

光インタコネクタ搭載ネットワークスイッチ

西村信治、工藤知宏*、西宏章*、原澤克嘉**、松平信洋**、坪重人**
RWC光インターコネクション日立研究室

東京都国分寺市東恋ヶ窪1-280 (株)日立製作所中央研究所内

Tel:0423-23-1111 (ext. 3560), Fax:0423-27-7740

E-mail:nisimura@crl.hitachi.co.jp

*RWCつくば研究センタ, RWCP並列分散システムアーキテクチャ研究室

**日立通信システム

要旨

分散環境に配置したパーソナルコンピュータを高速光ネットワークにて接続することにより高性能な並列計算機システムが構築できる。これを実証する為、大容量光インタコネクションと高速スイッチLSI (LVDS-I/O; 784ピン-BGA; ワンチップCMOS-LSI) を組み合わせた大容量8×8スイッチシステムを (RHINET-2/SW) 開発した。RHINET-2/SWは分散環境での高速並列計算処理に必要な、大容量 (64 Gbit/s) パケットデータスイッチ機能を小型装置で実現する。光インタコネクタの使用により、ノード間接続距離は最大100mを実現している。

和文キーワード:光インタコネクション、並列計算機、ネットワークスイッチ

Network switch implemented with optical interconnection

Shinji Nishimura, Tomohiro Kudoh*, Hiroaki Nishi*,
Katsuyoshi Harasawa**, Nobuhiro Matsudaira**, and Shigeto Akutsu**

RWC Optical Interconnection Hitachi Laboratory

c/o Central Research Lab., Hitachi, Ltd.,

1-280 Higashi-koigakubo, Kokubunji-shi, Tokyo

Tel:+81-423-23-1111 (ext. 3560), Fax:+81-423-27-7740

E-mail:nisimura@crl.hitachi.co.jp

* RWCP Distributed and Parallel System Architecture Tsukuba Laboratory, RWC Tsukuba research center

** Hitachi Communication Systems Inc.

Abstract

A network based parallel computing system can be constructed by optically interconnected PC network. We have developed a high-throughput, compact network switch (the RHINET-2/SW) for a distributed parallel computing system. Optical interconnection modules and one-chip CMOS SW-LSI (LVDS-I/O, 784-pin BGA package) are integrated on a compact RHINET-2/SW circuit board. The RHINET-2/SW switch enables large-throughput (64 Gbit/s), packet data switching for a high-speed parallel processing in a distributed computing environment. Optical interconnection provides long-transmission-length (< 100 m) interconnection between the nodes.

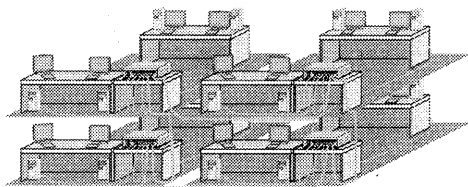
Keywords:Optical interconnection, Parallel computer, Network switch

1 はじめに

パソコン (PC) やワークステーション (WS) の単体性能は近年著しく向上しており、ハイエンドの超並列計算機のノード性能に匹敵もしくは上回るようになってきている。このため高速光ネットワークで PC・WS 間を接続し並列計算システムを構成した場合、非常に高い計算機能力を得られる可能性がある。

我々は本可能性を実証するため分散並列計算機システム RHiNET-2 (RWCP high-performance network-2) の開発を進めている。RHiNET-2 システムは、長距離を広帯域接続できる光インターコネクションを PC ノード間接続に採用し、ビル内フロア内に分散配置された計算機類を自由に接続して高性能な並列計算処理を実現する事を目標としている (図 1 参照) [1-3]。高性能 PC がオフィスや研究機関等に数多く配置されるようになった現状を考えると、一箇所に集められた計算機だけでなく、ビル・フロア内などに広く分散配置された計算機群を効率良く使用することで、「オフィスを丸ごとスパコンとして使う」ような形態の並列コンピュータシステムが実現可能になると考えている。

RHiNET-2/SW: 8×8 crossbar, 64-Gbit/s data throughput



8.8-Gbit/s (800-Mbit/s × 11-bit, 1-bit clock)
parallel optical interconnection
図 1 RHiNET システムの概念図

RHiNET-2 システムにおいては、各 PC は PCI バスベースの NIC (network interface card) を搭載し、8.8Gbps 光インターコネクションを介して 8×8 の専用クロスバースイッチ (RHiNET-2/SW: 図 2 参照) に接続される。本システムに用いた光インターコネクションは 800Mbps×12 ビットの並列同期伝送 (内訳: データ 11 ビット、クロック 1 ビット) を実現し、接続距離は最大 100m と電気ケーブルの 10 倍を有する。また搭載されたスイッチ LSI には、大容量の SRAM (512kbyte) をオンチップ搭載し、メモリアクセス時間を最小化を図っている。本稿では、RHiNET-2

システムの中心装置である 8×8 スイッチボード (RHiNET-2/SW: 通信容量 64Gbps) に関して、構造と動作試験結果に関して述べる。

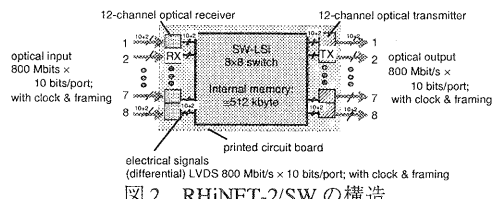


図 2 RHiNET-2/SW の構造

2 RHiNET-2 における光インターコネクション

RHiNET-2 は並列コンピュータシステムである。この為、通信容量の拡大ばかりでなくノード間通信のレイテンシーを最小限にする必要がある。また、通常の PC をノードとして使用しているため、光素子の消費電力や物理的サイズも小さいものが要求される。

RHiNET-2 においては、リボンファイバを用いた並列同期転送を採用した [4-7]。用いた光モジュールは日立製作所製のインタコネクトモジュール [MDS4212A/MDR4212A] である [4-6]。本モジュールは 12 チャンネルのレーザとホットダイオードを並列駆動する事により、小型パッケージ (3.9cc) で 8.8Gbps の大容量光データ接続を可能にしている (800Mbit/s × データ 11 チャンネル + クロック 1 チャンネル)。発振波長 1310nm の端面発光レーザアレイを使用し、12 芯シングルモードリボンファイバ (MPO [multi-path push-on] コネクタ付) と組合せる事により、伝送中の波形劣化を抑えて、並列同期信号のチャンネル間スキューを最小化している。電気インタフェースは LVDS と P-ECL インタフェースを搭載し、コンピュータ内部回路と容易に接続できる。ファイバ接続距離は最大 100m を実現しており (制限要因はファイバスキュー)、電気ケーブル (限界距離: 約 10m) よりも自由度の高いノード間接続が可能である。また、搬送クロックを並送する並列同期伝送の採用により、受信側でフレーム同期やクロック生成等の複雑な信号処理を必要とせず、同期系で構成されるコンピュータシステム同士を接続するのに適したコンパクトで低遅延なデータ転送を実現している。

1Gbps クラスの伝送速度を有する標準技術 Gigabit Ethernet [8, 9]・Fiber Channel においては、送受各々において 1 本のファイバと 1 つの波長信号を用いたシ

リアルデータ転送を採用している。しかし、シリアルデータ転送にて、より高速なデータ転送レートを実現する為には厳密な温度管理が必要となり、結果として送受信モジュールの大型化と消費電力の増大をもたらす。この為、現在のデバイス技術では、厳密な温度管理なしに実現できるデータレートは、最大 2.5 Gbps 前後と考えられている。それ故、2.5 Gbps を超える伝送レートを、装置間接続等の近距離において小型装置サイズで実現するには、複数のデータ転送系を用いた並列データ転送がもっとも有利であると考えられる。

3 RHINET-2/SW システムの構造

(1) スイッチ LSI

RHINET-2/SW に使用したスイッチ LSI は 8 入力 8 出力のクロスバススイッチ機能を搭載している。各ポートは 10 ビットの packets データと 1 ビットのフレーム同期信号、1 ビットのクロック信号の 12 ビット構成を有し、8Gbps/port の通信容量を実現できる。前節の高速動作と CMOS 化に対する要求を満たすため、全ての高速電気インタフェースを高速 CMOS-LVDS (low-voltage differential signaling) で統一した。I/O ピンは 800Mbit/s の差動信号線を 192 組 (384 本) 入出

力する構成とし、単一チップで 64Gbps の大容量スループットを実現している (これは CMOS のネットワークスイッチ LSI としては世界トップクラスの性能である)。0.18 μm プロセスの CMOS-ASIC の使用により、大容量 SRAM (512kbyte) のオンチップ搭載を実現している。LSI パッケージは 784pin BGA [ball grid array] である。

図 3 に LSI の内部ブロック図、図 4 はスイッチ LSI のダイのレイアウトを示す。800Mbps \times 10bits の入力信号は、入力端で 1:8 DEMUX にて 100Mbps \times 80bits に展開された後、ECC デコーディングとエラスティックバッファを通過する。搬送波クロックに同期した入力信号はエラスティックバッファ内で内部クロック (BASECLK: 100MHz) に再同期される。エラスティックバッファからの入力信号 (80bit) はスイッチコア部 (SW-core) にて行路切換処理される。スイッチコアからの出力信号は ECC エンコーディング処理を経て、8:1 MUX にて 800Mbps \times 10bits の信号に多重化され、フレーム同期信号 (AO) と搬送波クロック (OUTCLK) を付加されて LSI から出力される。LVDS 入出力バッファ部は 800Mbps を実現するため、高速特性を要求される。

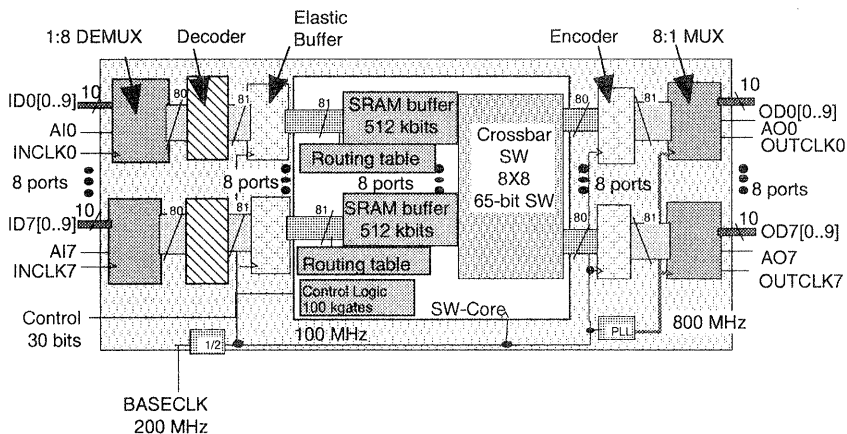


図 3 LSI の内部ブロック図

Routing Table (8 ports): 2 bit * 4096- for single cast+ 4 bit * 64- for multicast
 ID0[0..9], ID1[0..9], ..., ID7[0..9]: input data (LVDS); 10 bits*800 Mbit/s

OD0[0..9], OD1[0..9], ..., OD7[0..9]: output data (LVDS); 10 bits*800 Mbit/s
 INCLK 0..7: input clock (LVDS); 400 MHz,
 OUTCLK 0..7: output clock (LVDS); 800 MHz
 AI0..7, AO0..7: framing (LVDS) 800 Mbit/s

RHiNET-2を含めた高速計算機システムのノード内の計算処理は電気信号にて処理される。このため、各ノードのデータ入出力部においては光・電気信号変換処理が必要となる（光信号のままのスイッチング技術 [11,12] も様々開発されているが、通信路の設定に必要な時間が最低でもマイクロ秒オーダーの時間を要し、LAN 装置には使いにくい）。この為、高速な光 I/O を有するネットワークスイッチには、高速電気・光インターフェースが必要となる。また同時に、メモリアクセスのボトルネックを回避するためにロジック回路と大容量メモリのオンチップ集積が必須となる事から、構成回路の CMOS 化が求められる。

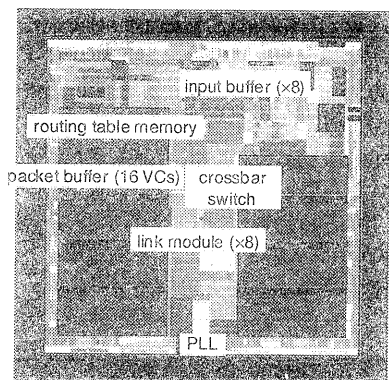


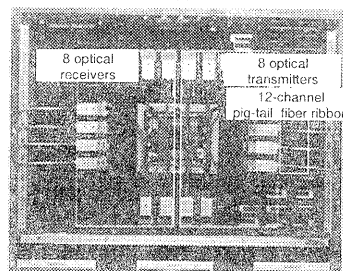
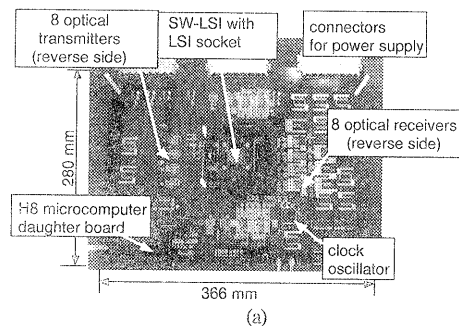
図4 LSIダイ内部のブロック配置図

(2) RHiNET-2/SW のボードレイアウト

RHiNET-2/SW においては、ボードの中心に CMOS スイッチ LSI が搭載され、その直近に光インタコネクションモジュール送受 8 対を高密度実装して 8×8 のクロスバースイッチを実現している (図 5, 6 参照) [12]。各ポートは 12bit 並列データを速度 800 Mbps で並列同期入/出力できる (ビット内訳: パケットデータ信号 10 bit、搬送クロック 1 bit、フレーム信号 1 bit、総データ容量: 8 Gbps/port)。ボードサイズは 280×366mm である。

LSI は高速デバイス専用のソケットに搭載されている。この LSI ソケットは、帯域 (DC-6GHz)、 $\text{path-inductance} < 1 \text{ nH}$ 、 $\text{capacitance}[\text{signal-to-signal}] < 1 \text{ pF}$ と優れた高周波特性を有する。スイッチ LSI の制御用に RS-232C インタフェースを持つ H8 マイコンのデータボードを搭載しており、ルーティングテーブル等の書換は RS-232C、H8 ボードを介して外付けの PC から行う。周波数 200MHz の内部クロック源と

して水晶発振器を 1 つ搭載している。LSI のピン配置およびボード内配線のレイアウトは、前述のテスト基板による実測結果を踏まえ、プリント基板配線の最短化 (最長 150mm) と交差数 0 の制約条件の元に最適化設計している。



(b)

図5 RHiNET-2/SW, マザーボードの部品レイアウト (a): LSI 搭載側, (b): 光モジュール搭載サイド

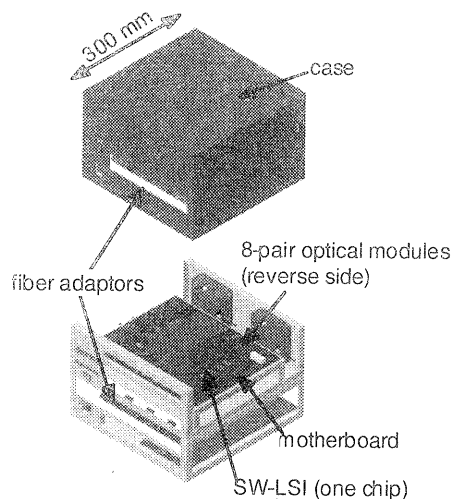


図6 RHiNET-2/SW の概観図

4 RHiNET-2/SW システムの評価結果

RHiNET-2/SWの各入力ポートに800Mbps×10bitのPRWS信号を入力し、出力信号のアイパターンの観測し、ビットエラーレート(BER)を測定した。図7に評価系の構成を示す。800Mbps×12チャンネル(クロックCLK, フレーム信号AI, データDI[9..0])の同期信号をデータジェネレータ(HP80000)にて生成し、光インターコネクションモジュール[MDS4212A]にて光信号に変換したのち、50mのリボンファイバを介してRHiNET-2/SWの入力ポートに入力した。RHiNET-2/SWの出力ポートからの光信号は、同じく50mのリボンファイバを通過の後、光イ

ンターコネクション受信モジュール[MDR4212A]にて電気信号に変換され、ロジックアナライザ(HP81200-LA)とオシロスコープ(HP54120)にて観測した。

ロジックアナライザの入力段の信号のアイパターンを見ると(図8参照)Tr/Tf:150ps、ジッタ<100psの良好な波形を得られた。またBERはデータレート800Mbpsの条件で 10^{-11} 以下の高信頼性を実現した。本評価結果はRHiNET-2/SWのデータI/O系の信頼性は、並列計算機システムを構成する上で要求されるものを十分満たしていると考ええる。

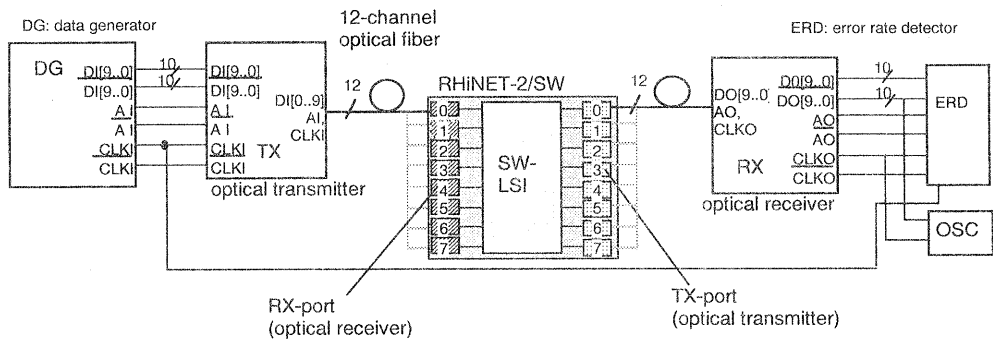


図7 測定系

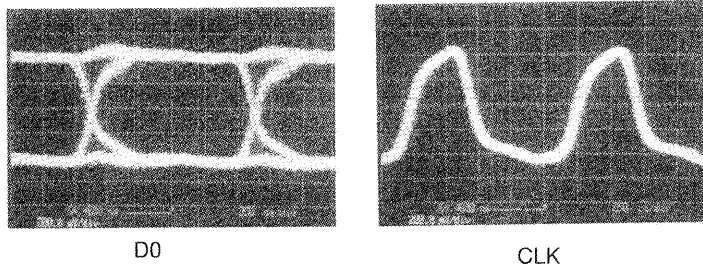


図8 アイパターンの測定結果
0データビット[D0], クロック[CLK]
250 ps/div, 200 mV/div

5 まとめ

RHiNET-2 システムはオフィス LAN 等の分散配置された計算機環境において、高性能な並列計算処理を実現できる。8.8Gbps 大容量光インタコネクション 8 対と 64Gbps 高速スイッチ LSI (LVDS-I/O; 784 ピン-BGA; ワンチップ CMOS-LSI) を組み合わせた大容量 8×8 スイッチボード (RHiNET-2/SW) 開発した。本装置は分散環境での高速並列計算処理に必要な、大容量 (64 Gbps) パケットデータスイッチングを、小型装置で実現可能とする。並列同期光インターコネクションを採用することにより、小型装置規模・低遅延でスキューの小さいデータ転送を実現している。高速動作素子をボード上に高密度実装するため、MWB™を用いた高速・高密度配線技術を開発し、RHiNET-2/SW ボードに応用した。全ての電気インタフェースは CMOS-LVDS に統一し、高速動作・低コスト化を容易にしている。本プロトタイプシステム RHiNET-2/SW において BER 測定による信頼性評価試験を実施した結果、データレート 800Mbps 動作時において BER 10^{-11} 以下の高信頼性を確認した。RHiNET-2/SW ネットワークシステムを用いることにより分散配置環境での高性能並列計算機能が実現可能である。

参考文献

- [1] T. Kudoh, J. Yamamoto, F. Sudoh, H. Amamo, Y. Ishikawa, and M. Sato: "Memory based light weight communication architecture for local area distributed computing", Innovative architecture for future generation high-performance processors and systems, IEEE Computer Society Press, pp. 133-139, (1997).
- [2] S. Nishimura, H. Inoue, H. Matsuoka, and T. Yokota: "Optical interconnection subsystem used in the RWC-1 massively parallel computer", IEEE Journal of selected topics on quantum electronics, vol. 5, pp. 360-367, (1999).
- [3] T. Yoshikawa and H. Matsuoka: "Calibration-free parallel optical-interconnection subsystem implemented by a GByte/s-array optical transceiver and a one-ship link LSIs (Invited)", Technical Digest of OC'98, pp. 524-527, Brugge, Belgium, (1998/6).
- [4] A. Takai, T. Kato, S. Yamashita, S. Hanatani, Y. Motegi, K. Ito, H. Abe, and H. Kodera: "200-Mb/s/ch 100 m Optical Subsystem Interconnections Using 8-Channel 1.3- μ m Laser Diode Arrays and Single-Mode Fiber Arrays", J. of Lightwave Technology, vol. 12, pp. 260-270, (1994).
- [5] 刀祢平高一朗, 三浦篤, 高井厚志, 上野聡, 内田勝己, 豊中隆司, 齊藤勝美: "800 Mbit/s/ch×12ch 光インタコネクト受信モジュール", 1999 電子情報通信学会ソサエティ大会, SC-4-2, (1999).
- [6] <http://www.hitachi.co.jp/Prod/tcd/hikari/tjn00634.htm>
- [7] 三好一徳: "並列光インタコネクション用 622Mbit/s×12ch 送受信モジュール", 電子情報通信学会第 1 回光インターコネクト情報処理研究会, OIP99-8, (1999).
- [8] HIPPI-6400 working drafts, T11.1 maintenance drafts of ANSI NCITS
- [9] IEEE802.3 Higher Speed Study Group <http://grouper.ieee.org/groups/802/3/>
- [10] Keishi Habara, Tohru Matsunaga, and Ken-ichi Yukimatsu: "Large-Scale WDM Star-Based Photonic ATM Switches", Journal of Lightwave Technology, Vol.16, No.12, pp.2191-2201, (1998/12).
- [11] S. Kitajima, H. Takano, and M. Kobayashi: "Traveling type optical cell buffer with small variation of output cell level", IEICE Trans. of Communications, vol. E82-B, pp. 281-287, (1999).
- [12] S. Nishimura, K. Harasawa, N. Matsudaira, S. Akutsu, S. Sasaki, T. Kudoh: "Network switch using optical interconnection for high performance distributed parallel computer using PCs", Proceedings of the Sixth International Conference on Parallel Interconnects, pp. 5-12, Anchorage, (1999/10).