

Title	SX-4 マルチノードシステムのハードウェア構成
Author(s)	幅田, 伸一
Citation	大阪大学大型計算機センターニュース. 1997, 103, p. 17-24
Version Type	VoR
URL	https://hdl.handle.net/11094/66191
rights	
Note	

Osaka University Knowledge Archive : OUKA

<https://ir.library.osaka-u.ac.jp/>

Osaka University

SX-4 マルチノードシステムのハードウェア構成

幅田 伸一

1. まえがき

SX-4 マルチノードシステムは、SX-4シリーズの最上位に位置するシステムです。SX-4シリーズの最先端アーキテクチャ、国際標準の超高速インタフェース、最先端CMOS LSI技術、高密度実装技術を採用することにより、プロセッサ当たり2ギガFLOPS、システムの最大演算性能1TFLOPSを実現しております。SX-4シリーズは、超高速演算性能を実現するために、

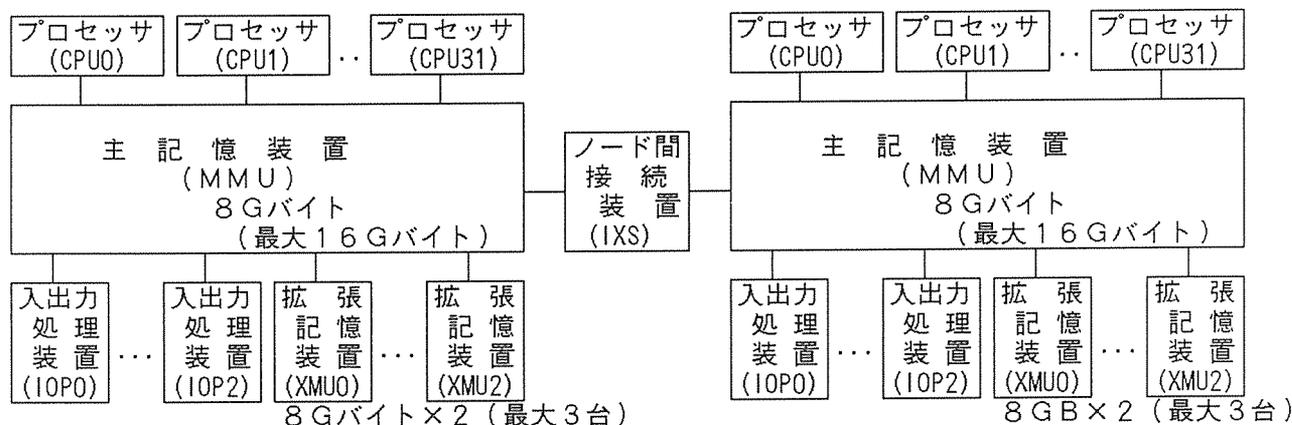
- (1) ベクトル型スーパーコンピュータとして、世界で初めて、共有メモリ型システム（シングルノードシステム）を最大16台接続するクラスタ型システム構成
- (2) 共有メモリ型のベクトル型スーパーコンピュータとして、最大32台のマルチプロセッサ構成
- (3) 種々の高速化技術を採用した高速スカラユニット
- (4) 32本のベクトル演算パイプラインの同時動作が可能な強力なベクトルユニット
- (5) ANSI準拠の超高速HIPPIチャンネル
- (6) 高速入出力処理装置
- (7) 高性能主記憶装置
- (8) 大容量高速拡張記憶装置
- (9) 超高速ノード間接続装置

などを備えております。

以下に、これらのハードウェアについて紹介いたします。

2. ハードウェア構成の要素

SX-4 マルチノードシステムの構成および構成要素のシステム内での位置づけを図1に示します。



SX-4 マルチノードシステムは、SX-4シリーズの最上位システムであり、SX-4 シングルノードシステムを最大16台接続するクラスタ型システムを構成します。SX-4 シングルノードシステムは、最大32台のプロセッサ（CPU：Central Processing Unit）を接続できる共有メモリ型マルチプロセッサシステムであり、プロセッサ、入出力処理装置、主記憶装置、拡張記憶装置から構成されます。

プロセッサは機能的にスカラユニット (SU: Scalar Unit) とベクトルユニット (VU: Vector Unit) とに分けられます。スカラユニットはプロセッサ上で実行される命令の解読と実行制御を行うとともに、スカラ演算パイプラインによりスカラ演算を高速実行します。ベクトルユニットは8セットのベクトル演算パイプラインと144Kバイトのベクトルレジスタを備え、ベクトル演算を高速実行します。

主記憶装置 (MMU: Main Memory Unit) はプログラムおよびそれらが使用するデータ類を格納し、ノードあたり最大16Gバイトの記憶容量があります。

入出力処理装置 (IOP: Input Output Processor) は主記憶装置と周辺処理装置の間に位置し、プロセッサの指示により、プロセッサと独立して入出力処理を高速実行し、ノードあたり最大3台まで接続できます。ANSI準拠の100Mバイト/秒の超高速チャネルであるHIPPI (High Performance Parallel Interface) を最大48本接続でき、超高速スター型のHIPPIネットワーク、UltraNet、超高速アレイディスクを接続できます。この他、16ビットSCSI-2 differential、FDDIなどをノードあたり最大16チャンネル接続できます。

拡張記憶装置 (XMU: Extended Memory Unit) は磁気ディスク装置と主記憶装置の間に存在するアクセスギャップ (アクセス時間の差) を埋めるもので、ノードあたり最大24Gバイトの記憶容量があります。

ノード間接続装置 (IXS: Internode Crossbar Switch) はSX-4シングルノードシステムを最大16台接続することができるパケット交換型スイッチネットワーク装置です。シングルノードシステム側に実装されるノード間接続アダプタ (IXSアダプタ) が提供するデータムーバ機能により、他のノードの主記憶装置上のデータを他のノードのプロセッサを使用せずに、アクセスすることができます。

3. プロセッサ (CPU)

SX-4のプロセッサは、SX-2およびSX-3で培われてきた超高速演算プロセッサ技術をさらに発展させ、最新のアーキテクチャおよびテクノロジーを採用することにより、高速なプロセッサを実現しております。

テクノロジーとして、およそ400万トランジスタを収容可能な超高集積度のCMOS LSIと空冷高密度実装技術を採用し、プロセッサを1枚のカードに実装した1カードプロセッサを実現しました。この最新テクノロジーと高度なパイプライン処理により、2GFLOPSの演算性能を備えた1カードプロセッサを実現しました。

SX-3で培われてきた共有メモリ型マルチプロセッサ方式を発展させ、SX-4シングルノードシステムは最大32台のマルチプロセッサシステムを構築可能とすることにより、最大64GFLOPSの性能を実現しています。

さらに、共有メモリ型マルチプロセッサ方式の限界を乗り越えるべく、ベクトル型スーパーコンピュータとしては初めて、クラスタ型システム方式を採用し、最大16台のシングルノードシステムを接続するマルチノードシステムにより、最大1TFLOPSの性能を実現しました。

プロセッサの構成を図2に示します。

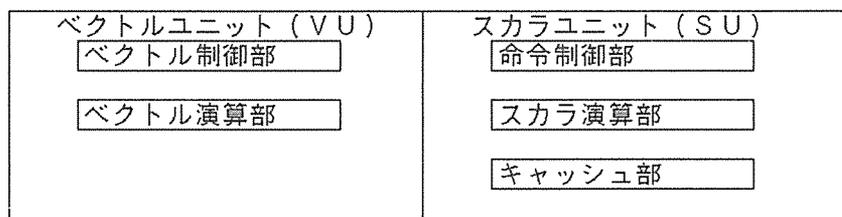


図2 プロセッサ (CPU) の構成

プロセッサは大きく分けて、スカラユニットとベクトルユニットの2つのユニットによって構成されます。

スカラユニットは、プロセッサ上で実行される命令の解釈と実行制御を行うとともに、スカラ演算パイプラインによりスカラ演算を高速実行します。

ベクトルユニットは8セットのベクトル演算パイプラインと144Kバイトのベクトルレジスタを備え、ベクトル演算を高速実行します。

3. 1 プロセッサ (CPU) のアーキテクチャ

SX-4のプロセッサは、SX-3のアーキテクチャを基に、パイプライン処理およびマルチプロセッサ機能の改善、マルチノードシステムにおける他ノードとのデータ転送/通信機能の強化を行っております。

スカラでは、SX-3のRISCアーキテクチャ、スカラ演算のパイプライン処理、ハードウェアによる命令の実行順序の並べ換え制御を改善するとともに、2セットの命令処理パイプラインを備えるスーパースカラアーキテクチャを採用し、高速スカラ処理を実現しています。

ベクトルでは、演算パイプラインの多重並列化により、最大8個の要素に対する浮動小数点加算と最大8個の要素に対する浮動小数点乗算を同時に実行可能とし、さらに、GFLOPSあたりのベクトルレジスタと主記憶装置間のデータ転送性能を改善し、超高速ベクトル処理を実現しています。

マルチプロセッサをサポートする機能としては、プロセッサ間的高速な通信を可能とする通信レジスタの容量を拡大し、マルチタスク方式による高速並列処理を実現しています。シングルノードシステム内の並列処理に使用する通信レジスタは主記憶装置に実装されており、マルチノードシステムにおけるノード間にまたがる並列処理に使用する通信レジスタはノード間接続装置に実装されています。

SX-4 マルチノードシステムでは、SX-3に対して以下のような命令を追加して、ノード間にまたがる並列処理を実現しています。

- ・ノード間データ転送命令
- ・グローバル通信レジスタアクセス命令

3. 2 スカラユニット

スカラユニット (SU) は、命令制御部、スカラ演算部、キャッシュ部から構成されています。

(1) 命令制御部

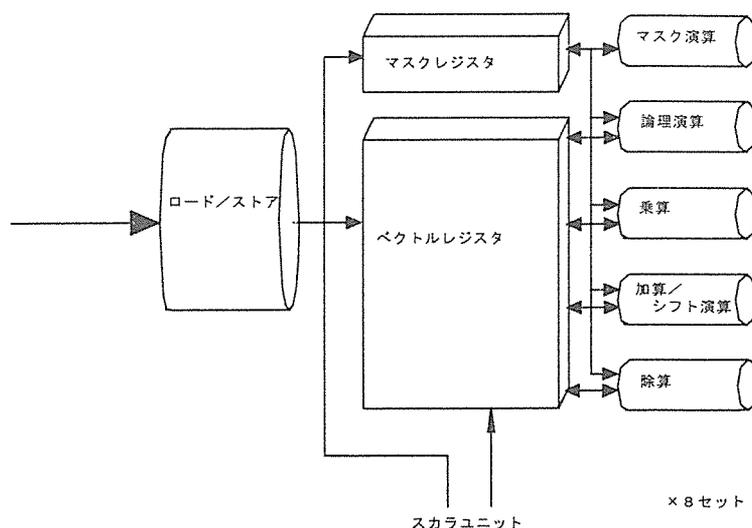
命令制御部は、高速大容量の命令バッファを持ち、プロセッサが実行すべき命令を解釈し、スカラユニットおよびベクトルユニットの各機能を有効に使用するように、命令の実行順序を制御する機能を持っています。

1) 命令制御回路

命令制御回路はパイプライン化され、命令を解釈し、使用するオペランドの準備状態、スカラレジスタへの書き込みバスの競合、ベクトル命令においては、使用するベクトルレジスタ、演算パイプライン、チェイニングのタイミングなどを管理し、スカラ性能およびベクトル性能を最大限に引き出す命令実行順序の制御を行います。

さらに、2セットの命令処理パイプラインからなるスーパースカラアーキテクチャの採用により、1マシンサイクルに2命令の解釈、最大4命令の実行起動を可能としています。

図3にスカラユニットの構成を示します。



この他に、命令制御部は以下のようなプログラムの実行状態を計数するカウンタを備え、プログラムの実行を高速化するために必要な支援情報を提供できるようにしています。

- ・ タイマレジスタ
- ・ 命令実行数カウンタ
- ・ ベクトル命令実行数カウンタ
- ・ ベクトル要素数カウンタ
- ・ ベクトル命令実行時間計測カウンタ
- ・ 浮動小数点演算カウンタ

2) 命令バッファ

命令バッファは128バイト×64ブロック(8Kバイト)の容量を持ち、4バイト命令を2048個分格納することができます。

繰り返しループ中の命令などは、命令バッファに格納されている確率が高く、主記憶装置から命令を読み出すことなく、命令バッファから命令を取り出して実行することができます。

命令バッファの中に実行すべき命令がない場合には、その命令を含む128バイトのブロックを命令キャッシュまたは主記憶装置からロードしますが、その場合、必要になっている命令語からロードを始め、その命令がロードされた時点でただちに解読処理を開始することにより、高速な命令実行処理を可能としています。

3) 高速分岐機構

分岐命令は、分岐する側としない側のいずれかにかたよる傾向があります。命令バッファは、分岐命令が前回どちら側の命令を実行したかを記憶する分岐履歴バッファを備え、この内容に基づいて前回と同じ側の命令を先行して取り出し、解読を進めることにより、分岐命令の高速化を図っています。

4) 大容量スカラレジスタ

スカラユニットは128個のスカラレジスタを備えています。このレジスタは汎用のレジスタであり、アドレス計算用のベース/インデックスレジスタ、固定/浮動小数点演算用の演算レジスタとして利用できます。

128個の大容量レジスタを備えることによって、レジスタ競合による演算パイプラインの乱れを少なくしたり、変数や中間結果のセーブ/リストアする頻度を下げ、メモリアクセス回数を減少させています。さらに、大容量のスカラレジスタにより、コンパイラのレジスタ割り当ての自由度が増し、ハードウェアを効率よく使用するオブジェクトコードの生成が可能となり、高速なスカラ演算を実現しています。

(2) スカラ演算部

スカラ演算部は、8バイトデータ演算を基本とし、2つの固定小数点算術/論理演算器、浮動小数点加算器、乗算器、除算器の5つの演算器から構成されています。これらの演算器は、ベクトル演算器と同様にパイプライン化されており、スカラ演算命令を連続、並列に実行することができます。

(3) キャッシュ部

プロセッサの性能向上には、演算器の性能向上とともに、メモリアクセス性能の向上が重要な要素となります。キャッシュ部には、主記憶装置上の使用頻度が高い命令の写しを保持する64Kバイトの容量をもつ命令キャッシュと、データの写しを保持する64Kバイトの容量をもつオペランドキャッシュがあり、命令制御部からの命令語あるいはオペランドのメモリアクセス要求に対し、高速に応答します。

3.3 ベクトルユニット(VU)

ベクトルユニットは、ベクトル演算部とベクトル制御部から構成されます。図4にベクトルユニットの構成を示します。

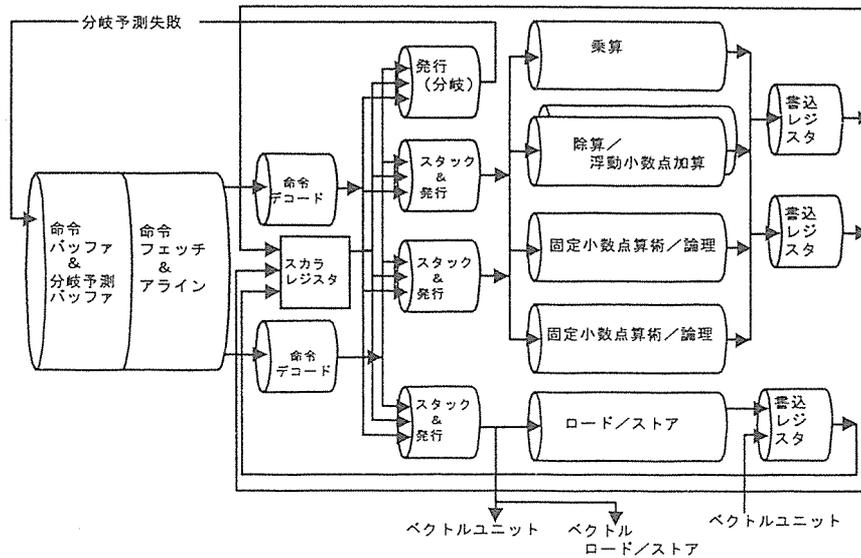


図4 ベクトルユニット (VU) の構成

(1) ベクトル演算部

ベクトル演算部は、加算器/シフタ、乗算器、除算器、論理演算器の4本のベクトル演算パイプラインからなるベクトル演算パイプラインセットを8セット、144Kバイトの容量をもつ72個のベクトルレジスタ、256ビット長のデータを保持する16個のベクトルマスクレジスタ、および、これらを制御する回路から構成されます。

1) ベクトル演算パイプライン

加算器/シフタ、乗算器、除算器、論理演算器からなるベクトル演算パイプラインセットは、合計4本のベクトル演算パイプラインを独立して動作できます。プロセッサは、ベクトル演算パイプラインセットを8セット並列にもつことにより、合計32本のベクトル演算パイプラインの同時動作を可能とし、ベクトル演算の高速化を図っています。

各々の演算パイプラインは、SX-3の標準形式と拡張指数データ形式の2つの浮動小数点データ形式に加え、新たにIEEEの浮動小数点データ形式をサポートしています。

2) ベクトルレジスタ

それぞれが最大256要素まで保持することができるベクトルレジスタを72個用意しています。ベクトルレジスタの総容量は144Kバイトとなります。主記憶装置からベクトルレジスタへのベクトルデータのロード、またはベクトルレジスタから主記憶装置へのベクトルデータのストアは、1セットのベクトル演算パイプラインあたり1マシンサイクルに2語(8バイト/語)の転送能力を持ち、8セットのパイプラインが同時に動作できるように構成しています。さらに、先行するベクトル命令によるベクトルデータの処理が完了するのを待たずに、最初の要素の演算結果がベクトルレジスタに書き込まれた時点で、その結果を入力とする後続ベクトル命令の実行を開始することができる自動チェイニング機構、データ転送前後のベクトルデータの格納読み出しを最適なタイミングで行う制御などを行っています。

3) ベクトルマスクレジスタ

ベクトルマスクレジスタを16個備えており、それぞれ256ビット長のデータを保持します。マスク付き演算、ベクトルの圧縮/伸張/マージなどのベクトル命令の制御に使用するデータを保持し、保持するデータの各ビットがベクトルレジスタの保持する各ベクトル要素に対応します。

4) ベクトル制御部

ベクトル制御部は、スカラユニット内の命令制御部から送られてくるベクトル命令の解釈、8セットの並列演算パイプラインをインタレース方式で動作させるための制御、8セットの並列演算パイプラインに対応するベクトルデータおよびマスクビットのアライン制御、ベクトルデータ長の制御を行います。

スカラレジスタとベクトルマスクレジスタ間のマスクビットの設定および読み出しも並列演算パイプラインを意識することなく実行できるようにしています。

4. 入出力処理装置 (IOP)

入出力処理装置は、主記憶装置と周辺処理装置の間に位置し、プロセッサの指示により、プロセッサと独立して入出力処理を高速実行します。入出力チャンネルとしては、ANSI X3T9.3に準拠した超高速チャンネルHIPPI、16ビットSCSI-2 differential、FDDIなどを提供します。

HIPPIはスーパーコンピュータの超高速周辺装置接続手段および超高速ネットワーク構築手段として広く使用されており、HIPPI用超高速周辺装置としてアレイディスク、ネットワークとしてHIPPIスイッチを使用したスター型ネットワーク、UltraNetなどが接続されます。HIPPIの特徴は以下の通りです。

- ① 100Mバイト/秒の単方向インタフェース (入出力には2チャンネル必要)
- ② 4バイト幅のECLレベルの対信号を使用し、最大25mのケーブルで1:1接続される。
- ③ 1Kバイト単位のブロック転送 (バーストと呼ぶ) を基本に、大量のデータを高速に転送する。
- ④ 垂直パリティおよび水平パリティにより、データの正当性を保証している。

入出力処理装置は、HIPPIを基本チャンネルとし、最大16本のHIPPIチャンネルを接続できます。ノードあたり最大3台の入出力処理装置が接続ができ、HIPPIチャンネルの最大本数は、ノードあたり48本となります。また、チャンネル拡張機構を接続することにより、16ビットSCSI-2 differential、FDDIなどのチャンネルをノードあたり最大16チャンネルまで接続することができます。

5. 主記憶装置 (MMU)

プロセッサ上で実行するプログラムおよびそれらが使用するデータを格納します。

表1に主記憶装置の諸元、図5に主記憶装置の構成を示します。

表1 主記憶装置の諸元

項目	諸元
記憶容量	256Mバイト ~ 16Gバイト
インタレース	32ウェイ ~ 1,024ウェイ
記憶素子	4Mビット同期型スタティックRAM
論理素子	CMOS-LSI
誤り訂正	1ブロック誤り訂正と2ブロック誤り検出
データ転送速度	16Gバイト/秒 ~ 512Gバイト/秒

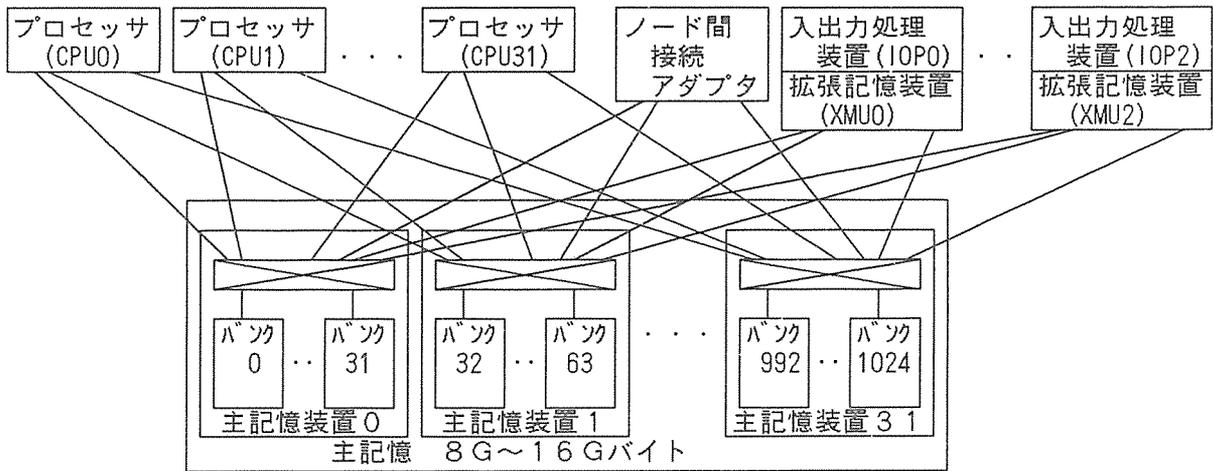


図5 主記憶装置 (MMU) の構成

主記憶装置は、36ポート、32バンク、最大512Mバイトの記憶容量をもち、ノードあたり最大32台の主記憶装置を使用し、最大記憶容量16Gバイト、最大インタレース数1、024ウェイ、最大ポート数1、152ポートの主記憶を構成します。1台のプロセッサは主記憶の32個のポートと接続され、1マシンサイクルに最大16語のベクトルデータをロードまたはストアすることができ、最大16Gバイト/秒の極めて高いデータ転送能力を実現しています。

ノード間接続アダプタ、入出力処理装置/拡張記憶装置も、プロセッサと同じように、主記憶の32個のポートに接続されます。

大容量高性能の主記憶装置を実現するため、記憶素子としてアクセス時間15ナノ秒、4MビットのシンクロナスSRAMを使用しています。

データに対しては、1ブロックの誤り訂正と2ブロックの誤り検出符号を採用し、信頼性の向上を図っています。

6. 拡張記憶装置 (XMU)

拡張記憶装置は、主記憶装置と磁気ディスク装置の間に位置する大容量半導体記憶装置です。拡張記憶装置は、最大8Gバイトの記憶容量を有し、4Gバイト/秒のデータ転送性能をもちています。ノードあたり最大3台の拡張記憶装置を接続することができ、最大48Gバイト、12Gバイト/秒のデータ転送速度を備える拡張記憶を構成することができます。

表2に拡張記憶装置の諸元を示します。

項目	諸元
記憶容量	2Gバイト ~ 8Gバイト
増設単位	2Gバイト
記憶素子	16Mビット・ダイナミックRAM
論理素子	CMOS-LSI
誤り訂正	1ブロック誤り訂正と2ブロック誤り検出
最大データ転送速度	4Gバイト/秒

拡張記憶装置の特徴を以下に示します。

- ① 記憶素子に16MビットのダイナミックRAMを使用し、高密度実装とともに、装置の大容量/小型化を実現しています。
- ② 論理素子にプロセッサと同じCMOS-LSIを採用し、処理速度の高速化を実現しています。

- ③ 電源供給モジュールを冗長構成とし、電源供給モジュールの障害が発生した場合でも、装置の電源断とならない方式を採用し、信頼性の向上を実現しています。
- ④ 内部処理データ幅を4バイト単位で複数のブロックに分割し、各々のブロックに対してECC (Error Cecking and Correction) による1ブロックの誤り訂正と2ブロックの誤り検出を行うことにより、ブロック間にまたがる複数ブロックエラーの自動訂正を可能とし、信頼性の向上を実現しています。
- ⑤ プロセッサに同期したクロックで動作することにより、データ転送速度の高速化を実現しています。

7. ノード間接続装置 (IXS)

ノード間接続装置は、SX-4シングルノードシステムを最大16台接続できるパケット交換型スイッチネットワークです。シングルノード側に実装されるノード間接続アダプタが提供するデータムーバ機能とノード間接続装置のネットワーク機能により、他ノードのプロセッサを使用せずに、他ノードの主記憶装置上のデータをアクセスすることができるノード間主記憶データ転送機能を提供します。ノード間主記憶データ転送機能には、以下のものがあります。

- ・ブロック転送機能
- ・2ディスタンス転送
- ・間接アドレス転送 (リストベクトル転送)

ノード間にまたがるマルチタスク方式の並列処理を高速化するため、上記ノード間主記憶データ転送機能の他に、ノード間接続装置はグローバル通信レジスタ (GCR: Global Communication Register) を備えております。グローバル通信レジスタに対する操作には、以下のようなものがあります。

- ・ロード/ストア
- ・テスト&セット
- ・フェッチ&インクリメント
- ・フェッチ&ディクリメント

ノード間接続装置とノード間接続アダプタ間は、超高速光インタフェースを使用し、8Gバイト/秒の極めて高いデータ転送性能を実現しています。光を使用することにより、ノード間接続装置とシングルノード間を接続する光ケーブル長は最大100mまで延長することができます。表3にノード間接続装置の諸元を示します。

表3 ノード間接続装置の諸元

項 目	諸 元
最大データ転送速度	8Gバイト/秒
最大ケーブル長	100m
GCRの容量	4K語 (32Kバイト)
論 理 素 子	CMOS-LSI
光 素 子	1.25Gbps 光モジュール

[執筆者紹介] NEC コンピュータ事業部第四技術部