livermore National Laboratory</u> United States	<u>ASCI White, SP Power3 375 MHz</u> IBM	8192	2000	7304	12288
5	<u>NERSC/LBNL</u> United States	<u>Seaborg - SP Power3 375 MHz 16 way</u> IBM	6656	2002	7304	9984
6	<u>Lawrence Livermore National Laboratory</u> United States	<u>xSeries Cluster Xeon 2.4 GHz - Quadrics</u> IBM/Quadrics	1920	2003	6586	9216
7	<u>National Aerospace Laboratory of Japan</u> Japan	<u>PRIMEPOWER HPC2500 (1.3 GHz)</u> Fujitsu	2304	2002	5406	11980
8	<u>Pacific Northwest National Laboratory</u> United States	<u>Cluster Platform 6000 rx2600 Itanium2 1 GHz Cluster - Quadrics</u> Hewlett-Packard	1540	2003	4881	6160
9	<u>Pittsburgh Supercomputing Center</u> United States	<u>AlphaServer SC45, 1 GHz</u> Hewlett-Packard	3016	2001	4463	6032
10	<u>Commissariat a l'Energie Atomique (CEA)</u> France	<u>AlphaServer SC45, 1 GHz</u> Hewlett-Packard	2560	2001	3980	5120

TOP500LIST-June2004

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>The Earth Simulator Center</u> Japan	<u>Earth-Simulator</u> NEC	5120	2002	35860	40960
2	<u>Lawrence Livermore National Laboratory</u> United States	<u>Thunder - Intel Itanium2 Tiger4</u> 1.4GHz - Quadrics California Digital Corporation	4096	2004	19940	22938
3	<u>Los Alamos National Laboratory</u> United States	<u>ASCI Q - AlphaServer SC45, 1.25 GHz</u> Hewlett-Packard	8192	2002	13880	20480
4	<u>IBM - Rochester</u> United States	<u>BlueGene/L DD1 Prototype (0.5GHz PowerPC 440 w/Custom)</u> IBM/ LLNL	8192	2004	11680	16384
5	<u>NCSA</u> United States	<u>Tungsten - PowerEdge 1750, P4 Xeon 3.06 GHz, Myrinet</u> Dell	2500	2003	9819	15300
6	<u>ECMWF</u> United Kingdom	<u>eServer pSeries 690 (1.9 GHz Power4+)</u> IBM	2112	2004	8955	16051
7	<u>Institute of Physical and Chemical Res. (RIKEN)</u> Japan	<u>RIKEN Super Combined Cluster</u> Fujitsu	2048	2004	8728	12534
8	<u>IBM Thomas J. Watson Research Center</u> United States	<u>BlueGene/L DD2 Prototype (0.7 GHz PowerPC 440)</u> IBM/ LLNL	4096	2004	8655	11469
9	<u>Pacific Northwest National Laboratory</u> United States	<u>Mpp2 - Cluster Platform 6000 rx2600 Itanium2 1.5 GHz, Quadrics</u> Hewlett-Packard	1936	2003	8633	11616
10	<u>Shanghai Supercomputer Center</u> China	<u>Dawning 4000A, Opteron 2.2 GHz, Myrinet</u> Dawning	2560	2004	8061	11264

TOP500LIST-June2005

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>DOE/NNSA/LLNL</u> United States	<u>BlueGene/L - eServer Blue Gene Solution</u> IBM	65536	2005	136800	183500
2	<u>IBM Thomas J. Watson Research Center</u> United States	<u>BGW - eServer Blue Gene Solution</u> IBM	40960	2005	91290	114688
3	<u>NASA/Ames Research Center/NAS</u> United States	<u>Columbia - SGI Altix 1.5 GHz, Voltaire Infiniband</u> SGI	10160	2004	51870	60960
4	<u>The Earth Simulator Center</u> Japan	<u>Earth-Simulator</u> NEC	5120	2002	35860	40960
5	<u>Barcelona Supercomputer Center</u> Spain	<u>MareNostrum - JS20 Cluster, PPC 970, 2.2 GHz, Myrinet</u> IBM	4800	2005	27910	42144
6	<u>ASTRON/University Groningen</u> Netherlands	<u>Stella - eServer Blue Gene Solution</u> IBM	12288	2005	27450	34406.4
7	<u>Lawrence Livermore National Laboratory</u> United States	<u>Thunder - Intel Itanium2 Tiger4 1.4GHz - Quadrics</u> California Digital Corporation	4096	2004	19940	22938
8	<u>Computational Biology Research Center, AIST</u> Japan	<u>Blue Protein - eServer Blue Gene Solution</u> IBM	8192	2005	18200	22937.6
9	<u>Ecole Polytechnique Federale de Lausanne</u> Switzerland	<u>eServer Blue Gene Solution</u> IBM	8192	2005	18200	22937.6
10	<u>Sandia National Laboratories</u> United States	<u>Red Storm, Cray XT3, 2.0 GHz</u> Cray Inc.	5000	2005	15250	20000

TOP500LIST-June2006

Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	<u>DOE/NNSA/LLNL</u> United States	<u>BlueGene/L - eServer Blue Gene Solution</u> IBM	131072	2005	280600	367000
2	<u>IBM Thomas J. Watson Research Center</u> United States	<u>BGW - eServer Blue Gene Solution</u> IBM	40960	2005	91290	114688
3	<u>DOE/NNSA/LLNL</u> United States	<u>ASC Purple - eServer pSeries p5 575</u> <u>1.9 GHz</u> IBM	12208	2006	75760	92781
4	<u>NASA/Ames Research Center/NAS</u> United States	<u>Columbia - SGI Altix 1.5 GHz, Voltaire Infiniband</u> SGI	10160	2004	51870	60960
5	<u>Commissariat a l'Energie Atomique (CEA)</u> France	<u>Tera-10 - NovaScale 5160, Itanium2</u> <u>1.6 GHz, Quadrics</u> Bull SA	8704	2006	42900	55705.6
6	<u>Sandia National Laboratories</u> United States	<u>Thunderbird - PowerEdge 1850, 3.6 GHz, Infiniband</u> Dell	9024	2006	38270	64972.8
7	<u>GSIC Center, Tokyo Institute of Technology</u> Japan	<u>TSUBAME Grid Cluster - Sun Fire X64 Cluster, Opteron 2.4/2.6 GHz, Infiniband</u> NEC/Sun	10368	2006	38180	49868.8
8	<u>Forschungszentrum Juelich (FZJ)</u> Germany	<u>JUBL - eServer Blue Gene Solution</u> IBM	16384	2006	37330	45875
9	<u>Sandia National Laboratories</u> United States	<u>Red Storm Cray XT3, 2.0 GHz</u> Cray Inc.	10880	2005	36190	43520
10	<u>The Earth Simulator Center</u> Japan	<u>Earth-Simulator</u> NEC	5120	2002	35860	40960

TOP500LIST-June2007

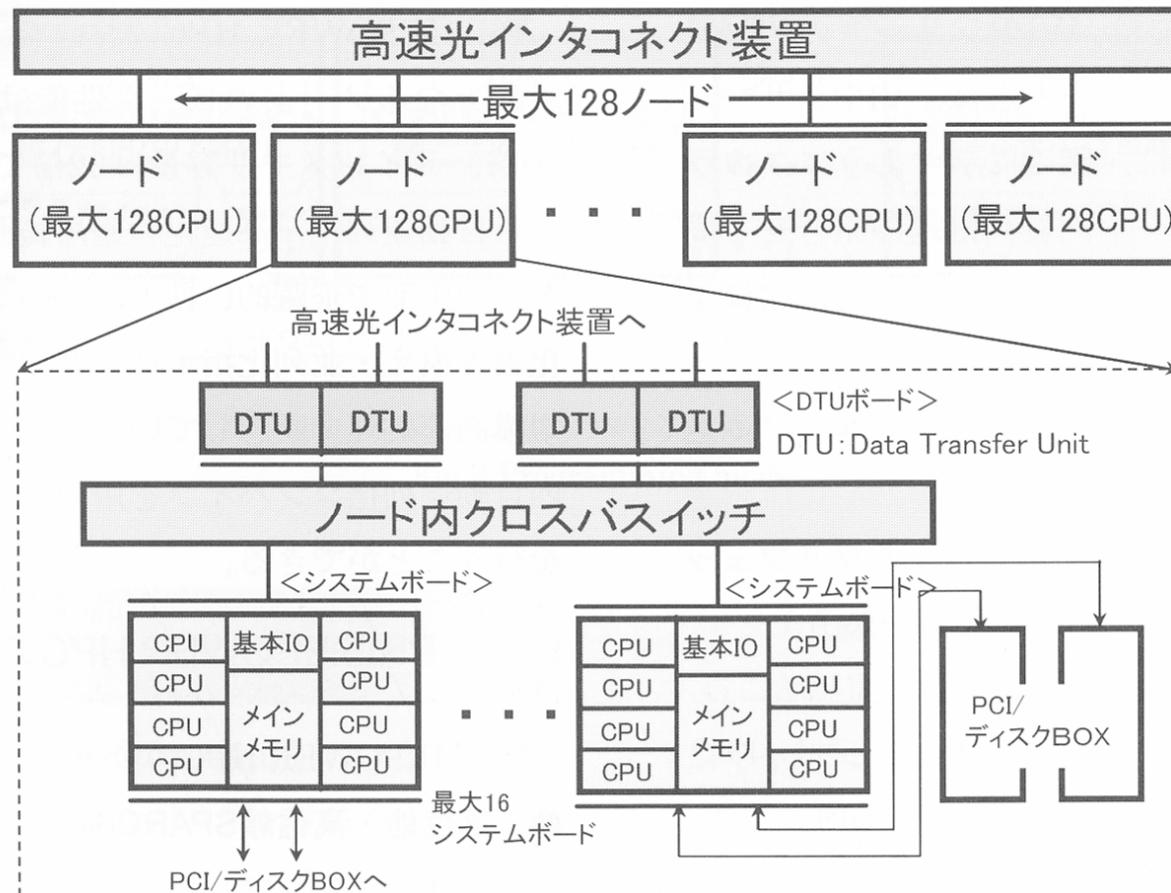
Rank	Site	Computer	Processors	Year	R _{max}	R _{peak}
1	DOE/NNSA/LLNL United States	BlueGene/L - eServer Blue Gene Solution IBM	131072	2005	280600	367000
2	Oak Ridge National Laboratory United States	Jaguar - Cray XT4/XT3 Cray Inc.	23016	2006	101700	119350
3	NNSA/Sandia National Laboratories United States	Red Storm - Sandia/ Cray Red Storm, Opteron 2.4 GHz dual core Cray Inc.	26544	2006	101400	127411
4	IBM Thomas J. Watson Research Center United States	BGW - eServer Blue Gene Solution IBM	40960	2005	91290	114688
5	Stony Brook/BNL, New York Center for Computational Sciences United States	New York Blue - eServer Blue Gene Solution IBM	36864	2007	82161	103219
6	DOE/NNSA/LLNL United States	ASC Purple - eServer pSeries p5 575 1.9 GHz IBM	12208	2006	75760	92781
7	Rensselaer Polytechnic Institute, Computational Center for Nanotechnology Innovations United States	eServer Blue Gene Solution IBM	32768	2007	73032	91750
8	NCSA United States	Abe - PowerEdge 1955, 2.33 GHz, Infiniband Dell	9600	2007	62680	89587.2
9	Barcelona Supercomputing Center Spain	MareNostrum - BladeCenter JS21 Cluster, PPC 970, 2.3 GHz, Myrinet IBM	10240	2006	62630	94208
10	Leibniz Rechenzentrum Germany	HLRB-II - Altix 4700 1.6 GHz SGI	9728	2007	56520	62259.2

日立スーパーコンピュータ

- ベクトルパラレル
- 1982 S-810 630MFLOPS
- 1987 S-820 3GFLOPS
- 1992 S-3000 32GFLOPS
- スカラパラレル
- 1995 SR2201 600GFLOPS
- 1999 SR8000 7.3TFLOPS
- 2003 SR11000 62TFLOPS
- (Power5(1.9GHz)、16PE/ノード、121.6GFLOPS/ノード、最大512ノード、多段クロスバネット:12GB/sx2(ノード当たり))

FUJITSU PRIMEPOWER HPC

VPPからSPPに切り替え



ノード: 共有メモリ、
スヌープ方式

SMP

512インタリーブ
メモリ(ノード内)

Fujitsu Vol.53, No.6,
2002, 特集 サーバ

図-1 PRIMEPOWER HPCシステムの構成
Fig.1-System configuration of PRIMEPOWER HPC.

- Fujitsu PRIMEPOWER HPC

表-1 PRIMEPOWER HPCノード諸元

項目	諸元
CPU	SPARC64 V
CPU周波数	1.3 GHz
最大CPU数	128
アドレススヌープ性能	133 Gバイト/秒
最大メインメモリ容量	512 Gバイト
最大メインメモリインタリーブ数	512ウェイ
最大PCIスロット数	320

表-2 PRIMEPOWER HPCシステム諸元

項目	諸元
最大ノード数	128
最大CPU数	16,384
最大論理性能	85.2 TFLOPS★
最大メインメモリ容量	64 Tバイト
ノード間結合方式	クロスバ
ノード間転送性能	1ノードあたり 最大16 Gバイト/秒×2 (入力/出力)

★ : Tera FLoating point Operation Per Second

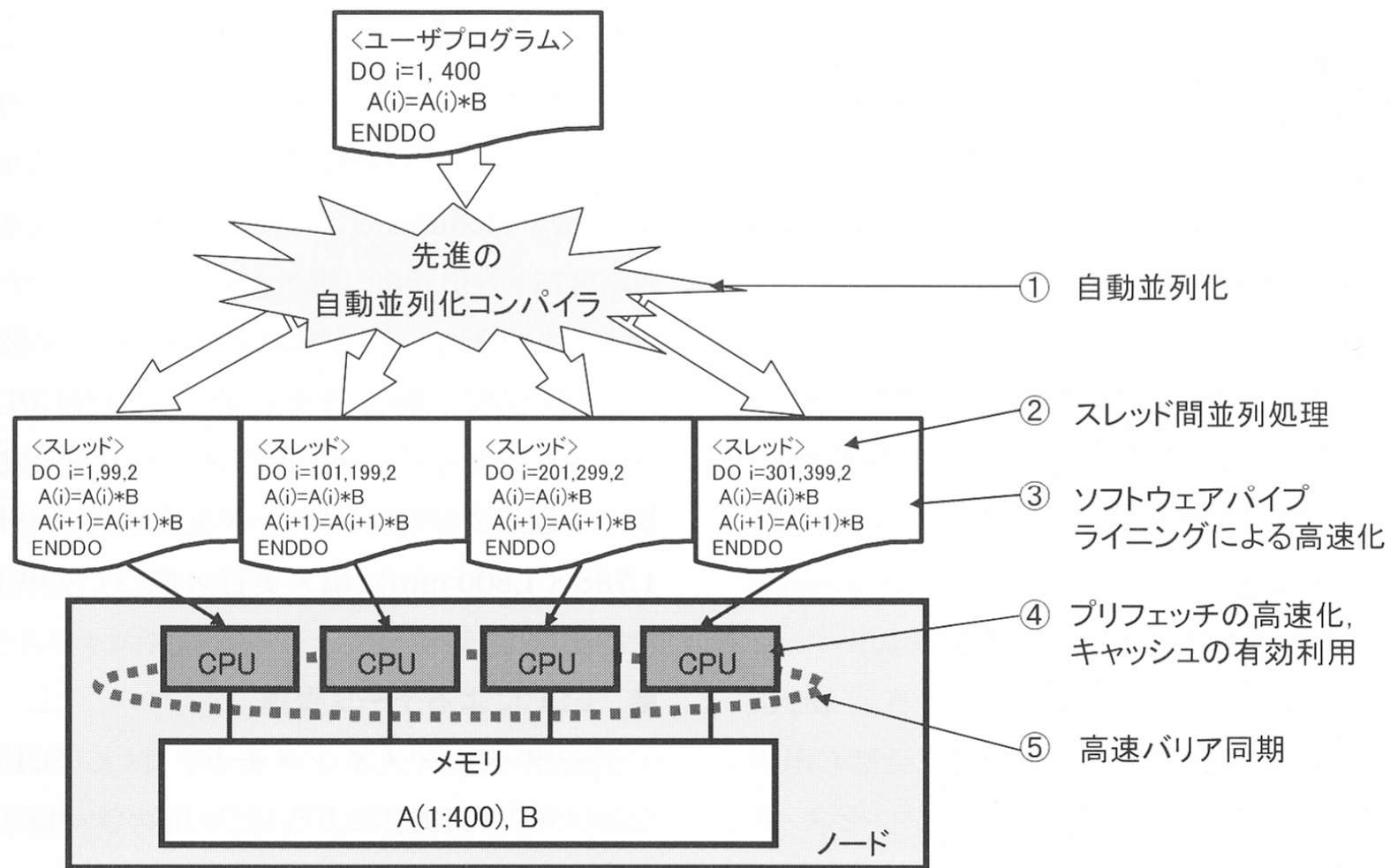


図-2 ノード内並列処理
Fig.2-Parallel execution method in node.

SPARC64 V

整数演算2台、
浮動小数点2台、
アドレス計算2台

5.2GFLOPS

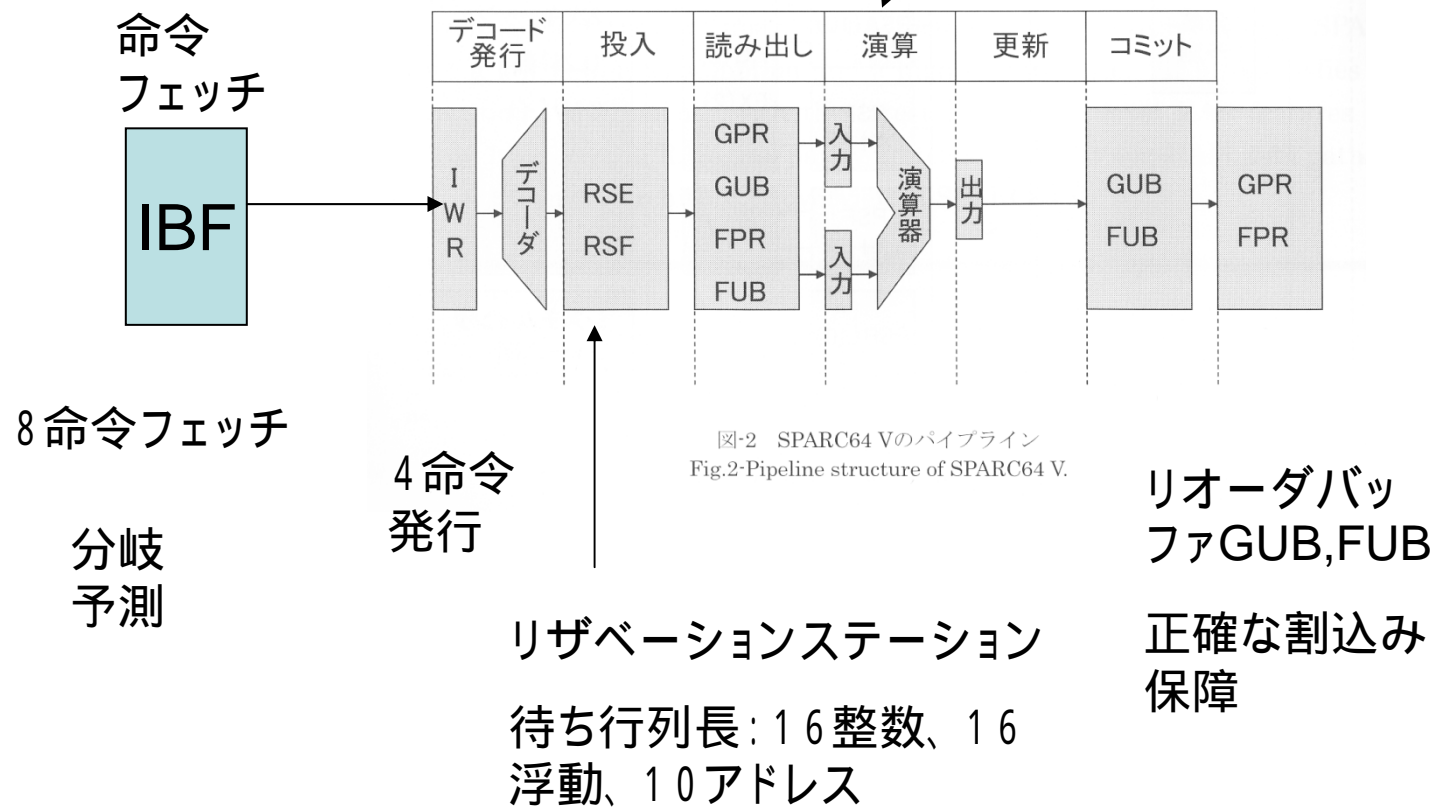


図-2 SPARC64 Vのパイプライン
Fig.2-Pipeline structure of SPARC64 V.

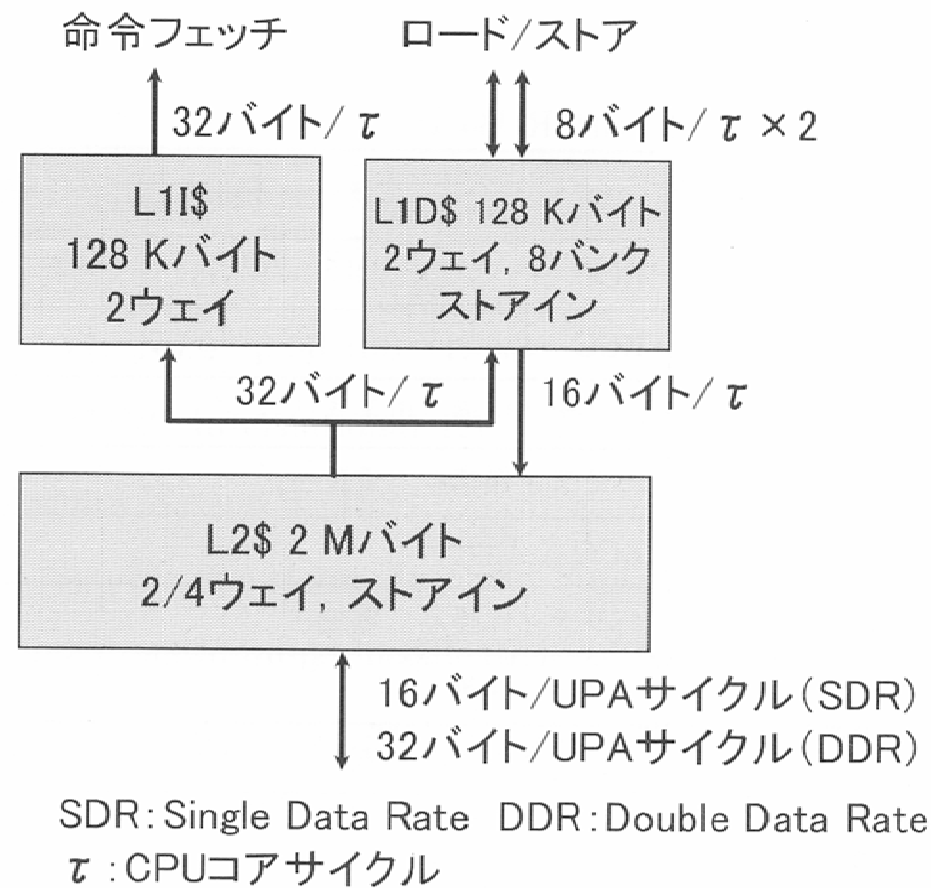


図-3 SPARC64 Vのキャッシュ

Fig.3-Cache structure of SPARC64 V.

- デザインルール: $0.13\ \mu\text{m}$
- TR数: 19,100万個
- 信号ピン数: 269
- チップサイズ: $17.8 \times 15.7\text{mm}$
- 動作周波数: 1.3GHz
- 消費電力: 50W

NECスーパーコンピュータ

- 機種 年 サイクル単体性能 最大性能 台数
- Cray-1 1976 12.5ns 160MF 160MF 1台
- SX-1/2 1984 6ns 1.3GF 1.3GF 1台
- SX-3 1989 2.9ns 5.5GF 22GF 4台
- SX-4 1994 8ns 2GF 1TF 512台
- SX-5 1998 4ns 8GF 4TF 512台
- SX-6 2001 2ns 8GF 8TF 1024台
- (CMOSシングルチップ、8PE/1ノード、最大128ノード、0.15 μm)
- SX-7 2002 1.8ns 11.4GF 23TF 2048台
- (32PE/1ノード、最大64ノード、0.15 μm)
- SX-8 2004 0.5ns 16GF 65TF 4096台
- (8PE/1ノード、最大512ノード、0.09 μm)
- SX-9 2007 0.3ns 102.4GF 839TF 8192台
- (3.2GHz、8-16PE/1ノード、最大512ノード、0.065 μm 、11層銅配線)

科学新聞

週刊

(金曜日発行)

発行所 科学新聞社

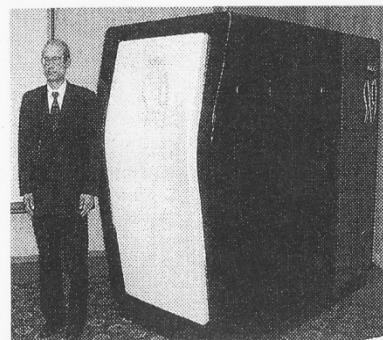
世界最速ベクトル型スパコン

単一コアで100GFLOPS超達成

SXの新モデル投入 NEC

NECは、世界最高速のベクトル型スーパーコンピュータ「SXシリーズ」モデル「SX-9」を製品化して、10月25日から世界で同時発売した。これは、単一チップ当たりで100GFLOPS(1ギガフロップス)毎秒10億回の浮動小数点演算性能)超という極めて優れた演算性能をもった、世界で最速の一チップベクトルプロセッサを実現したもの。1ノード当たり16個のCPUを搭載して演算性能1・6TFLOPS(1テラフロップス)毎秒1兆回の浮動小数点演算性能)を、さらにそれを最大512ノード接続した、総

積販売台数を誇り、特に、気象・気候解析をはじめ航空宇宙、環境、流体解析などで高い評価を得ている。



SX-9(シングルノードシステム)の説明をする丸山執行役員常務

SX-9は、従来のSX-8の13ノード(1・66TFLOPS)に匹敵する性能を、1ノード(1・6TFLOPS)で実現する、驚異的な演算性能をもったスパコンである。

心臓部のプロセッサは、基本的に従来のSXベクトルアーキテクチャを継承。これに演算機の追加、ベクトルパイプライン数増強などのアーキテ

クチャ改良を加えた。

また、これまでの壁を破る高速化、低消費電力化など最先端技術を採用した65nm(ナノ)CMOS、11層銅配線LSIにより3・2GHzという高周波数化を達成して、単一コアチップとしては、世界初の102・4GFLOPSという演算性能、256GBバイト/秒のメモリバンド幅を実現した。同社の丸山好一・執行役員常務は、SX-9製品発表会の席で「例えば気象予測計算分野でSX-9はスカラ型の機能より3・4倍、物性計算分野でも4・7倍優れている」と、ベクトル型SX-9の利点をアピールした。

世界トップ」と説明した。レンタル価格は月額298万円(税込)で、今後3年間で700システムの販売を見込んでいる。すでに受注もあり、20ノードのSX-8Rを導入した大阪大学サイバーメディアセンターが、10ノードのSX-9を来夏に追加導入する。

さらに、設置面積と消費電力において、従来のSX-8の13ノード対し、同等の性能のSX-9の1ノードはいずれも約4分の1と小型・省エネ化している。コストパフォーマンスでは、従来機の約6倍である。同社の西川岳・第一コンピュータ事業本部長は「HPC(ハイ・パフォーマンス・コンピューティング)分野でも、従来の速さだけではなく、最近のECOの観点も重視されるようになってきており、世界のスパコンランキングトップ500では補完的リストとして、世界で最もエネルギー効率のよいスパコンランキングが発表されているが、SX-9は電力性能比なら



Empowered by Innovation

NEC

NEC SUPERCOMPUTER
SX SERIES MODEL

SX.9

世界最速性能で、未知への扉を開く。

ベタフロップスコンピューティングを視野に入れたスーパーコンピュータ **SX9** 登場!



果てしない宇宙の深淵を探る。生命科学の謎を解明する。

ナノテクノロジーの新たな領域を切りひらく。

スーパーコンピュータは、いまや社会の発展になくてはならない存在です。

もっと、高速に、効率的に、そしてもっと使いやすく。

そんな科学の最前線のニーズに応じて、いまハイパフォーマンスマシンが誕生しました。シングルコア性能世界一。

マルチノードではベタフロップスに迫る演算性能を実現。

圧倒的な高速性能と使いやすさを兼ね備えた先進のベクトルスーパーコンピュータSX-9。

世界最速性能*で、サイエンスの未知の扉を開きます。

NEC SUPERCOMPUTER
SX SERIES MODEL
SX-9

SX-9の特長① 容易に超高速で大規模な計算が可能

CPUはプロセッサ当たり100GFLOPSを超える世界最高速の性能と256GB/秒の広いメモリバンド幅を実現。シングルノードシステムでは、1.6TFLOPSを超える演算性能と最大1TB/秒の大規模共有メモリにより、容易に高速演算が可能です。さらに、マルチノードシステムでは、最大512ノードで839TFLOPSの超高速演算性能とノード当たり最大128GB/秒×2のノード間超高速通信により、高いシステムトータル性能が得られ、より大規模・長時間スケールの現象の計算が可能になります。

GFLOPS: Giga Floating-point Operations Per Second (1秒当たり10億回の浮動小数点演算)
TFLOPS: Tera Floating-point Operations Per Second (1秒当たり1兆回の浮動小数点演算)

SX-9の特長② 拡張性の高いシステムを容易に管理

従来のSXシリーズで磨きあげられてきたスーパーコンピュータ用オペレーティングシステムSUPER-UXを提供。大規模マルチノードシステムへの対応などの機能として、リソース管理、ジョブ管理、チェックポイント・リスタートなど様々な機能を提供し、最大512ノードの構成においても容易なシステム管理が可能です。また、SX-9の性能を最大限に発揮するためのコンパイラを提供しており、総合プログラム開発環境を利用した効率的なソフトウェア開発が可能です。

SX-9の特長③ 高密度実装技術と低消費電力技術によるTCOの削減

1チップベクトルプロセッサとメモリモジュールを、高密度実装技術によりコンパクトな装置にパッケージング。最先端の低消費電力回路設計技術を用い、電力効率に優れたベクトルプロセッサの性能当たり電力をさらに改善。このため、設置環境・消費電力も大幅に改善されています。

SX-9の特長④ SXシリーズの応用ソフトウェア資産継承

従来のSXシリーズでサポートしているさまざまなアプリケーションソフトウェアを、そのまま継承可能。SX-9の超高速演算性能を活用できます。

* 単一コア当たり1秒間の最高演算性能 (2012年11月時点)

あらゆるニーズに応える 新アーキテクチャ

最先端のLSIテクノロジー、及び高密度実装技術によりSX-9は世界で初めて100GFLOPSを超えた世界最速の1チップベクトルプロセッサ。ノード演算性能1.6TFLOPS、及びメモリ容量1TBの超大規模共有メモリアーキテクチャによるシングルノードシステム、そして最大128GB/秒の超高速インターコネクトにより512ノードを接続し、839TFLOPSの総合演算性能を有する共有・分散メモリアーキテクチャの超大規模マルチノードシステムを実現。NECはSX-9を通して、あらゆるHPCニーズにお応えいたします。

世界最速の 1チップベクトルプロセッサ

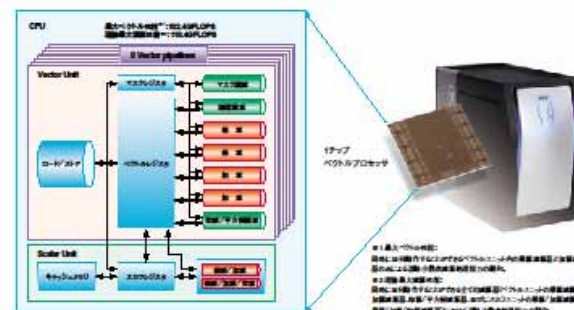
● 超高速演算を可能にした強力なプロセッサ
HPCにおいて並列度が増加することは、プログラミングの難しさを引き起こすだけではなく、並列処理におけるスケラビリティが得られなくなる問題があります。「使いやすい」、かつ「高性能」なスーパーコンピュータをご提供し続けるために、NECは単一チップの演算性能にこだわります。SX-9は100GFLOPS超という驚異的な演算性能を誇る世界最速1チップベクトルプロセッサを実現。これにより、演算処理において必要となる「並列度」を低く抑えることを可能とし、



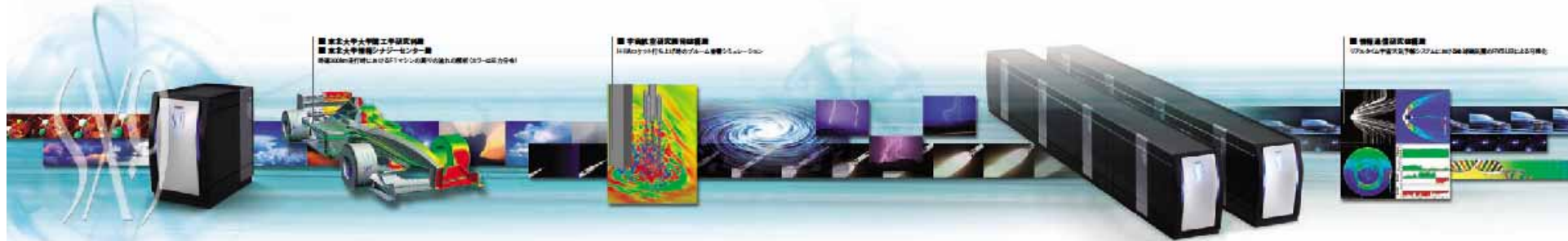
単に理論演算性能が高いだけではなく、実務において「使いやすい」高い実効性能を実現します。演算処理において高い実効性能を実現するためには、プロセッサの演算性能とバランスのとれたメモリバンド幅も重要な要素となります。SX-9は、1チップ当たり256GB/秒のメモリバンド幅を実現しています。

● 超高速プロセッサを実現する ベクトルアーキテクチャ・テクノロジー

最先端技術の結晶ともいえるスーパーコンピュータSX-9。その心臓部であるプロセッサは、基本的に従来のSXベクトルアーキテクチャを継承しつつ、演算器の追加、ベクトルバイパス増強などのアーキテクチャ改良を行っています。また、これまでの「壁」を打破する超高速、低消費電力化の各種最先端テクノロジーを採用した65nm CMOS、11層銅配線LSIにより3.2GHzという高周波数化を実現。これにより、単一コアチップとしては世界初となる102.4GFLOPSの演算性能、及び256GB/秒のメモリバンド幅を実現しています。



もっと使いやすく、さらに高性能に。世界最速1チップベクトルプロセッサが、未知の世界を解明する。



ハードウェア Hardware

高度なシミュレーションが求められるあらゆる分野で、その能力をいかんなく発揮するために、NECが求めたのは、「使いやすさ」と「高性能」の両立。この課題を世界最速1チップベクトルプロセッサの開発によって解決しました。単一チップで100GFLOPS超の演算性能と256GB/秒のメモリバンド幅が生む、使いやすく高い実効性能が、科学の探求を強力に支援します。

● これまでの「壁」を軽々と超える最新テクノロジー
HPCのみならず、計算機システムが現在直面している壁、すなわち「メモリバンド幅の壁」、及び「消費電力の壁」、スーパーコンピュータには、これらの壁を打破することが求められています。SX-9はメモリアーキテクチャの改良、及び最先端テクノロジーによりメモリバンド幅の壁を打破し、256GB/秒という驚異的なメモリバンド幅を実現しています。また、電力効率の高い10Gbps-SerDes、Multi-Vib、クロックゲーティング、チップ内電力センサを用いた電力制御などの最先端低消費電力化技術を採用することにより、高性能かつ低消費電力のプロセッサを実現。

電力効率の高いスーパーコンピュータ環境をご提供いたします。

使いやすいと高性能を併せ持つ 大規模SMPノード

SX-9シングルノードシステムは、最大16CPU(ノード演算性能1.6TFLOPS)、最大1TBの共有メモリを装備するAモデルと、最大8CPU(ノード演算性能800GFLOPS)、最大512GBのBモデルを用意し、ニーズに合わせて使いやすく、かつ高性能なSMPによるHPC環境を実現します。SX-9は従来のSXシリーズとの互換性を確保しつつ、演算性能、メモリバンド幅、メモリ容量、

入出力性能などのトータルバランスを追求。実績のあるベクトル型スーパーコンピュータの資産をスムーズに継承・移行することが可能です。

スケラブルなHPCを実現する マルチノードシステム

SX-9マルチノードシステムは、シングルノードシステムを構成要素として、片方向最大128GB/秒/秒という超高速ノード間接続装置(DXS)によりノード間を接続。最大512ノード、8192台のプロセッサによるマルチノードシステム構成が可能です。超高速な1チップベクトルプロセッサ、共有メモリアーキテクチャによるシングルノード、及び超高速

ノード間接続装置による大規模な共有・分散メモリシステムによりスケラブルなHPC環境を実現します。SX-9マルチノードシステムは、512ノードの最大構成時においても、シングルシステムイメージでの運用環境を提供。512ノード、8192台のプロセッサによる共有・分散メモリシステムは839TFLOPSを実現し、超大規模、かつ使いやすいHPCシステムをご提供いたします。

高い演算性能と バランスのとれた高いI/O性能

使いやすく、かつ高性能なシステムであるためには、演算性能とバランスのとれた高いI/O性能

も重要な要素です。SX-9はシングルノードあたり最大64GB/秒という超高速の総合I/O性能を有しています。

多彩なRAS機能による高信頼性

SX-9は、高集積設計により部品点致を大幅に削減。これによりハードウェアの信頼性を格段に向上させています。さらに、NECのメインフレーム開発で培われた各種高信頼性技術、高可用性技術を搭載することにより、信頼性、可用性の向上を図っています。メモリには誤り検出訂正符号(ECC)を採用。回路部には二重化などに

よる誤り検出機能を組み込んでいます。また、各装置内にエラー箇所を検出するビルトイン診断機能(BID)を備え、速やかな故障箇所の指摘や回復・再構築処理を実行します。さらに、障害情報の自動収集、サービスセンターへの自動通報や、センタからの遠隔保守により、迅速な故障診断と容易な予防保守を実現。システムの信頼性・可用性・保守性を総合的に高めています。

周辺装置やネットワーク環境については、SX-9ホームページをご覧ください。
<http://www.nec.com/jp/hpc/sx9/product/hardware.html>

パワフルで柔軟なSUPER-UX

オペレーティングシステムとしてUNIX System Vに準拠した「SUPER-UX」を提供。高速入出力機能、通信機能など、スーパーコンピュータにふさわしい機能を備えています。

● マルチプロセッサ/マルチノードサポート

マルチプロセッサと並列処理をさらに強化してサポート。また、マルチノードシステムにおいてもシングルスレッドシステムイメージ (SSI) を実現。より使いやすく快適な利用環境を提供します。

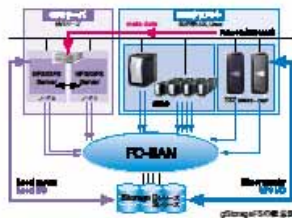
● リソースブロック機能

システムの持つ多数のCPUや大規模なメモリ

資源の分割管理が可能。特にノード内では、分断損のまったく発生しない、効率のよい資源管理を自動的に行うことができます。

● gStorageFS

gStorageFS (GFS) は、複数ノード、あるいは異機間において、高速にファイルを共有できる機能を提供します。GFSを利用することにより、GFSクライアントでもローカルディスクに近い高性能のアクセスが可能です。ユーザはGFSを意識することなくNFSと同様に利用できます。SXシリーズ、TX7シリーズ (OS: Linux) やEM64T (OS: Linux)、SGI Altix などとの間でもファイル共有が可能です。また、GFSサーバのgStorage NVゲートウェイ



シリーズは、ファイルサーバ化して完全な二重化構成をとっており、高信頼、高可用な運用を提供します。

強力なバッチ処理環境 / ジョブスケジューリング機能

● NQSII

さまざまなシステム構成に対応可能な次世代バッチシステムです。クライアント機能、運用管理機能、ジョブ実行機能を物理的に分離することでSSIを実現。統一したユーザインタフェース、管理機能の一元化によるシンプルな運用管理、ジョブが使用可能な計算リソースの最適化が行えます。

● ジョブスケジューリング機能

ジョブのスケジューリングを行うバッチスケジューラはNQSII本体とは独立に構成され、柔軟なジョブスケジューリングを実現しています。基本的なスケ

ジューリング機能を持つ標準スケジューラは、投入順に処理を行うFCFS (First Come First Served) 方式のスケジューリングを提供します。

● JobManipulator

JobManipulatorは、計画的なリソース管理機能によってシステム稼働率の最大化を実現する高機能バッチスケジューラです。ジョブ実行に必要な計算リソースを計画的に割り当て、占有利用を可能にするバックファイル・スケジューリングを提供すると共に、ユーザ、グループ、および組織単位に、公平なジョブの優先度制御を実現するフェアシェアスケジューリング機能や、ジョブの実行開始時刻と必要リソース量を保証する事前予約機能をサポートしています。

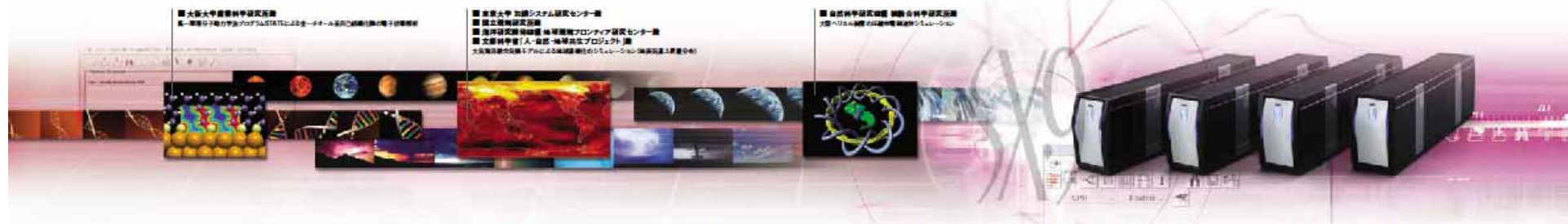
柔軟で快適な運用管理を支援

● チェックポイント・リスタート機能

実行中のプログラムを任意の時点で中断させ、後に再開させることができます。これにより、緊急ジョブの最優先実行や長時間ジョブの分断実行など、柔軟な運用が可能です。マルチノードシステムではNQSIIにより、マルチノードMPIプログラム

のチェックポイント・リスタート機能および実行中のジョブを他の資源の空いているノードへ移すマイグレーション機能が利用できます。

ハードウェアのポテンシャルを最大限に引き出す柔軟なシステム運用と充実したプログラム開発ツール。



ソフトウェア Software

さまざまな分野で求められる高速化処理に応えるために、オペレーティングシステムとして業界標準に準拠した「SUPER-UX」を搭載。大規模マルチノードシステムへの対応、効率をアップさせる運用管理など、柔軟な機能を備えています。また、ベクトル処理、並列処理のすべてのプログラミングモデルに対応。ライブラリ、ツールなども充実し、最適な開発・利用環境を提供します。

先進の言語処理体系とソフトウェア開発環境

SX-9の性能を最大限に引き出すためのさまざまなライブラリやツールを提供しています。

● ベクトル化・並列化の中核コンパイラ

FORTRAN90/SX、C++/SX

FORTRAN90/SXは、従来のSXシリーズで実績を増ってきた高度な最適化、自動ベクトル化、自動並列化機能を装備し、SXシリーズ向けに最適設計されたコンパイラです。最新のFortran95仕様をフル実装し、SXシリーズのベクトル化・並列化の中核を担います。C++/SXコンパイラは、C言語およびC++言語のISO規格をサポートするC/C++

コンパイラです。FORTRAN90/SXと共通のバックエンドをもち、高度な自動ベクトル化、自動並列化および最適化機能を提供します。

● 共有メモリ並列処理機能OpenMP

コンパイラの強力な自動並列化機能に加え、OpenMP APIにより、通常のプログラムから容易に並列処理を利用できます。FORTRAN90/SX、C++/SXでOpenMP API Ver2.5を提供します。

● メッセージ・パッシング・インターフェースライブラリMPI/SX、MPI2/SX

MPI/SXおよびMPI2/SXは、シングルノード内であれば共有メモリの特長を活かし、また、マルチノードシステムであればIXSのハードウェア性能を

最大限に引き出すことで、低レイテンシ、高スループットのデータ転送を実現しています。MPI/SXは、MPI-1.2仕様の機能を提供し、MPI2/SXでは、MPI-2仕様の機能を提供します。

● データ・パワフル言語HPF/SX V2

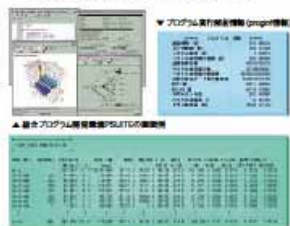
HPF/SX V2は、HPF2.0基本仕様を準拠し、主要なHPF公認拡張仕様、HPF/JA拡張仕様をサポートします。

研究者のためのパワフルなプログラム開発ツール

● GUIベースのプログラム開発環境PSUIITE

SXシリーズでは、スーパーコンピュータと組みあ

接続されたネットワーク上で、GUIによる編集、翻訳、実行、デバッグ、チューニング作業が可能な統合プログラム開発環境PSUIITEを提供します。



● 性能解析ツールprognif、ftrace機能

FORTRAN90/SXやC++/SXが標準で備える機能で、環境変数やコンパイルオプションによって、プログラム全体やサブルーチン単位、ユーザ指定範囲の性能情報を出力することが可能で。

NEC提供アプリケーションソフトウェア

● リアルタイム可視化システムRVSLIB

流体解析や構造解析などの大規模並列シミュレーション実行中に、計算途中結果をSXシリーズ上でリアルタイムに可視化するシステムです。

● 数学ライブラリ

豊富な機能を持つ数値計算ライブラリASL、統計計算ライブラリASLSTAT、ならびにBLAS/LAPACK/ScalLAPACKなどを含むライブラリ集MathKisanを提供しています。いずれもSXシリーズ向けに最適化されています。

さらに詳しいソフトウェア開発・利用環境についてはSX-9ホームページをご覧ください。
<http://www.nec.co.jp/hpc/sx9/product/software.html>

HPC

High Performance Computing

SX Series

English

トピックス

イベント情報

事例紹介

お問い合わせ

SX シリーズ

SX-9

SX-8R

SX-8i

ソフトウェア情報

NECアプリケーション

サポートソフトウェア

TX7 Series

Express5800 Series

HPC販売推進本部



人類は未来を切り拓く力を手に入れた。

TOP

装置諸元

多彩な応用分野

ハードウェア

ソフトウェア

カタログダウンロード

装置諸元

Specifications

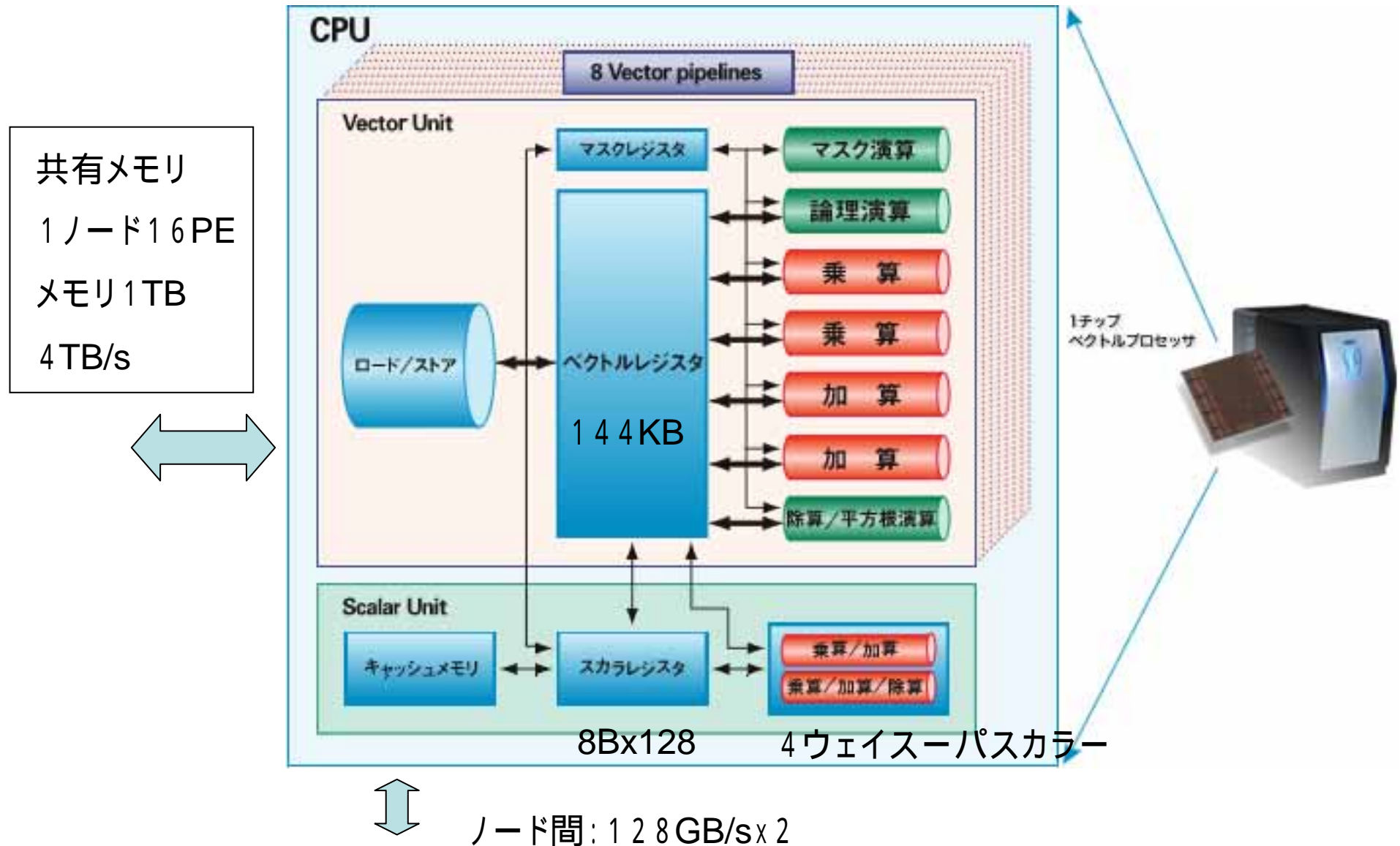
シングルノードシステム

モデルグループ	SX-9	
モデル名	A	B
中央処理装置 (GPU)		
CPU台数	8~16	4~8
理論最大演算性能※1	947.2G~1,894.4GFLOPS	473.6G~947.2GFLOPS
最大ベクトル性能※2	819.2G~1,638.4GFLOPS	409.6G~819.2GFLOPS
ベクトルレジスタ	144Kバイト×8~16	144Kバイト×4~8
スカラーレジスタ	64ビット×128×8~16	64ビット×128×4~8
主記憶装置 (MMU)		
メモリアーキテクチャ	共有メモリ	
容量	512Gバイト、1Tバイト	256Gバイト、512Gバイト
最大転送性能	4Tバイト/秒	2Tバイト/秒
入出力機構		
最大チャンネル数	32チャンネル※3	16チャンネル
最大データ転送性能	128Gバイト/秒※3	64Gバイト/秒

マルチノードシステム

モデルグループ	SX-9
ノード数	2～512※4
中央処理装置(CPU)	
CPU台数	32～8,192
理論最大演算性能 ※1	3.8T～969.9TFLOPS
最大ベクトル性能※2	3.3T～838.9TFLOPS
ベクトルレジスタ	144Kバイト×32～8,192
スカラーレジスタ	64ビット×128×32～8,192
主記憶装置(MMU)	
メモリアーキテクチャ	共有・分散メモリ
容量	1T～512Tバイト
最大転送性能	2,048Tバイト/秒
入出力機構	
最大チャンネル数	16,384チャンネル※3
最大データ転送性能	64Tバイト/秒※3
ノード間接続装置	
最大データ転送性能	128Gバイト/秒×2(双方向)/ノード

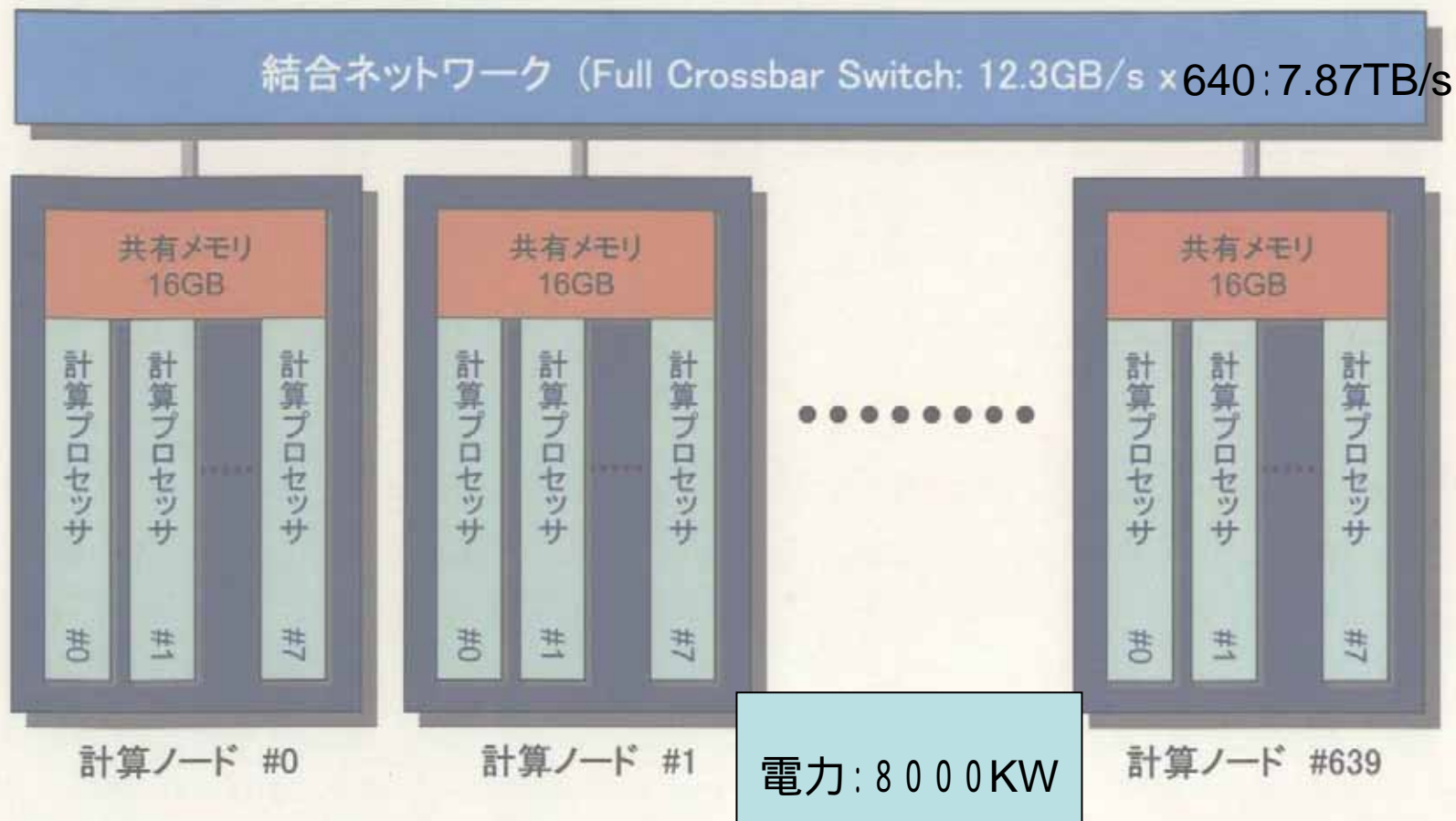
SX-9 1プロセッサ:102.4GF 最大839TF 8192台
(3.2GHz、8-16PE/1ノード、最大512ノード、0.065 μ m、11層銅配線)



地球シミュレータ

地球シミュレータの全体構成

- 総計算ノード数: 640
- ピーク性能: 40TFLOPS
- 主記憶容量: 10TB
- 総プロセッサ数: 5120
- 計算プロセッサのピーク性能: 8GFLOPS
- 計算ノードのピーク性能: 64GFLOPS
- 計算ノードの主記憶容量: 16GB



計算プロセッサ(AP)の構成

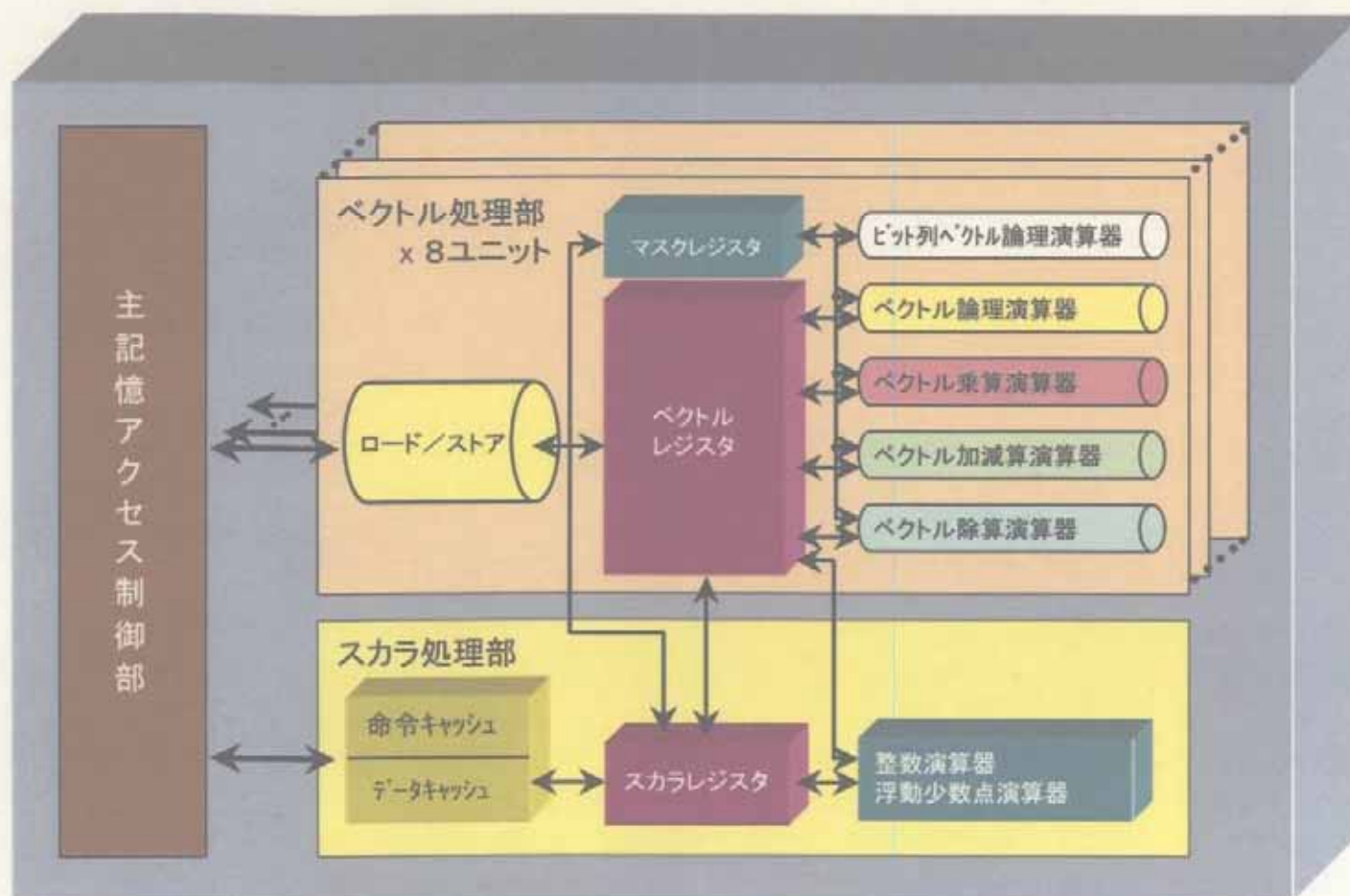
○ ベクトルユニット:8セット

- ◆ 6種のベクトルパイプライン
- ◆ 256要素のベクトルレジスタ: 72個
- ◆ 256ビットのマスクレジスタ: 17個

○ 主記憶アクセス制御部

○ スカラユニット

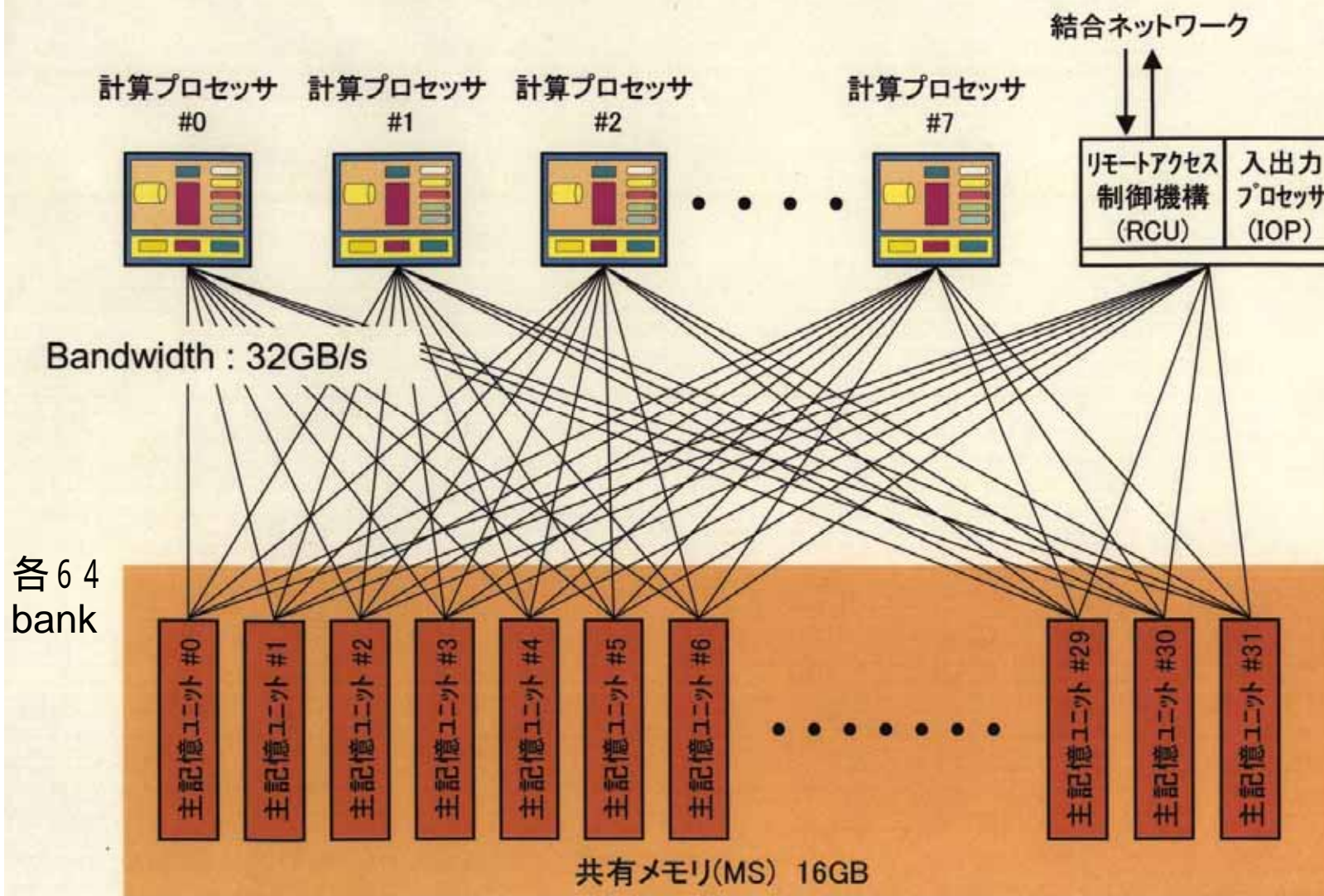
- ◆ 4-ウェイ スーパースカラ
- ◆ 64KB 命令キャッシュ
- ◆ 64KB データキャッシュ
- ◆ 128個の汎用レジスタ



1チップLSI: 8Gflop

- ◆ 0.15 μ m CMOSテクノロジ + 銅配線
- ◆ 20.79mm x 20.79mm
- ◆ 5,700万トランジスタ
- ◆ 5185 ピン
- ◆ クロック周波数
500MHz(1GHz)
- ◆ 消費電力
135W(Typ.)

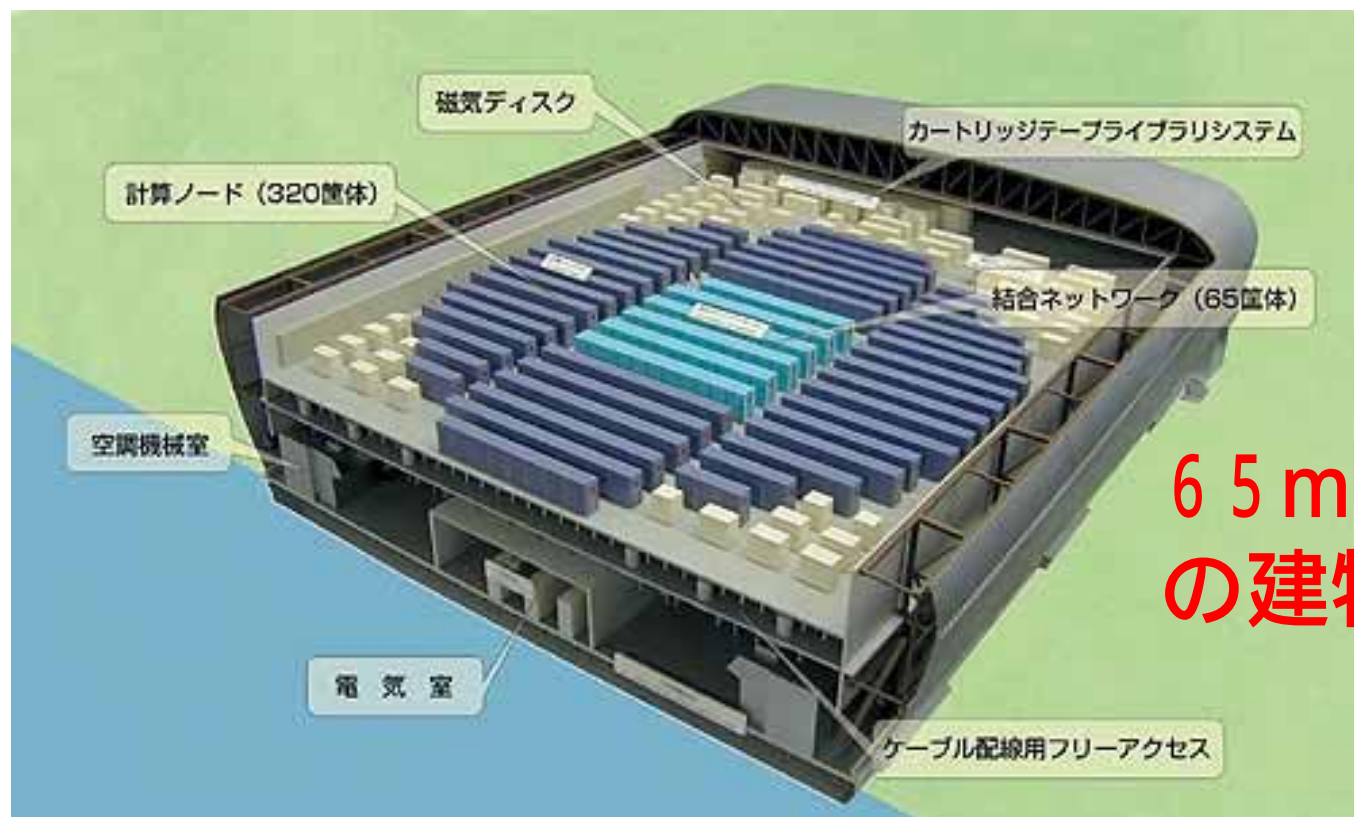
計算ノード(PN)の構成



ESRDC@JAERI

全体で 2048 バンク、24nsec/バンク

「地球シミュレータ」とは、どんなコンピュータか



65m * 50m
の建物

8台のスーパーコンピュータからなる計算ノードを、高速のネットワークで640台つないだものです。

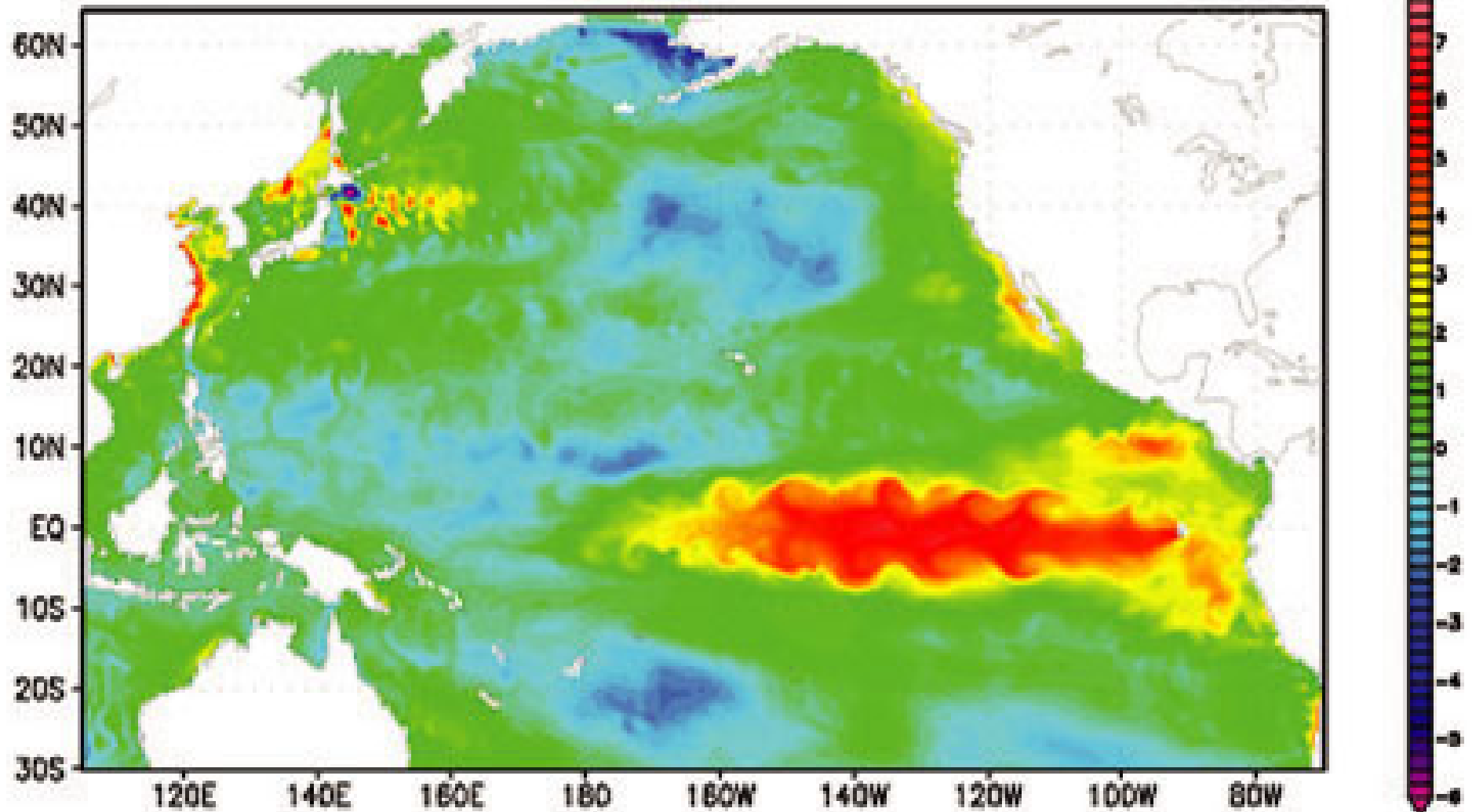
(総計5120個のスーパーコンピュータから構成)

完成時には世界最速のコンピュータになると予想されます。

平成14年3月からの利用開始を目指して開発中です。

地球を10km四方に分割

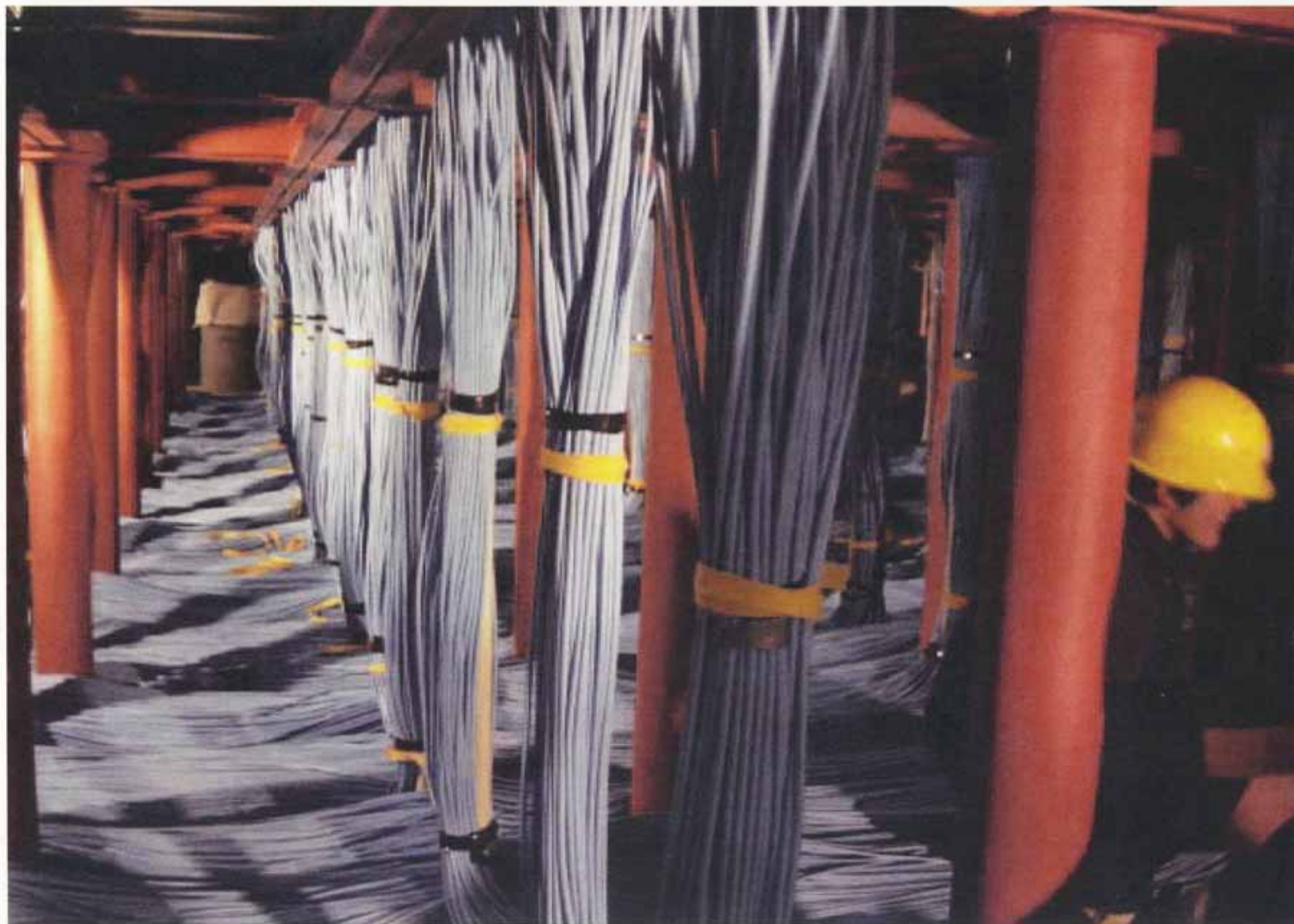
SST diff. between Dec/27/1997 and Dec/30/1984



地球シミュレータ施設（航空写真）



計算ノード・結合ネットワーク間ケーブル敷設作業 (平成13年2月～平成13年5月)





計算ノード・結合ネットワーク間ケーブル敷設完了
(平成13年5月)



地球シミュレータ 設置完了（平成14年1月）



BlueGene/L: IBM

- ・ 2005年稼動予定
- ・ 65,536 プロセッサ
- ・ 360 TFLOPS
- ・ メッセージパッシング
- ・ 3D-トラス: $64 \times 32 \times 32$
適応ルーティング, 仮想チャネル4本
- ・ ブロードキャスト, リダクション: トリー

表 2.5 ASCI プラットフォームの概要

名称	Red	Blue Pacific	Blue Mountain	White	T30
設置研究所	Sandia	Lawrence Livermore	Los Alamos	Lawrence Livermore	Los Alamos
メーカー	Intel	IBM	SGI	IBM	未定
使用 MPU	9,536×Pentium II Xeon	5,856×Power PC	Origin2000 MIPS R10000	8,192×Power 3-II	未定
目標性能	1.8Tflops メモリ 606GB Disk 容量 40TB	3.1Tflops メモリ 2.6TB Disk 容量 75TB	3.1Tflops メモリ 2.5TB Disk 容量 75TB	10.2Tflops メモリ 2.5TB Disk 容量 75TB	30 + Tflops
実績	3.2Tflops (’99/10 月)	3.9Tflops (’98/10 月)	3.1Tflops (’98 年)	— — —	— — —

(注意) 性能はピーク性能値である。

米 計算機のスパコン 日 最速でしのぎ

用途拡大で躍起

毎秒35兆回↓360兆回:1000兆回へ

科学技術計算に利用されるスーパーコンピュータ(スパコン)で、日米間の世界最速争いに拍車がかかっている。スパコンはヒトのゲノム(全遺伝情報)を活用した医薬研究などに用途が広がってきた上、国防上も重要な役割を果たすためだ。最速マシンを日本製に奪われた米国側がナンバーワン奪回へ躍起となっている。

「日本は科学技術計算」での講演で世界最速を
で新時代を切り開いたと誇る日本の「地球シミュ
とで称賛されるべきだ」レータ」(ES)を引き
が、米国が新時代に遅れ 合いにし、最速の座を取
てはならない」。エー リ戻すと強調した。
ブラハム米エネルギー長 二〇〇二年に完成した
官は今月十日、ワシントン ESは、海洋科学技術セ

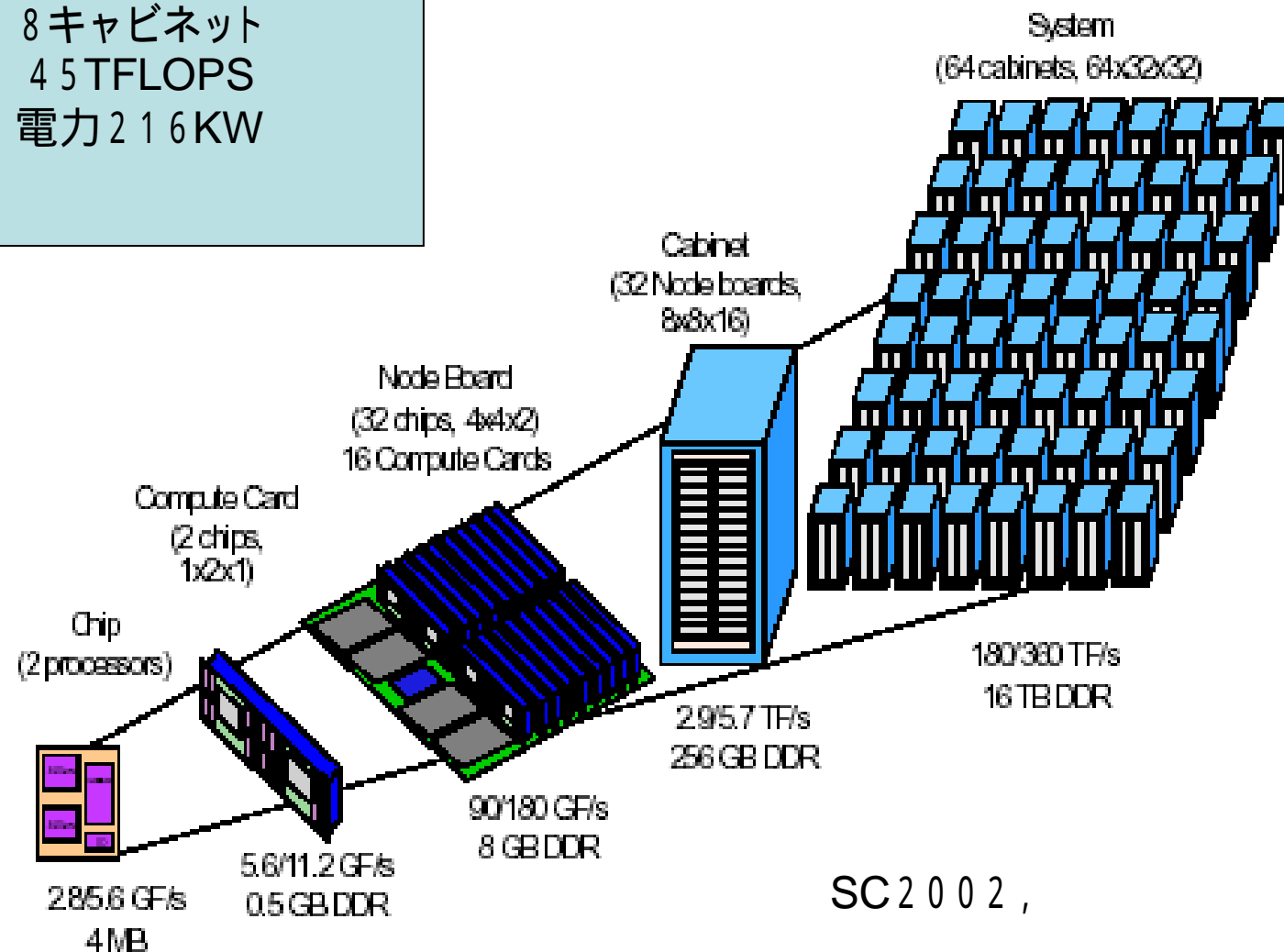
ンターなどが開発し、N 値で毎秒三五・八六〇兆
ECが製造に当たった。 回の計算能力を持つてい
平和利用を目的とする用 る。それまでの最速は、
途は気候変動予測、地殻 核兵器の備蓄管理などを
変動の解明などで、実測 目的に米国で用いられる

IBM製マシンの毎秒七 〇兆回規模の計算能力を
・二二六兆回だっただけ 備えたスパコン開発計画
に「衝撃的な数字」(米 を一九九九年から推進
ニューヨーク・タイムズ 中。〇一年には計画の一
紙だ。 環として、米エネルギー

米国の科学者らが今月 省と共同で「ブルー・ジ
中旬にまとめた最新のス ン(青い遺伝子/L」
パソコン上位五百機リス トと呼ぶスパコン開発を打
ち出している。目指す能 力は毎秒三百六十兆回。
首位を堅持、二位の米ヒ ューレット・パッカード
既に小型試作機が完成
製マシンの三倍近い能力 し、IBMは最新ランキ
となっている。 ングで試作機が七十三位

ただ、IBMは遺伝子 に入ったとアヒール。〇
からつくられる、たんば 五年の完成時にはトップ
く質の構造解析などに利 になる」と首位奪還を予
用するため、毎秒一〇〇 告している。

8 キャビネット
4.5 TFLOPS
電力 21.6 KW



SC2002,

www.sc-conference.org/sc2002/

PowerPC 440

2コア、内1つは
通常通信に使用

Figure 1: BlueGene/L packaging.

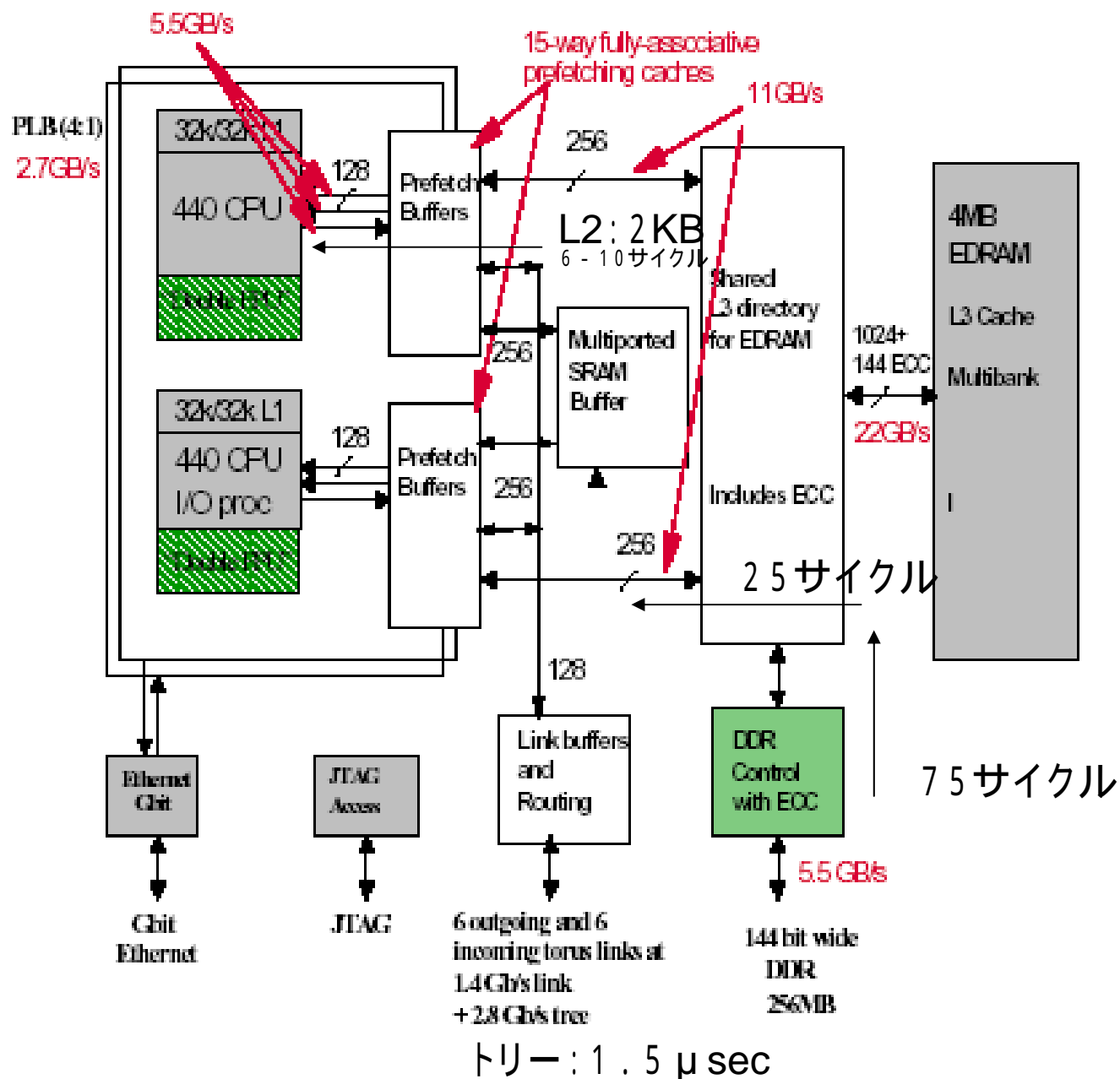


Figure 3: BlueGene/L node diagram. The bandwidths listed are targets.

目指せ世界最速タイトル奪還

次世代スパコン開発へ

2005.7.25
京都新聞

文科省

文部科学省は二十五日までに、最大演算速度が一〇ペタフロップス（一秒間に一京回＝一京は一兆の一万分）の次世代スーパーコンピュータ「京速計算機システム」の設計、開発に着手することを決めた。

二〇〇六年度概算要求に研究費数十億円を盛り込む。総事業費は八百億～一千億円に上る見込み。一〇

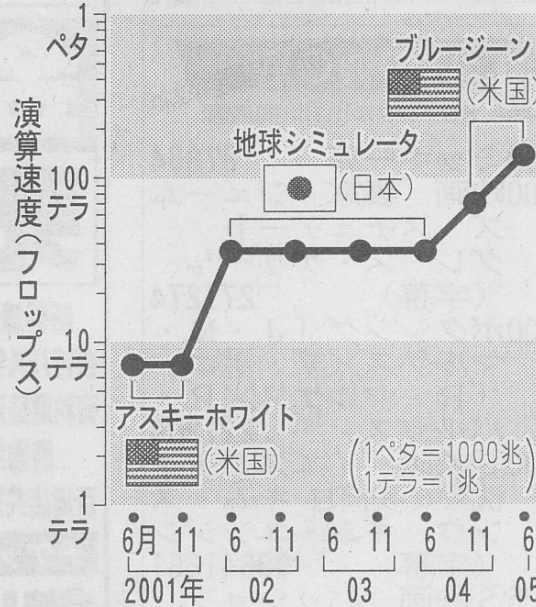
1秒1京回

年完成を目指し、米国のス京速計算機が完成すれば、パソコンが持つ世界最速のタイトル奪還に挑む。この約七十三倍の処理速度を持つことになる。

また、京速計算機システムの運用を担う「先端計算科学技術センター（仮称）イータ」は、〇四年までの設立方法や建設場所について調査研究も始める。現在の世界最速スパコンは、米ローレンスリバモア国立研究所のスパコン「ブルー・ジュン」で、一三六・八テラフロップス（一秒間に約百三十六兆八千億回）は激化している。米国も一〇年をめどに数

ハテフロップス（一秒間に約百三十六兆八千億回）は激化している。競争

世界最高速スパコンの変遷
(独マンハイム大などが年2回調査、対数グラフ)



フロップス コンピューターの計算方法の一つである「浮動小数点演算」を、一秒間に何回実行できるかを示す単位。コンピュータの処理速度の目安として使われる。例えば、1ギガフロップスなら1秒間に10億回、1テラフロップスなら1秒間に1兆回の浮動小数点演算能力を持つ。

次世代スーパーコンピュータ (京(けい)速コンピュータ)

地球シミュレータの250倍の性能: 10 PFLOPS

プロセッサはどうする

ネットワークはどうする

省電力はどうする: 30 MW?

使いやすいソフトウェアはどうする

アプリケーションはどうする

お金は: 1100億円(5年間)

開始: 2006

完成: 2012

場所: 神戸ポートアイランド、2007.3.決定

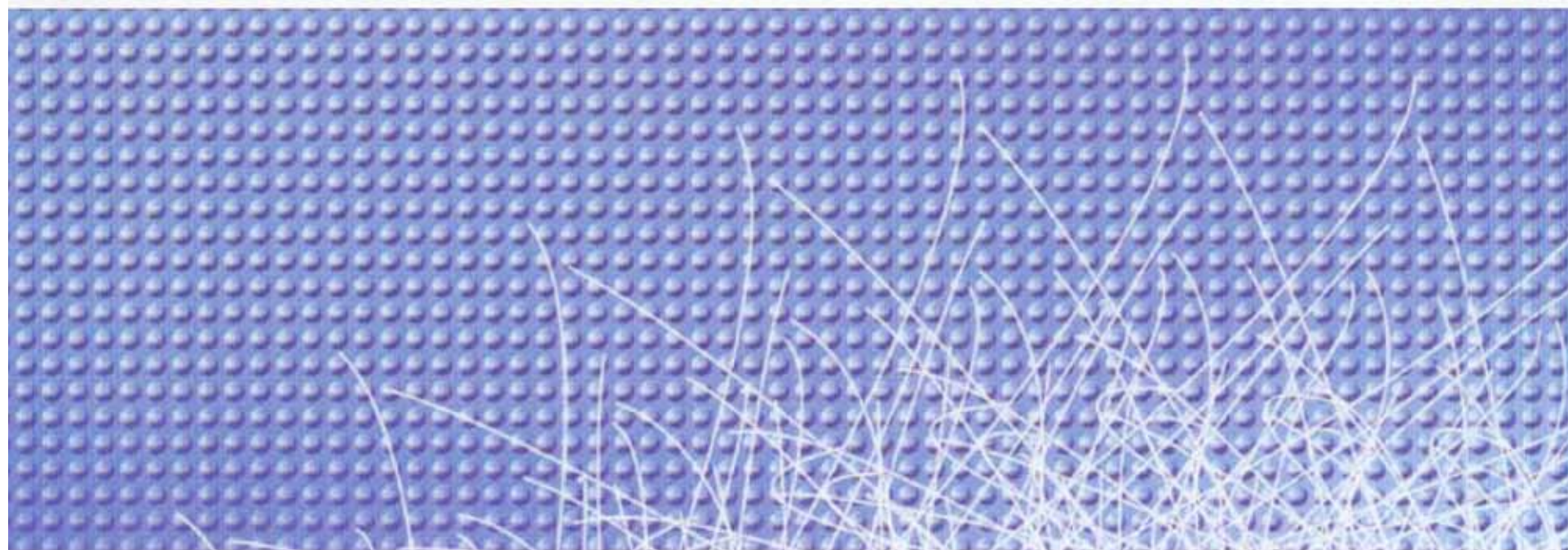
体制: 理化学研究所統括

富士通 + 日立 + NEC

次世代スーパーコンピュータの 開発と利用



独立行政法人 理化学研究所



広範な分野での利活用
-次世代スーパーコンピュータが
拓く世界-

スーパーコンピュータによるシミュレーションは、実験、理論と並ぶ研究開発の第3の手法として益々重要になっています。スーパーコンピュータは、自動車や飛行機の設計・製作のための構造解析や流体解析、天気予報のための気象シミュレーションなど様々な分野で使われており、今後の科学技術と産業の発展に不可欠です。

ものづくり

- 自動車の安全設計等
- 迅速な新製品開発



ナノテクノロジー

- 新物質設計
- 酵素・触媒反応の解明



防災

- 地震動予測
- 津波被害の予測



航空・宇宙

- ロケットエンジン設計
- 航空機開発



ライフサイエンス

- 新薬開発
- 治療・診断技術開発



地球環境

- 気候変動予測
- エルニーニョ現象の影響予測



原子力

- 原子力施設丸ごと解析
- 核融合炉の開発



天文・宇宙物理

- 宇宙の創成過程解明
- 銀河や惑星の形成過程解明



スーパーコンピュータの性能は急速に進歩しています。

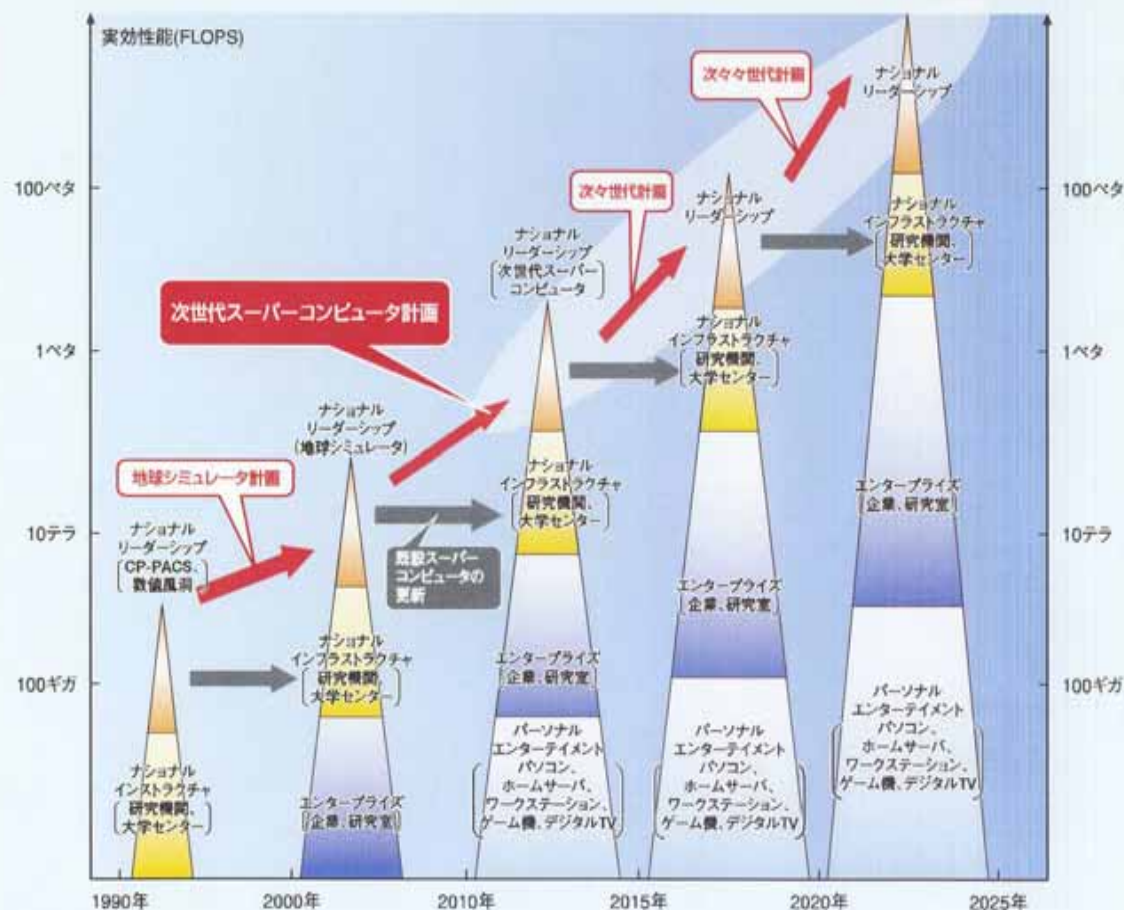
この30年間でスーパーコンピュータの計算能力は200万倍以上にもなっています。スーパーコンピュータの開発には半導体技術、光通信技術、ネットワーク技術、品質管理技術など、エレクトロニクスに関する総合的かつ高度な技術力が必要です。これらは家電や情報通信など多くの分野で国民生活と産業に不可欠な技術の基盤となります。

●スーパーコンピュータの性能推移

ピーク性能 (FLOPS: 1秒間の浮動小数点演算回数)



これまで、国家プロジェクトで開発された最先端のスーパーコンピュータシステム(ナショナルリーダーシップシステム)の技術は、全国の計算機システムに展開されてきました。更に、スーパーコンピュータを将来に亘って継続的に開発することにより、我が国の研究開発基盤の強化と技術の維持向上に大きく貢献することが期待されます。



国の総合科学技術会議は、我が国の科学技術及び産業の競争力の基盤となる「次世代スーパーコンピュータの開発・利用」プロジェクトを「国家基幹技術」と位置づけ、集中的に投資して推進することとしています。

プロジェクトの概要

理論、実験と並び、現代の科学技術の方法として確固たる地位を築きつつあるスーパーコンピューティングを更に発展させるため、長期的な国家戦略を持って取り組むべき重要技術(国家基幹技術)である「次世代スーパーコンピュータ」を2012年の完成を目指して開発します。

今後とも我が国が科学技術・学術研究、産業、医・薬など広範な分野で世界をリードし続けるべく、

- ①世界最先端・最高性能の「次世代スーパーコンピュータ^(注)」の開発・整備 (注) 10ペタFLOPS級
- ②次世代スーパーコンピュータを最大限利活用するためのソフトウェア(=グランドチャレンジ・アプリケーション)の開発・普及
- ③次世代スーパーコンピュータの共同利用と学術情報ネットワークを介した全国のスーパーコンピュータの重層的な利用環境の整備(=サイバー・サイエンス・インフラストラクチャの整備)
- ④次世代スーパーコンピュータを中核とする世界最高水準のスーパーコンピューティング研究教育の拠点の形成

を文部科学省のイニシアティブにより、開発主体である理化学研究所を中心に産学官の密接な連携の下、一体的に推進します。

【開発スケジュール(予定)】

		2006年度	2007年度	2008年度	2009年度	2010年度	2011年度	2012年度
						稼働▲	完成▲	
システム	演算部	概念設計 / 詳細設計 / 試作・評価 / 製造・据付調整						
	制御フロントエンド (トータルシステムソフトウェア)		基本設計	詳細設計	製作・評価		性能チューニング・高度化	
	共有ファイル		基本設計	詳細設計	製造・据付調整			
ソフトウェア プラットフォーム	次世代ナノ統合シミュレーション	開発・製作・評価					実証	
	次世代生命体統合シミュレーション		開発・製作・評価					実証
施設	計算機棟		設計	建設				
	研究棟		設計	建設				
運用		方針・体制の検討				準備活動	運用	

次世代スーパーコンピュータの開発・整備

理化学研究所は、プロジェクトの中核機関として、スーパーコンピュータ開発をリードする最高水準の汎用システムである「次世代スーパーコンピュータ」を開発・整備します。

【システム構成概要】

【世界最速のシステム】

- 1秒間に1京(ケイ=10の16乗)回の計算性能=10ペタFLOPS級

【汎用システム】

- 科学技術・産業で用いられる多様なアプリケーションやこれまで不可能だった複雑かつ大規模なシミュレーションが実行可能

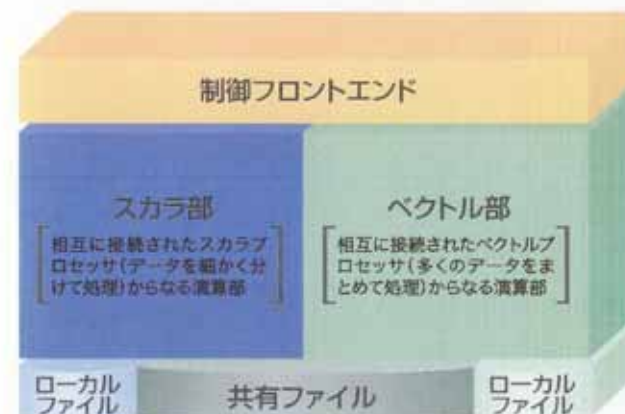
【革新的なシステム】

- 先端技術の積極的導入により、画期的な省電力、省スペースを実現。理化学研究所とメーカー3社(富士通、日本電気、日立製作所)との共同により、日本の技術の総力を結集して開発

【システムの基本的な構成】

多様なシミュレーションの実行に最適な計算環境を提供する複合汎用スーパーコンピュータシステム

- シミュレーションの特性に応じた最適な演算部で計算を実行
- スカラ部とベクトル部の併用により、従来困難だった複雑かつ大規模なシミュレーションも実行可能



【建設地】

- 兵庫県神戸市中央区港島南町7丁目(ポートアイランド第2期内)
ポートアイランド南駅より徒歩約1分(JR新神戸駅から約25分、神戸空港から約10分)



ライフサイエンス分野

次世代生命体統合シミュレーションソフトウェアの研究開発

【概要と目標】

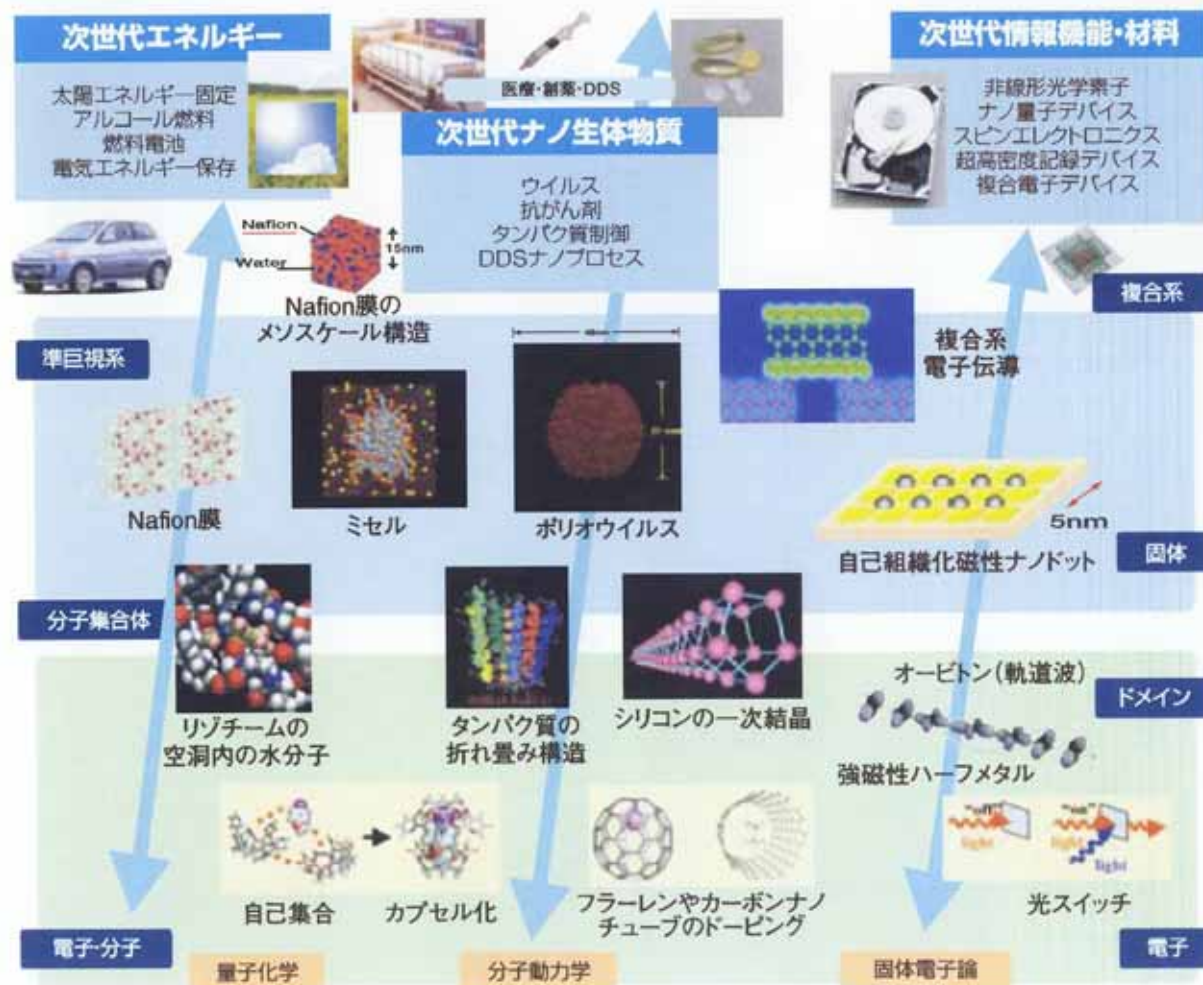
理化学研究所は、次世代スーパーコンピュータの性能を十分に発揮させ、分子から全身まで生体内で起こる種々の現象を統合的に理解するためのシミュレーションソフトウェアの研究開発を進めます。シミュレーションによる予測を基礎とした研究手法をライフサイエンスの分野に提供し、生命現象の統合的な理解による科学的進展、産業応用等を通じたヘルスサイエンスへの貢献を目指します。



ナノテクノロジー分野

次世代ナノ統合シミュレーションソフトウェアの研究開発

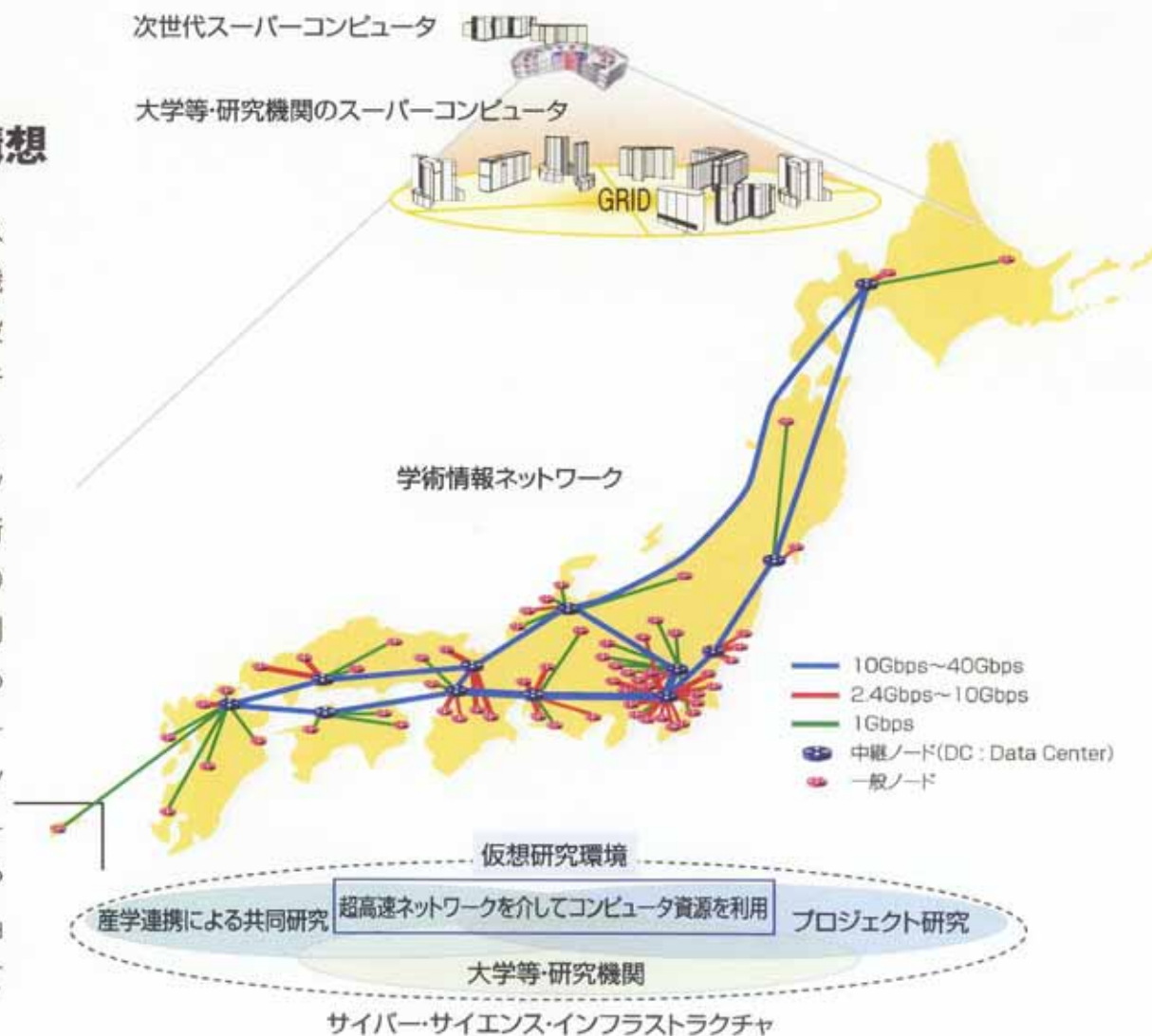
分子科学研究所は、最先端の知的ものづくりを実現する新材料開発のためのアプリケーション・ソフトウェアの開発を行う拠点です。開発されたソフトウェアが次世代スーパーコンピュータで十分に性能を発揮するよう、理化学研究所のハードウェア開発と密接な連携を図っています。



利用環境の整備

サイバー・サイエンス・インフラストラクチャ構想

サイバー・サイエンス・インフラストラクチャ構想とは、大学や研究機関が有しているコンピュータ等の設備、基盤的ソフトウェア、学術コンテンツ及び学術データベース、人材、研究グループそのものを超高速ネットワーク上で共有する最先端学術情報基盤であり、大学等との連携の下、国立情報学研究所や全国共同利用情報基盤センター等が中核となって推進しています。次世代スーパーコンピュータは、サイバー・サイエンス・インフラストラクチャ上で、スーパーコンピュータを保有する大学や研究機関等と連携しながら、ナショナルリーダーシップシステムとしての機能を発揮していきます。



「最先端・高性能汎用スーパーコンピュータの開発利用」 プロジェクトの実現に向けて（案）

文部科学省研究振興局 2005 年 10 月 26 日発表

開発主体：理化学研究所

平成 18 年度-平成 22 年度

1154 億円（平成 18 年度 40 億円）

米国：2009 年 1PFLOPS コンピュータ開発予定

科学新聞 2005 年 11 月 4 日

平成 17 年 8 月 10 日

文部科学省研究振興局

文部科学省ホーム
ページより

最先端・高性能汎用スーパーコンピュータの開発利用(案)

目的: 世界最先端・最高性能のスーパーコンピュータ「汎用京速計算機」システムの開発・整備及び利用技術の開発・普及

趣旨及び効果: 理論、実験と並び、現代の科学技術の方法として確固たる地位を築きつつあるスーパーコンピューティング(シミュレーション(数値計算)やデータマイニング、解析等)について、今後とも我が国が世界をリードし続けるため、

(1) スーパーコンピュータを最大限利活用するためのソフトウェア等の開発・普及

(2) 世界最先端・最高性能の汎用京速^(注)計算機システムの開発・整備

(注)京速=10ペタFLOPS

(3) 上記(2)を中核とする世界最高水準のスーパーコンピューティング研究教育拠点(COE)「先端計算科学技術センター(仮称)」の形成により研究水準向上と世界をリードする創造的人材の育成を総合的に推進。

世界最高性能の科学技術計算環境を実現し、複雑で多様な現象の系全体のシミュレーションや高度なデータマイニング、解析等を、幅広い分野で行い、「知的ものづくり」や「科学的未来設計」を実問題で可能とし、先端的スーパーコンピューティングにおける国際的なリーダーシップを確立。科学技術・学術や産業の競争力強化、安全・安心な社会の構築に貢献。

また、世界の英知を結集し、世界水準の人材育成を行い、シミュレーションにおける我が国の国際的な地位を確立する。

概要: 平成18年度は、世界最先端・最高性能の汎用京速計算機システムの開発・整備の前提であるシステム全般の設計・研究開発等に着手する。

1. ソフトウェア(OS、ミドルウェア、アプリケーションソフトウェア)等の設計・研究開発
2. ハードウェア(計算機システム及び超高速インターコネクション)の設計・研究開発
3. 「先端計算科学技術センター(仮称)」の形成に関する調査研究

体制: 国の責任で設備の整備から運用まで一体的に推進する。また装置の開発・運用を行うに当り、産学官の様々な組織から最も適したところを選択し、そのポテンシャルを活用する。

事業期間: 平成18年度～24年度

先端計算科学技術センター(仮称)を
平成22年度末までに開所



スーパーコンピューティング研究教育拠点(COE)の形成

科学技術・学術の発展と産業競争力強化に貢献
(以下を例に、様々な科学技術・学術・産業分野を対象)

材料～製品丸ごと設計



ナノ分野

生命体シミュレーション



ライフ分野

自動車開発



ものづくり分野

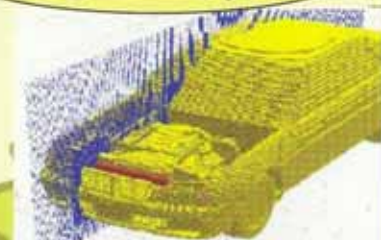
津波影響予測



防災分野 等

広汎な分野での利活用 - 次世代スパコンが拓く世界 -

ものづくり



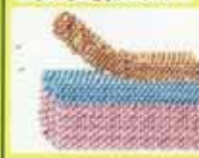
自動車開発

提供: 日産自動車(株)

ナノテクノロジー

物質設計

触媒設計



提供: (独)物質・材料研究機構



提供: (独)物質・材料研究機構

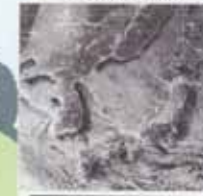
防災

津波被害予測



提供: 東北大学

雲の解析



提供: 気象研究所

原子力



原子炉
丸ごと解析

提供: 日本原子力研究所



レーザー
反応解析

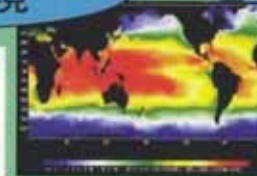
提供: 日本原子力研究所

ライフサイエンス



地球環境

エルニーニョ
現象の影響
予測



提供: (独)海洋研究開発機構

天文・宇宙物理

銀河形成解明

惑星形成解



提供: (独)理化学研究所



提供: 国立天文台

オーロラ
発生解明



提供: (独)海洋研究開発機構

航空・宇宙

ロケット
エンジン設計

航空機開発



提供: (独)宇宙航空研究開発機構



提供: (独)宇宙航空研究開発機構

先端計算科学
技術センター
(仮称)

汎用京速計算機が目指すグランドチャレンジ（例）

世界最高水準の科学技術創造立国を実現するため、国際競争力を支える新産業創造等の政策目標の実現をも視野に入れ、ナノテクノロジー／ライフサイエンス分野を革新する汎用京速計算機のグランドチャレンジを明示して戦略的に研究開発を進める。

<ナノテクノロジー分野アプリケーション>

次世代ナノ統合シミュレーション

電子・原子・分子から、ナノスケールの分子複合デバイスに至るまで、ナノ材料を丸ごと解析することにより、次世代ナノ材料（新半導体材料等）の創出などの実現を目指す。

<ライフサイエンス分野アプリケーション>

次世代生命体統合シミュレーション

遺伝子から全身の血流まで、人体丸ごと解析することにより、テーラーメイド医療や創薬などの実現を目指す。

研究開発スケジュール（案）

年度		平成17年度	平成18年度	平成19年度	平成20年度	平成21年度	平成22年度	平成23年度	平成24年度
開発項目	評価等	★ 研究開発チーム発足	計画本格化判断★ (設計仕様、開発体制、立地・運用方針等)			研究開発状況評価★ (システム性能・機能等)		COE形成、運用評価★ (利用状況、研究成果、人材育成状況等)	
ソフトウェア	システムソフトウェア	NAREGI ^(※4) (平成15年度より)	異機種統合ソフトウェア設計・製作			異機種統合ソフトウェア評価			
			グリッドミドルウェア設計・製作			グリッドミドルウェア評価			
	グランドチャレンジアプリケーション	(※4)	次世代ナノ統合シミュレーション設計・製作			次世代ナノ統合シミュレーション評価			
	革新的シミュレーションソフトウェアの研究開発 ^(※1)		次世代生命体統合シミュレーション設計・製作			次世代生命体統合シミュレーション評価			
	次世代高精度・高分解能シミュレーション技術の開発 ^(※3)		革新アプリケーション評価						
ハードウェア	要素技術開発 ^(※2)								
	大規模処理計算機部		設計			実装技術設計・評価		製作	システム強化
	逐次処理計算機部		設計			実装技術設計・評価		製作	システム強化
	特定処理計算加速部		設計			実装技術設計・評価		製作	
	異機種間接続超高速インターコネクション		設計			実装技術設計・評価		製作	
	遠隔可視化装置					実装設計・評価		製作	
	ファイルシステム					設計		製作	システム強化
その他	立地調査、建屋建設、付帯設備整備等		検討	設計	建設	付帯設備整備			

「最先端・高性能汎用スーパーコンピュータの開発利用」以外のプロジェクトを示す。

プロジェクト部分に該当。

※1:「次世代IT基盤構築のための研究開発」の研究開発領域の一つ。

※2: 科学技術振興機構「戦略的創造研究推進事業」の一戦略目標下の研究領域として、「情報システムの超低消費電力化を目指した技術革新と統合化技術」を設定。

※3: 科学技術振興機構「戦略的創造研究推進事業」の一戦略目標下の研究領域として、「マルチスケール・マルチフィジックス現象の統合シミュレーション」を設定。

※4:「超高速コンピュータ網形成プロジェクト(National Research Grid Initiative)」。平成15年度よりグリッドミドルウェアとナノシミュレーションソフトウェアの開発を進めている。

かってデータフローコンピュータ
という魅惑的な研究があった

1970年代後半から80年代前半
スーパスカラ方式のルーツ

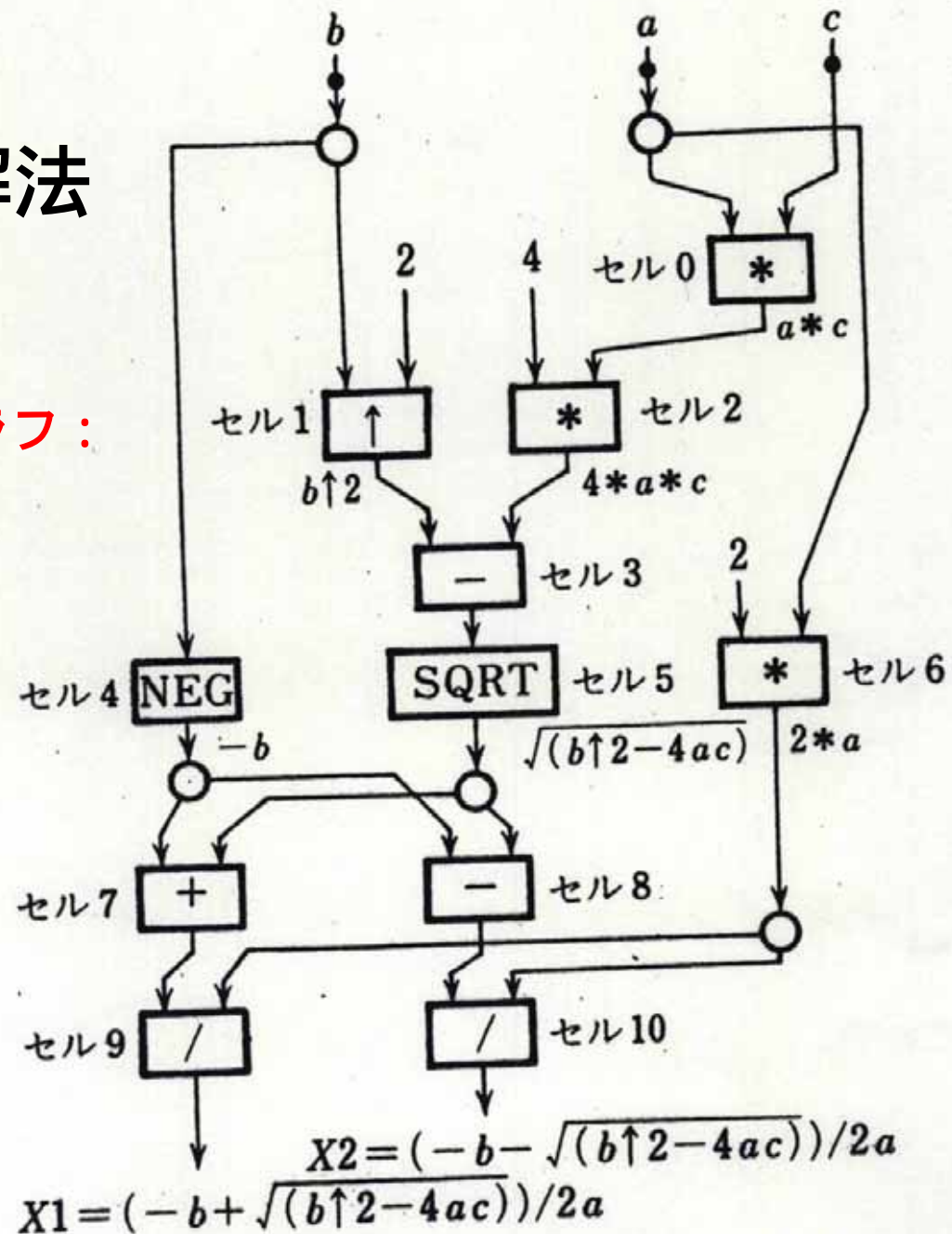
特徴

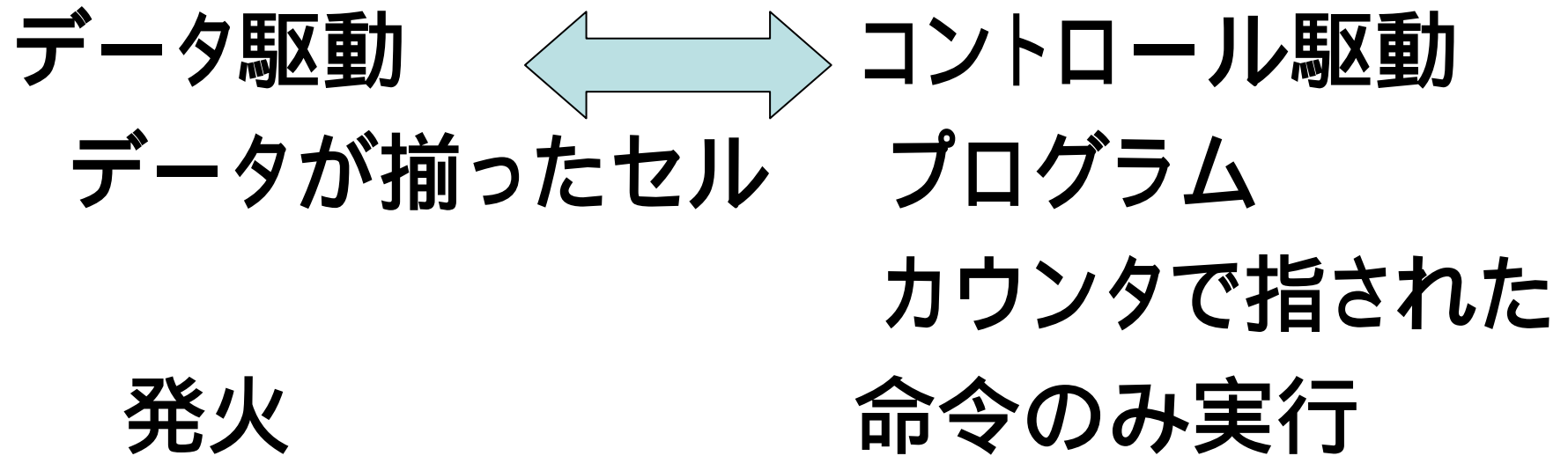
- ・プログラムカウンタがない
- ・プログラムに内在するすべての並列性を実行時に取り出せる

VLIW方式のルーツ: 水平型マイクロ
プログラム

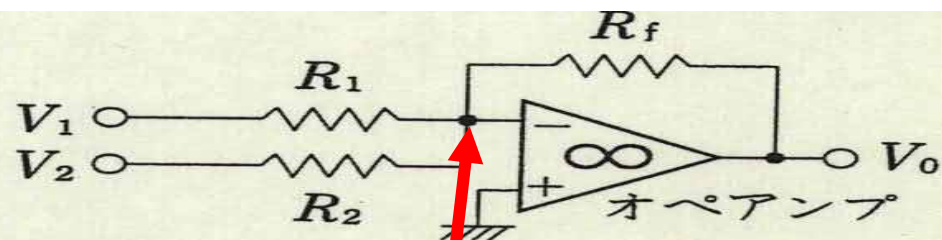
2 次方程式の解法

データ依存を示すグラフ：
データフローグラフ



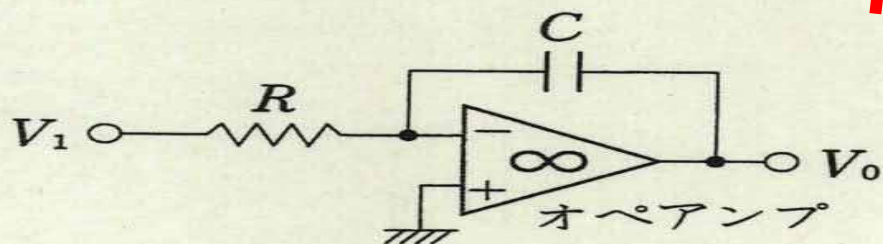


結果を次のセルに渡す

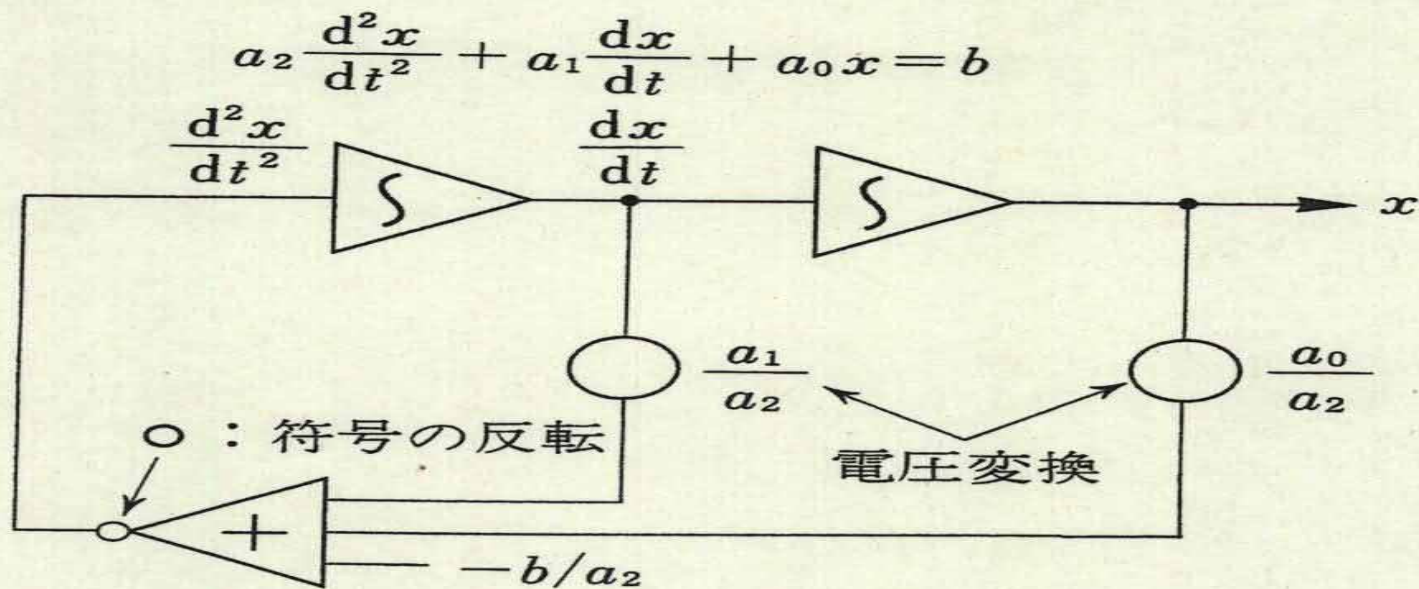


(a) 加算器 $V_0 = -R_f \left(\frac{V_1}{R_1} + \frac{V_2}{R_2} \right)$

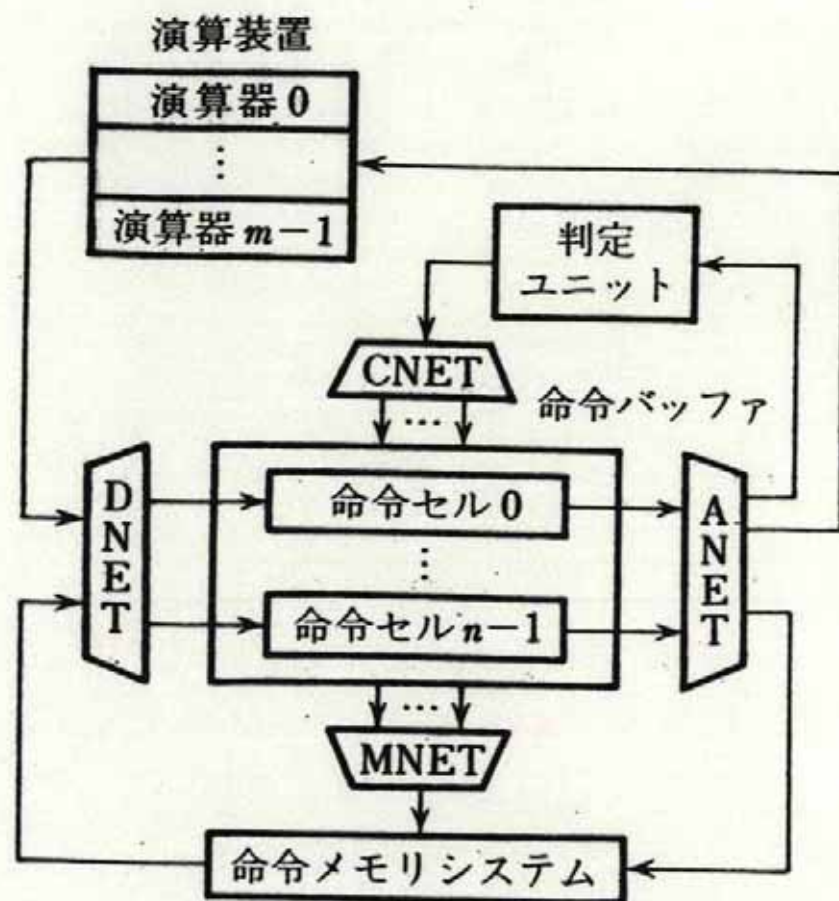
仮想接地: 0 V



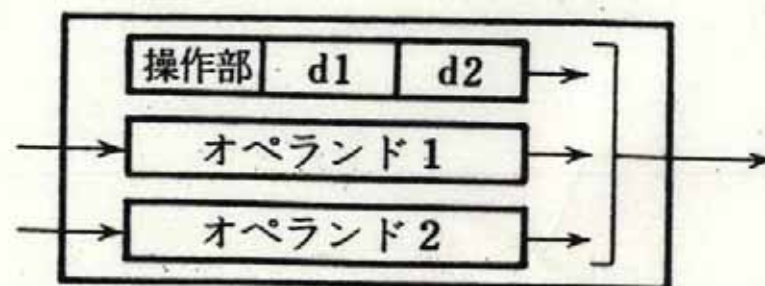
(b) 積分器 $V_0 = -\frac{Q}{C} = -\frac{1}{CR} \int V_1 dt$



(c) 常微分方程式を解く回路



(a) ハードウェア構造. ANET:アービトレーションネットワーク, DNET:ディストリビューションネットワーク, CNET:制御ネットワーク, MNET:メモリコマンドネットワーク



(b) 各命令セルの構造. d1, d2: デスティネーションセルを指定.

図 5.23 データ駆動型計算機の構成 (J. B. Dennis, *et al.*: 'A Preliminary Architecture for a Basic Data-Flow Processor', *Proc. of 2nd Annual Int. Symp. on Computer Architecture*, pp, 126-132 (1975) による)

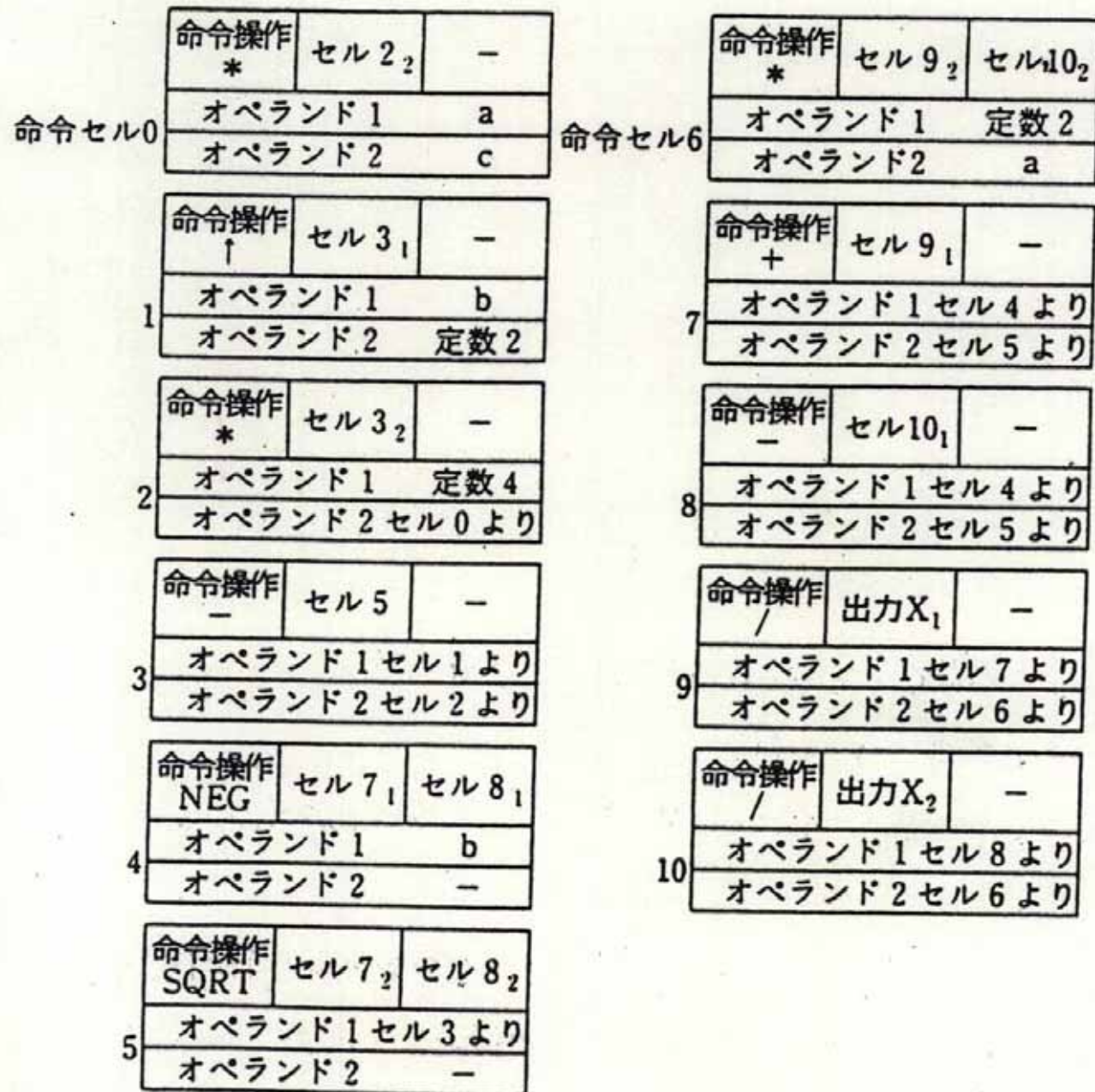
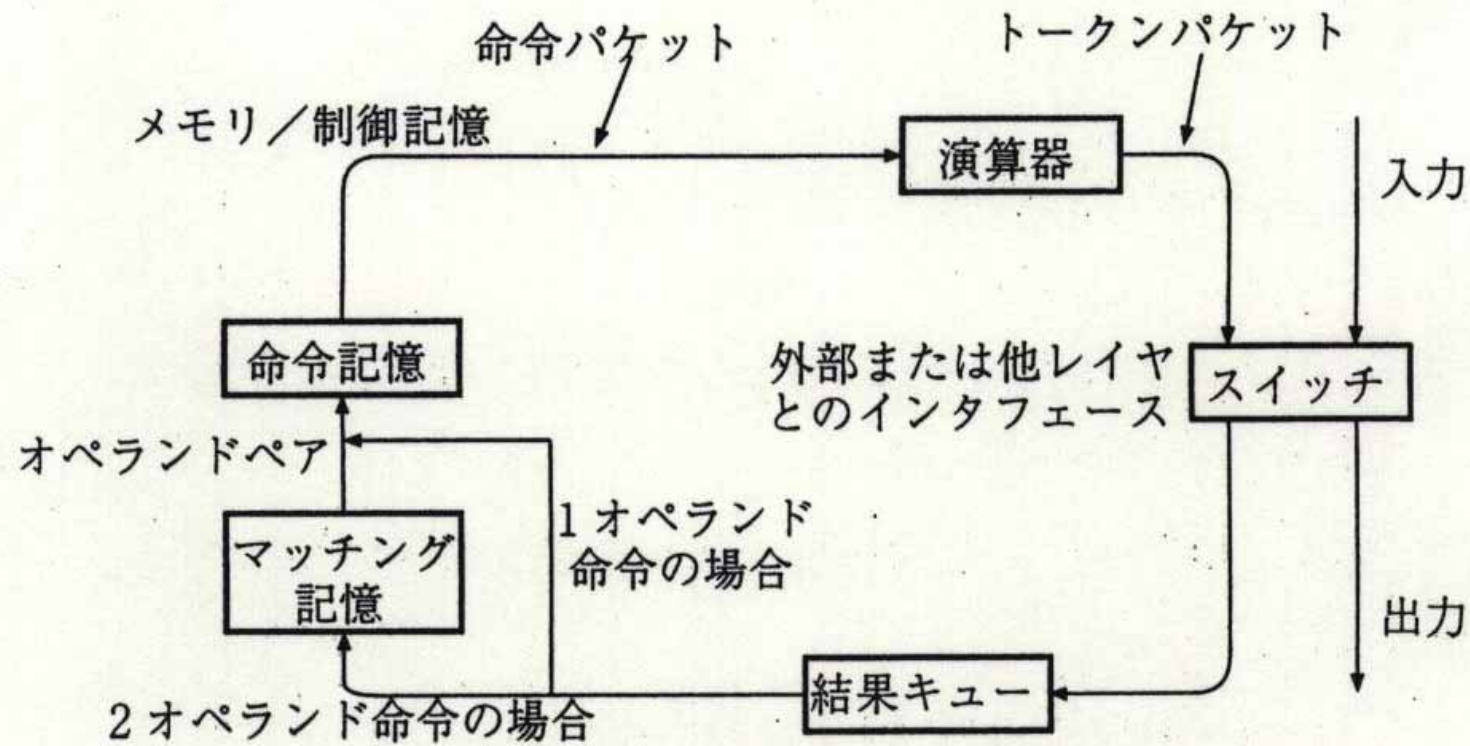


図 5. 24 2 次方程式の根を求めるデータフロープログラムに対するパケット群



命令パケット

タグ	命令	あて先1 (d1)	あて先2 (d2)	オペランド1	オペランド2
----	----	--------------	--------------	--------	--------

トークンパケット

タグ	あて先 (d)	結果
----	------------	----

(a) 循環パイプラインの構造

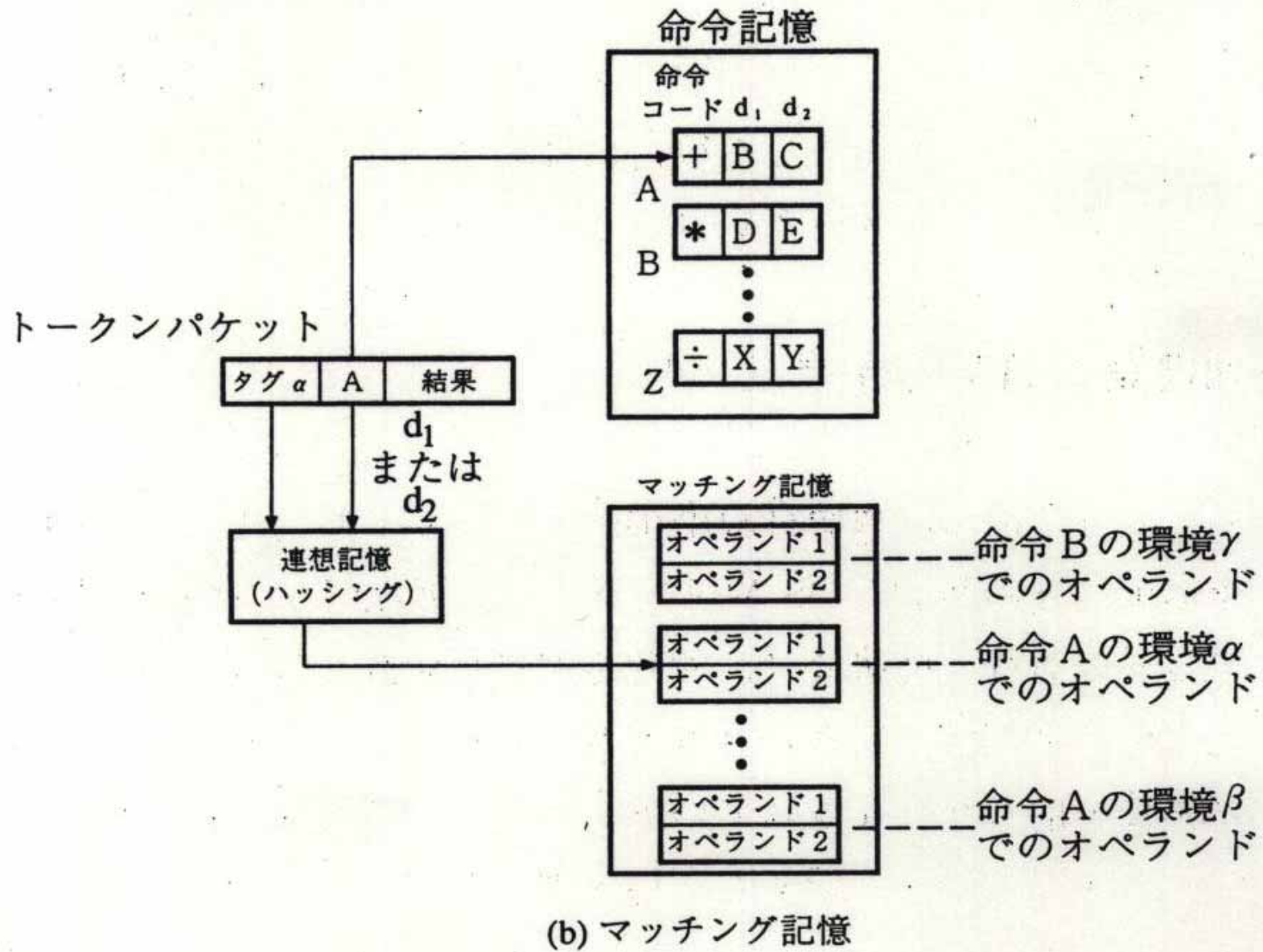


図 5.25 循環パイプライン型データ駆動コンピュータ

(4) データ駆動方式のマルチスレッド化

従来のデータ駆動方式の欠点

細粒度並列処理

メモリバンド幅

マッチング記憶

構造体メモリ

データ駆動方式：処理と通信を統合化

粗粒度化

スレッド起動・切換えの高速化

- ・スレッドごとの専有レジスタ
- ・スレッド起動の高速化
- ・スレッドの多重実行

5 . 7 並列言語と コンパイラ

5 . 7 . 1 並列言語

(1) 逐次言語

(2) 逐次言語の並列処理拡張

ループ並列化の記述

```
DOALL  I=1,N
```

```
A(I+1)=A(I)+1.0
```

```
END  DOALL
```

ベクトル演算 / 配列演算の記述

A、B、Cがベクトルの場合、

$C=A+B$ の表現

プロセス生成・終了、同期操作の指定

FORK-JOINなどのプロセスの生成や終了の
明示

データ分割の指定

ベクトルや配列に対してそのデータ分割
法を明示

通信文の明示的利用

Send/Receiveなどプロセス間での通信を

明示

HPFのアプローチ

Fortranプログラム内でデータ分割をユーザで
指定

所有者計算規則 (owner computes rule) による
目的並列プログラムを生成

SEND/RECEIVE文の挿入

メッセージ交換型マルチプロセッサを対象

Fortran-D、HPF(High Performance Fortran)、
Vienna Fortranなど

MPIのアプローチ

メッセージ通信のプロトコルの標準化

これらを利用したメッセージ交換

各プログラム（プロセス）：Cなど通常の逐次
言語で記述

PVM、MPIやLindaなどが代表例

ネットワーク上のワークステーションにも適用
手軽な並列処理の利用形態

（ 3 ）並列処理言語

データフロー言語のような関数型言語

メッセージ交換モデルをベースとしたOccam

5.7.2 HPF の概要

1992年：HPFF (High Performance Fortran Forum)
で標準化

ユーザが与えるコンパイラへの指示：

ディレクティブ

ディレクティブ：!HPF\$で表示

プロセッサ構成指定

!HPF\$ PROCESSORS P(16)

!HPF\$ PROCESSORS R(4,4)

「仮想」プロセッサ

データ分割指定

配列データの分割：各次元方向に、

BLOCK分割とCYCLIC分割

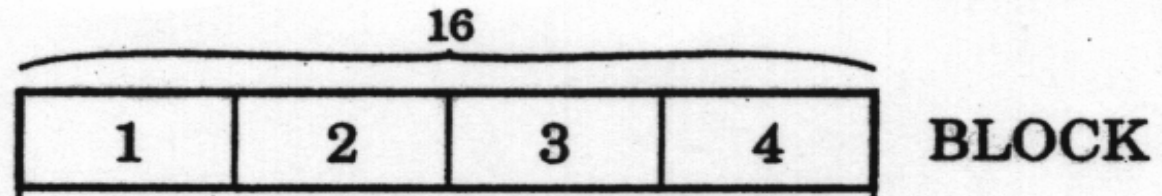
データ相互配置指定

!HPF\$ ALIGN A(I,J) with B(J,I)

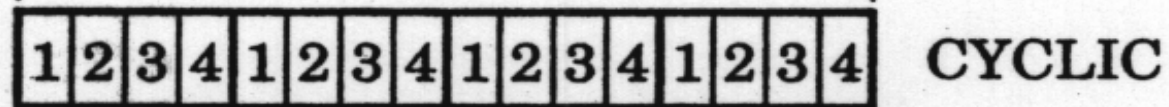
B(J,I)が(BLOCK,BLOCK):データ分割、

A : B の転置行列

B の転置元と A への転置先が同一

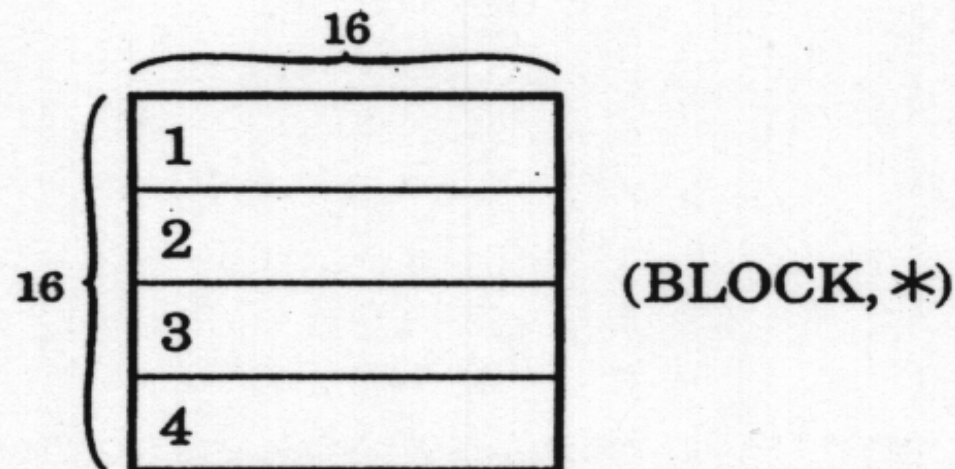


(i) !HPF\$ DISTRIBUTE A(BLOCK) ONTO P(4)



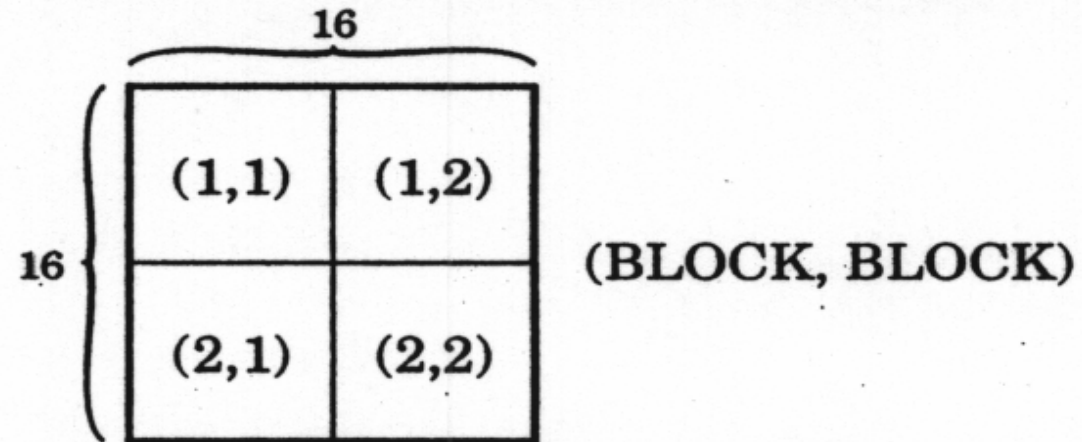
(ii) !HPF\$ DISTRIBUTE A(CYCLIC) ONTO P(4)

(a) 1次元の例

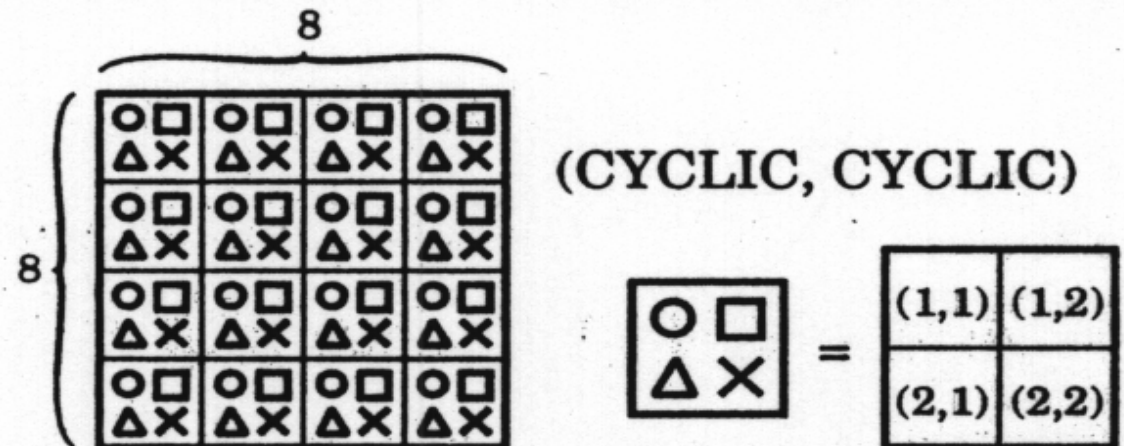


(i) !HPF\$ DISTRIBUTE A(BLOCK, *) ONTO P(4)

(ii) !HPF\$ DISTRIBUTE A(*, BLOCK) ONTO P(4)



(iii) !HPF\$ DISTRIBUTE A(BLOCK, BLOCK) ONTO R(2,2)



(iv) !HPF\$ DISTRIBUTE A(CYCLIC, CYCLIC) ONTO R(2,2)

(b) 2 次元の例

図 5. 26 データ分割の例

プロセッサに配置

`!HPF$ ALIGN A(I) WITH B(I+1)`

$B(I+1)$ と $A(I)$ が同一プロセッサに
配置される。

5.7.3 並列化コンパイラ

(1) ベクトルコンパイラとの差異

ベクトル化コンパイラでは、

データ依存関係の解析

ループ変換

データ依存関係にサイクル構造

ベクトル化できないループ

ループ変換

ループ分割、スカラエクспанション、

ループ交換、ストリップマイニング、

ウェーブフロント

マクロ演算機構

総和演算 / 内積演算

サイクリックリダクション

多重ループの高速化

マルチプロセッサとベクトルプロセッサの

根本的差異

ベクトルプロセッサ

集中メモリ

少数のパイプライン演算器

マルチプロセッサ

分散メモリ方式

データの局所的配置

処理と通信のオーバーラップ実行

一括通信

(2) データ分散配置

DO 10 I=2,100

DO 20 J=1,100

A(I,J)=A(I-1,J)+1.0

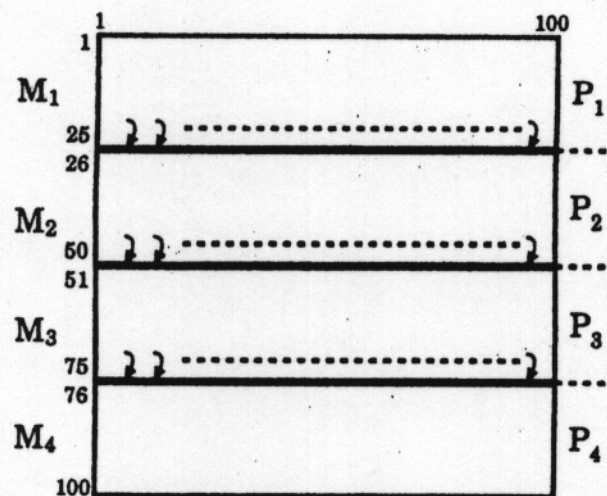
20 CONTINUE

10 CONTINUE

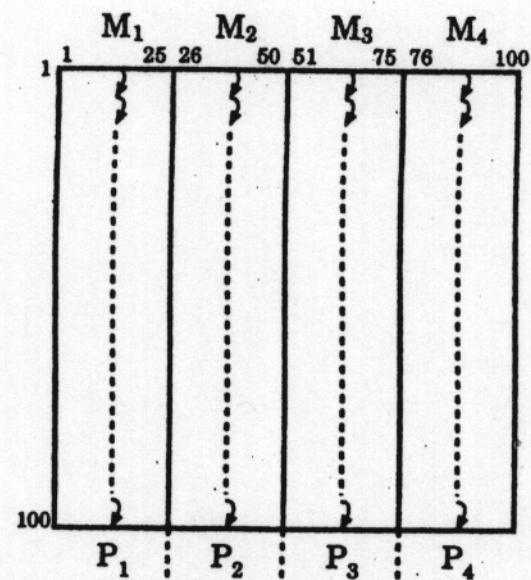
(3) 所有者計算規則

プロセッサ P_i のメモリ M_i に対する代入文(すなわち代入文での左辺のデータ格納先が M_i である)の計算(代入文の右辺の計算)はプロセッサ P_i でなされる。

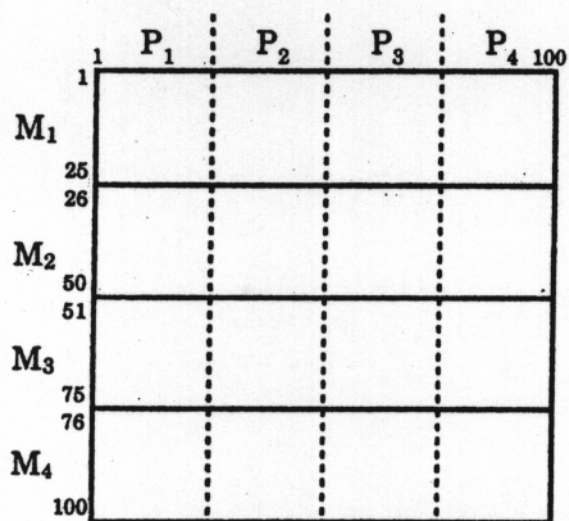
A=B+Cの場合



(a) プロセッサ行分割/メモリ行分割



(b) プロセッサ列分割/メモリ列分割



(c) プロセッサ列分割/メモリ行分割

DO 10 I=1,16

DO 20 J=1,16

A(I,J)=B(I,J)+C(I,J)

20 CONTINUE

10 CONTINUE

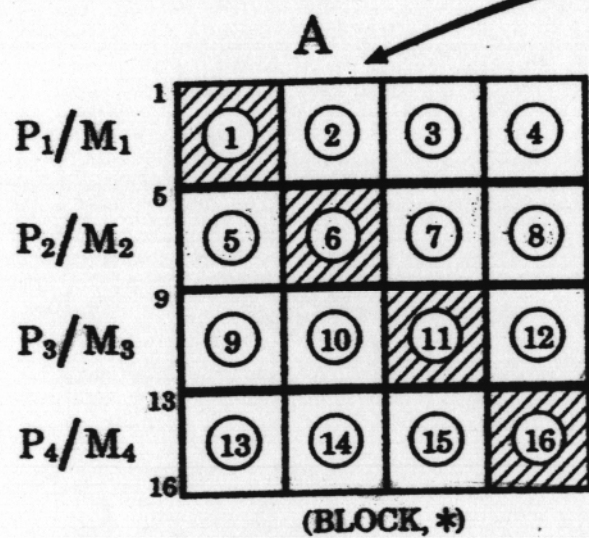
A: (BLOCK, *), B、C: (*, BLOCK) で分割

単純な並列化コンパイラ：

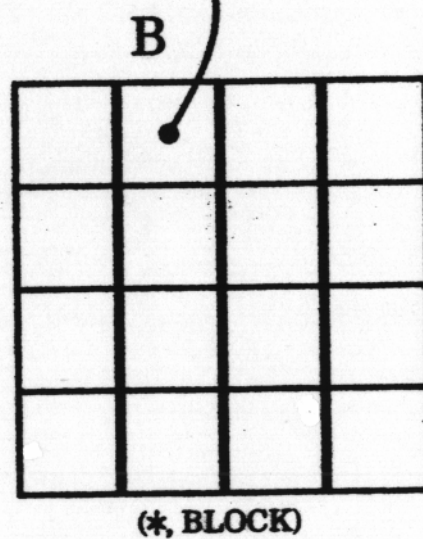
プロセッサP1のプログラム

• B(1:16, 1:4)、C(1:16, 1:4) のP2、P3、P4へのSEND

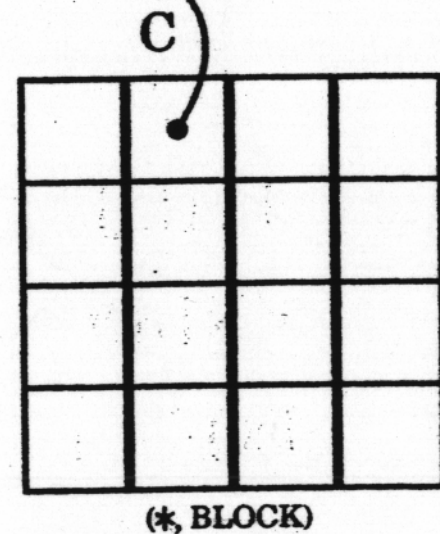
DO J=1,4



▨ : ローカル処理



P_1/M_1 P_2/M_2 P_3/M_3 P_4/M_4



P_1/M_1 P_2/M_2 P_3/M_3 P_4/M_4

D0 I=5,8

SEND P2 B(I,J)、C(I,J)

ENDDO

D0 J=1,4

D0 I=9,12

SEND P3 B(I,J)、C(I,J)

ENDDO

D0 J=1,4

D0 I=13,16

SEND P4 B(I,J)、C(I,J)

ENDDO

- B(1:4、5:16)、C(1:4、5:16)のP2,P3,P4からのRECEIVE

DO I=1,4

DO J=5,8

T1(I,J),T2(I,J)=RECEIVE P2 B(I,J)、C(I,J)

ENDDO

DO I=1,4

DO J=9,12

T1(I,J),T2(I,J)=RECEIVE P3 B(I,J)、C(I,J)

ENDDO

DO I=1,4

DO J=13,16

T1(I,J),T2(I,J)=RECEIVE P4 B(I,J)、C(I,J)

ENDDO

- 自データでの計算

DO I=1,4

DO J=1,4

A(I,J)=B(I,J)+C(I,J)

ENDDO

- RECEIVEデータによる計算

```
DO I=1,4
```

```
DO J=5,16
```

```
A(I,J)=T1(I,J)+T2(I,J)
```

```
ENDDO
```

(4) 通信の最適化

通信ベクトル化(message vectorization)

```
DOALL
```

```
A(I,J)=A(I-1,J)+(I+1,J)+
```

```
A(I,J-1)+A(I,J+1)-4*A(I,J)
```

```
END DOALL
```

通信合体法 (message coalescing)

通信集団化法 (message aggregation)

通信と処理のオーバラップ

検査官 / 執行官方式

たとえば、

$$A(I) = B(b(I)) + C(I)$$

収束演算の場合：

$b(I)$ がしばらくの間変更されないで使用
される場合

検査官の処理は 1 度でよく、毎回データ

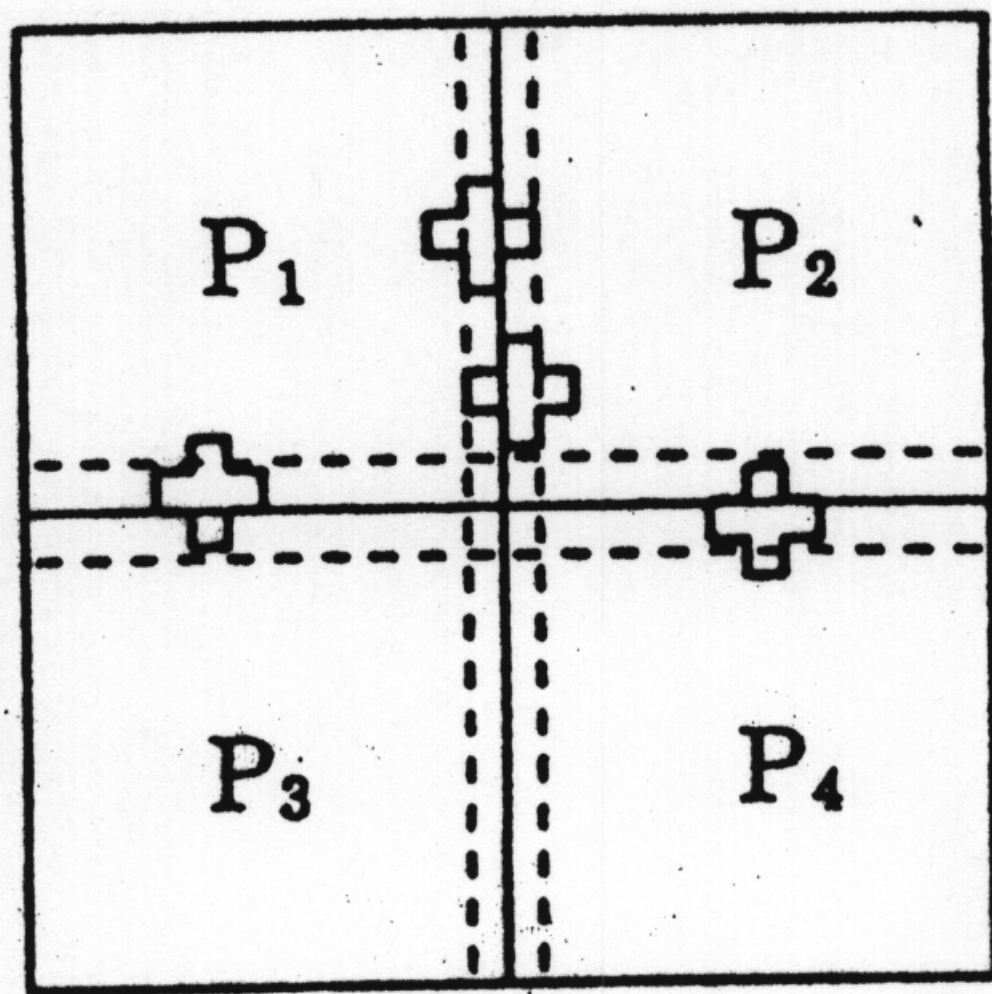
要求側からリード要求を発行する方式
より効率よい

(5) 今後の課題

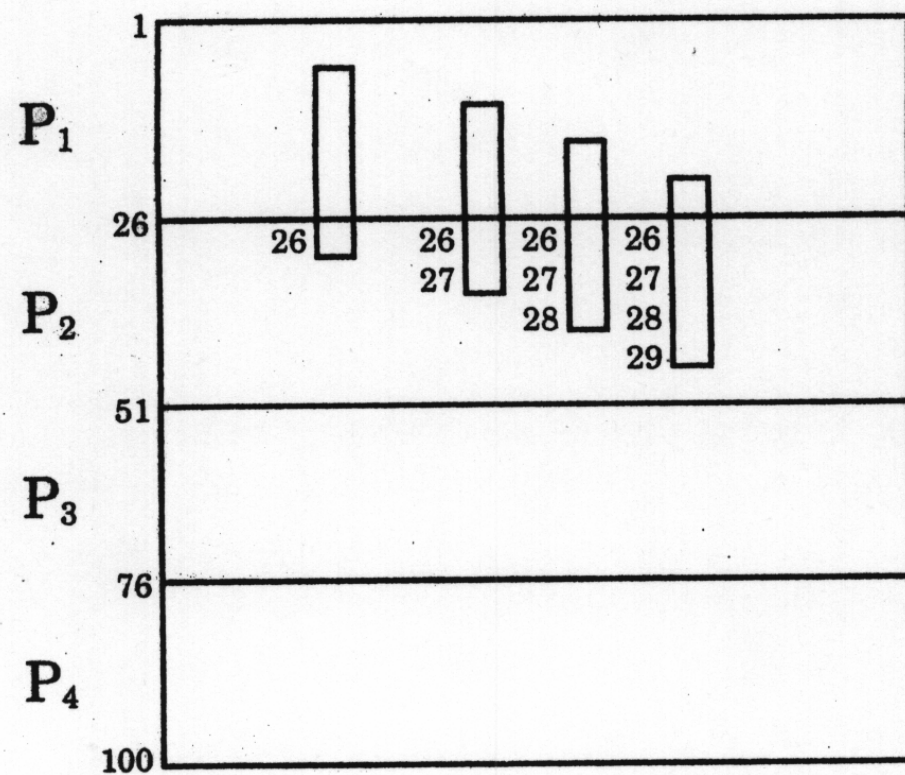
並列化コンパイラ：研究初期段階

人工知能の格好の例題

データベースの構築などによる知識の蓄積と例
題による解法



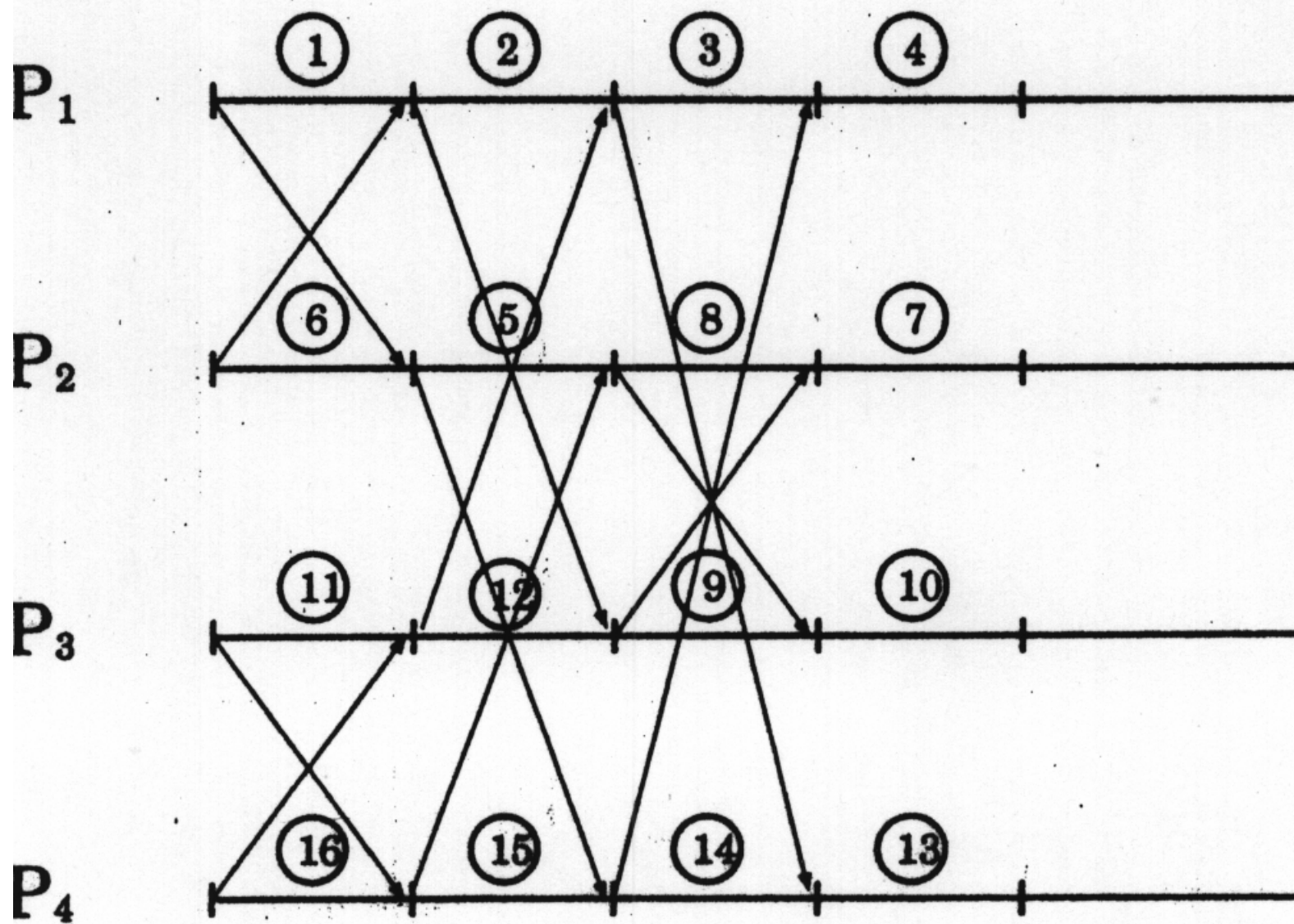
⊕ : ラプラス演算

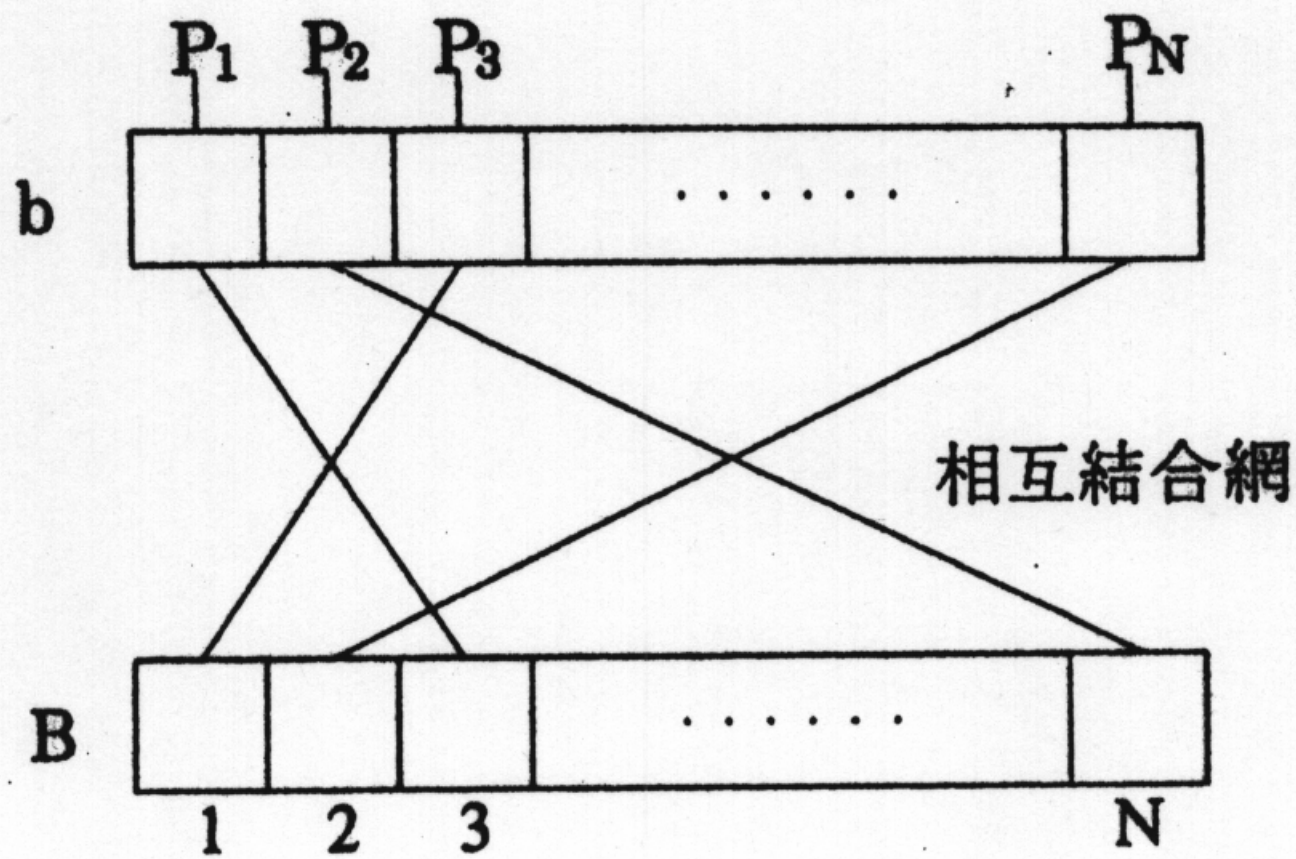


```

DO 10 J = 1, N
DO 20 I = 1, N-4
  A(I, J) = B(I+1, J) + B(I+2, J) +
            B(I+3, J) + B(I+4, J)
20 CONTINUE
10 CONTINUE

```





5 . 8 マルチプロセッサの その他の課題

マルチプロセッサのスケジューリング、
負荷分散、
デバッグなどの研究課題