

Maestro2 ネットワークインターフェースの開発

小野雅晃¹

筑波大学システム情報工学等支援室（装置開発班）
〒305-8573 茨城県つくば市天王台 1-1-1

概要

クラスタ型コンピュータ Maestro2 用のネットワークインターフェース (NI) を開発したので報告する。NI はパーソナルコンピュータ(PC)の 64 ビット幅の PCI スロットに挿して使用するイーサカードのような通信用のカードである。NI の PCI インターフェースの最大スループットは 64 ビットデータバス幅、66MHz 動作であるため 533Mbyte/sec で、PCI としては最大のスループットである。インターフェースは LVDS を使用し、物理的な最大スループットは片方向 3.2Gbps(Giga Bit Per Second)である。

1. はじめに

システム情報工学研究科、コンピュータサイエンス専攻の和田研究室では、クラスタ型のコンピュータに使用するために、高速、高機能のアドイン・カード及びスイッチから構成される通信システムを開発している。

5 年ほど前には、200Mbps の IEEE1394(i-LINK)を使用した第一世代の通信システム Maestro1 を開発した。

今回は、Maestro2 として第 2 世代の通信システムを開発した。使用した通信インターフェースは 3.2Gbps の LVDS である。通信システムは前述したように通信カードとスイッチで構成されるが、今回はそのうちのネットワーク・インターフェース(NI)について発表する。

2. NI の構成

図 1 に NI の写真、図 2 に NI のブロック図を示す。

NI は、300MHz 動作の PowerPC プロセッサを搭載したインテリジェントな PCI カードである。NI には CPLD(Complex Programmable Logic Device) FPGA(Field Programmable Gate Array)、PowerPC プロセッサ、FLASH ROM(Flash Read Only

Memory)、SDRAM(Synchronous Dynamic Random Access Memory)、LVDS(Low Voltage Differential Signalling)インターフェース用 IC が搭載されている。CPLD は主に FPGA のコンフィギュレーションの制御を行う。コンフィギュレーションの方式はスレーブセレクトマップで、コンフィギュレーション・データを FLASH ROM から 8 ビットずつ FPGA に転送できるため、迅速に FPGA をコンフィギュレーションすることが出来る。また、コンフィギュレーション・データが間違っていて、FPGA が立ち上がらない場合には、ウォッチドック機能により、正常なコンフィギュレーション・データを FPGA にロードして、エラーの回復を図ることも出来るように設計されている。(現在この機能はまだ実装されていない) CPLD は Xilinx 社の XC95144XL を使用し

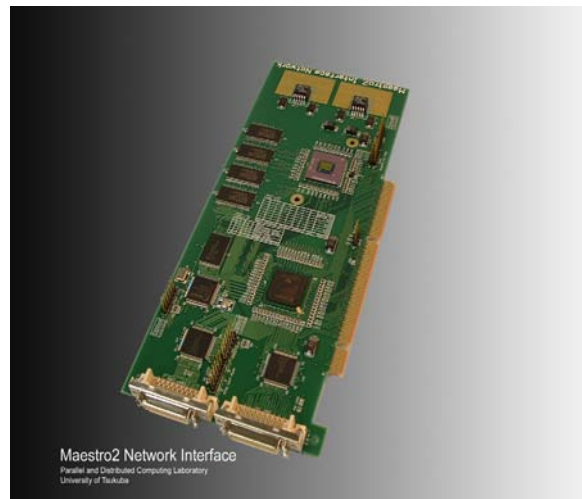


図 1 NI の写真

ている。

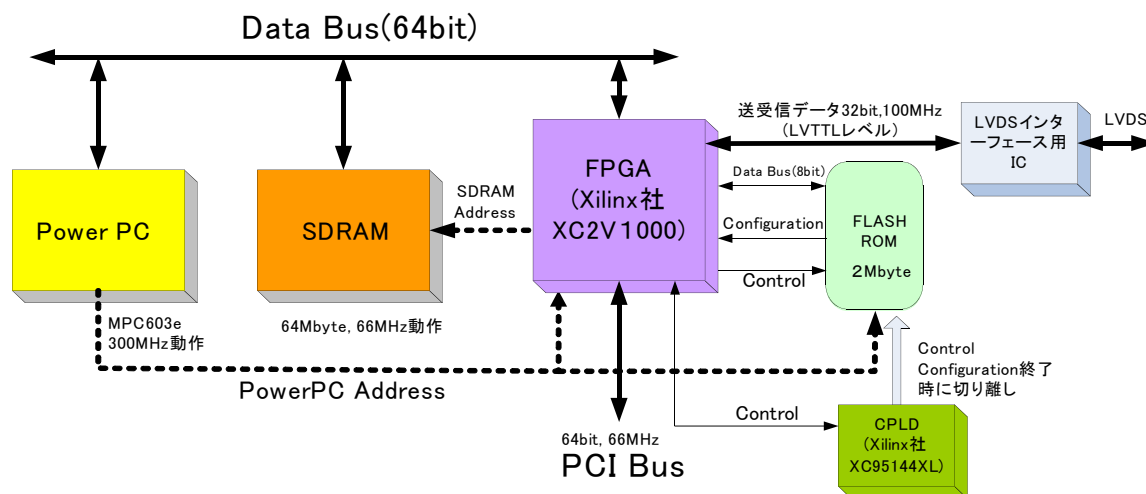


図 2 NI のブロック図

¹ E-mail: ono@sie.tsukuba.ac.jp; Tel: 029-853-5195

FLASH ROMは主に、FPGAのコンフィギュレーション・データを格納する目的で使用される汎用のFLASH ROMである。電源ON時には、CPLDの制御下に入り、FPGAにコンフィギュレーション・データを送る。それが終了すると、CPLDの制御から切り離され、FPGAの制御下に入る。

FPGAは、各デバイスの制御回路、DMAコントローラ、制御用レジスタが実装されているNIの心臓部である。NI上のPowerPCプロセッサやホストPCからのリソースの要求に応じて、データ転送の仲立ちをする。FPGAにはXilinx社のVirtex2シリーズのうちから100万ゲート相当のXC2V1000を使用した。

SDRAMはPowerPCプロセッサとFPGAにバスが直結されている。データバス幅は64ビット、動作周波数66MHzで動作している。

LVDSインターフェース用ICは、ナショナル・セミコンダクタ社のDS90CR481とDS90CR484を使用している。これらのICは、LVTTTL48ビットの入出力をLVDSデータ8ビット+クロック1ビットに変換する。FPGAからのLVTTTL信号は32ビットしかないので、48ビットの内32ビットを使用し、後の16ビットは使用していない。これは主にFPGAのピン数が足りないためである。FPGAからLVDS用ICに出力する周波数は100MHz、バンド幅は32ビットの100MHzなので、3.2Gbpsとなる。LVDSインターフェースで他のPCに挿入されたNIや独自仕様のスイッチと接続される。LVDSインターフェースで接続されるのは、NIや独自仕様のスイッチに限られ、汎用のボードやスイッチなどとは接続できない。

NIはPCIバスでホストPCと接続される。このPCIバスは64ビット幅、動作周波数は64MHzである。NIへのスレーブ・アクセスは連続的にデータ転送(バースト転送)をすることが出来ないで、全て単一転送となる。マスタ・アクセスはバースト転送をすることが出来る。その際の最大転送レートは、533Mbyte/sec(約4.26Gbps)となる。

尚、NIの製作で私が担当した箇所は、CPLD設計とFPGA設計の2/3ほどである。NIの回路設計は山際氏²と共同で担当した。NIのファームウェア、ホストPCのソフトウェアは、山際氏と青木氏³が担当した。

3. CPLDの設計

CPLDは電源ON時に、FPGAのコンフィギュレーションの制御を行う。CPLDはコンフィギュレーション・データをEEPROMに格納しているため、電源ONですぐに動作できる。しかし、FPGAは内部の回路情報をRAMに格納しているため、電源をOFFすると内部の回路情報はクリアされてしまう。そのため、電源ON時には、毎回ROMから内部回路の情報(コンフィギュレーション・データ)を読み込まなければならない。このROMには、FPGAコンフィギュレーション専用のROMを選択するのが普通であるが、本NIでは、あえて汎用FLASH ROMを使用することとした。

それは次のようなメリットがあるからである。

1. PCIバス経由でホストPCから、コンフィギュレーション・データを書き込むことが出来る。

2. ホストPCから書き込んだコンフィギュレーション・データが不正でも、回復させることが可能である。

電源ONのFPGAコンフィギュレーション時には、FPGAはモード選択ピンの設定によって、スレーブセレクトマップモードが選択される。スレーブセレクトマップモードはFPGAの外からクロックと8ビット幅のコンフィギュレーション用のデータをFPGAに入れる方式である。データが8ビット幅と広いので、高速にコンフィギュレーションすることが出来る。

CPLDはFLASH ROMにアドレスを供給して、FPGAのコンフィギュレーション用のデータを出力させる。FPGAとのハンドシェイクを行って、FPGAにデータが入力できたら、FLASH ROMのアドレスを+1して、FLASH ROMに次のデータを出力させる。このように、CPLDはFPGAのコンフィギュレーションを行う。

FPGAコンフィギュレーションが終了すると、CPLDはFLASH ROMの制御信号線を全てハイ・インピーダンス状態にして、FLASH ROMを切り離す。FLASH ROMの制御信号線はFPGAにも接続されているので、それから後はFPGAがFLASH ROMを使用することができる。FPGAはPCIコンフィギュレーション・アクセスに対して、FLASH ROMの領域をメモリ領域にアドレス・マップするので、ホストPCから読み書きすることが出来るようになる。さらに、NIのローカル・プロセッサのPowerPCからも書き込むことが出来るようになる。FPGAコンフィギュレーション専用のROMでコンフィギュレーションする場合には、NIボード上のJTAG端子にコンフィギュレーション用ツールを接続しなければならない。しかし、本NIで採用した方式を使用すれば、ホストPCから直接コンフィギュレーション・データを書き換えることが出来る。また、その他のデータもFLASH ROMに入れておくことが出来る。現在、イーサネットのMACアドレスなどを入れて、利用している。

FPGAに書き込むコンフィギュレーション・データが間違っていると、NIが立ち上がらなくなる。こうなると、遠隔操作でリカバリすることが出来なくなるため、ホストPCが遠い所にある場合はその場所まで行かなければならない。そのような場合のために、FLASH ROMに2つのコンフィギュレーション・データを格納しておく。新しく書き込んだデータが間違っていて、NIが立ち上がらない場合に、古く確実に立ち上がったデータを書き直すことで、NIを確実に立ち上げるようにすることも出来る。現在は、この機能はCPLDに実装していない。

4. FPGAの設計

FPGA内部モジュールの構成を図3に示す。矢印の線は、モジュール間のデータや制御信号のつながりを示す。

FPGAは、MLX32、PCI、PCI Target Controller、PowerPC Interface、Internal Registers、Network Buffer Link Interface(NBL Interface)、SDRAM Controller、DMA Controller、System Arbiter、FLASH ROM Controllerの各モジュールから構成される。

MLX32は、Maestro2独自の通信プロトコルを実装する。OSI階層のトランスポート層でデータの完全性を保障し、バースト・データ転送を続けて行う機能を持っている。このモジュールは山際氏が設計した。

PCIモジュールは66MHz、64ビット幅のPCIバスのインターフェース回路である。PCIスレーブ・アクセス(ホ

² Instituto de Engenharia de Sistemas e Computadores (INESC-ID), Portugal

³ 筑波大学 システム情報工学研究科

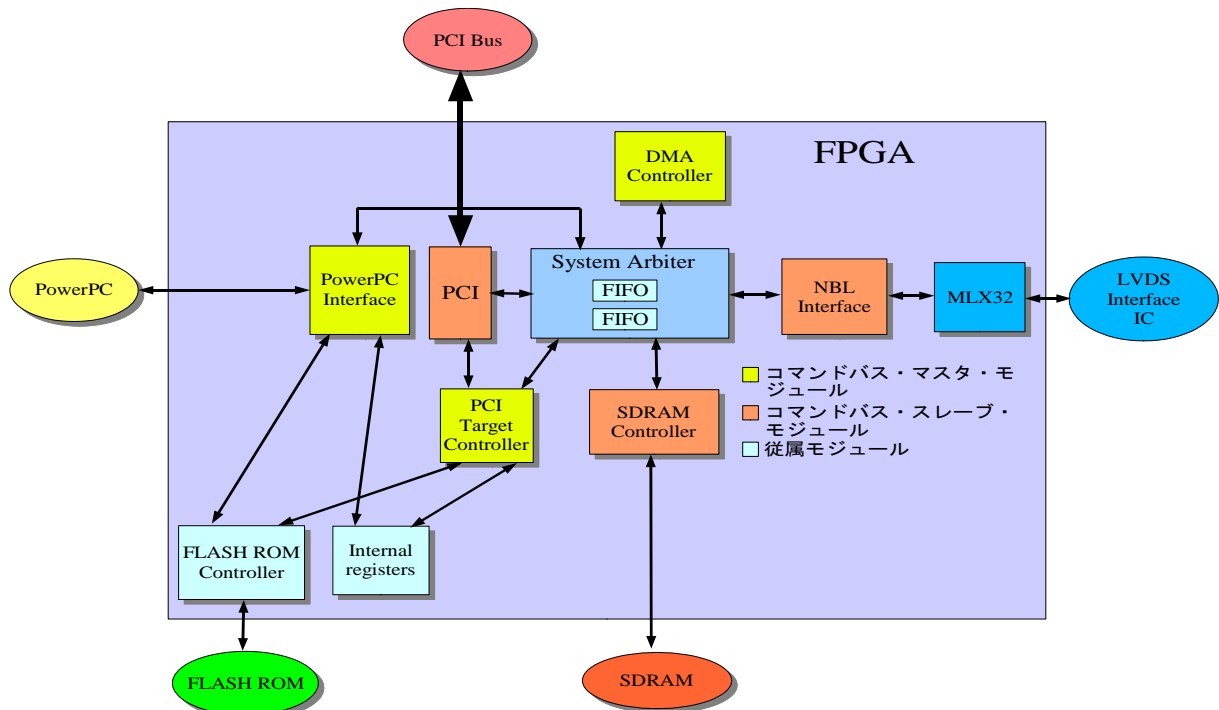


図 3: FPGA の内部モジュール

スト PC から NI へのデータ・アクセス)は、FPGA 回路内でマスタとなる必要があるので、PCI Target Controller に渡され、スレーブ・アクセスを行う。その他、PCI モジュールは NI からホスト PC へのデータ・アクセスを実行する。PowerPC Interface は、PowerPC603e とのバス・インターフェース回路である。バス幅は 64 ビット、バスの動作周波数は 66MHz となっている。Internal Registers は、FPGA の内部の状態などを保持しておく回路である。NBL Interface は、MLX32 とのインターフェース用の回路である。SDRAM Controller は、66MHz 動作の SDRAM を制御する回路である。DMA Controller は、DMA コントローラーが 2 チャンネル実装されていて、各モジュール間の DMA 転送を制御する。詳細はセクション 5 で記述する。

FPGA 回路内には、自分からその他のリソースを要求できるマスタとなれるモジュール(コマンドバス・マスタ・モジュール、以下マスタ・モジュールとする)が 3 つある。それは、PCI Target Controller、PowerPC Controller、DMA

Controller である。マスタ・モジュールの要求で、データ転送を行うのが、コマンドバス・スレーブ・モジュール(以下スレーブ・モジュール)である。DMA Controller は、マスタ・モジュールではあるが、スレーブ・モジュールを 2 つ使用し、それらの間でデータ転送を行わせる特殊なマスタ・モジュールである。スレーブ・モジュールの一方をソース(データを読み出し、出力するモジュール)、もう一方をディスティネーション(データを受け取り、書き込むモジュール)として、データ転送を行う。DMA Controller 自体はデータ・ストリームを送ったり、受けたりはしないで、制御だけを受け持つ。

マスタ・モジュールがどのようにスレーブ・モジュールとデータ転送をするかを説明する。図 4 に FPGA 内部コマンドバスを示す。

1. マスタ・モジュールが、System Arbiter にリソースを要求する。(図 4 の req)

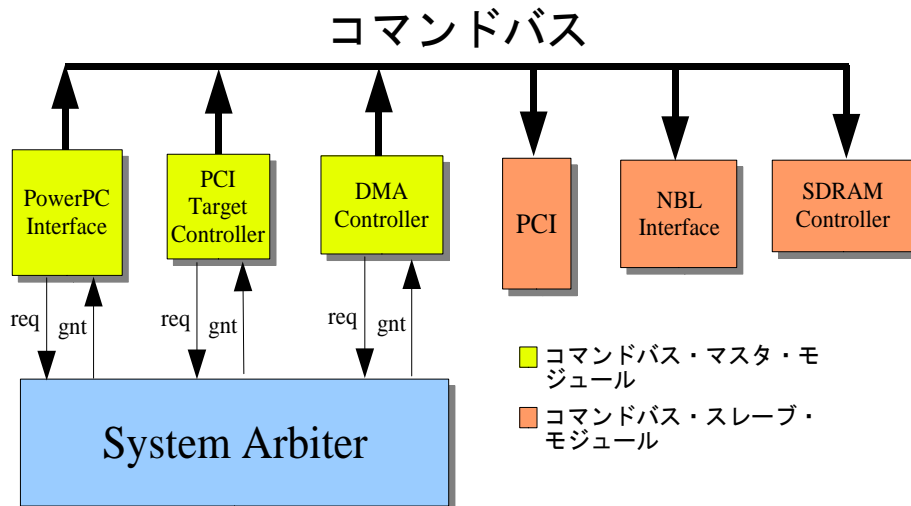


図 4: FPGA 内部のコマンドバス

- リソースが未使用で確保することが出来たら、System Arbiter はマスタ・モジュールに使用を許可する。(図4の gnt)
- System Arbiter は、マスタ・モジュールとスレーブ・モジュール間のデータバスと制御バスを接続する。(セクション4.2で詳述する)
- マスタ・モジュールは、そのリソースのコントローラにコマンドを送りコネクションを確立し、使用を開始する。(図4のコマンドバス)

従属モジュールは、Internal Registers と FLASH ROM Controller の2つで、PowerPC Controller と PCI Target Controller からだけアクセスされる。

4.1 コマンドバス

コマンドバス(図4)はグローバル・コマンドバス、ローカル・コマンドバスとコマンド・イネーブルに分けることが出来る。コマンドバスでコマンドを送るのは、前述したとおり、マスタ・モジュールである。また、コマンドを受けるのはスレーブ・モジュールである。

グローバル・コマンドバスには、各モジュールに共通の項目が設定される。例えば、モジュールのID、ソースまたはディスティネーション、バースト長などである。

ローカル・コマンドバスは、そのスレーブ・モジュールに特有の機能が設定される。例えば、SDRAM Controller ならば、アドレスと、バースト転送の最中にリフレッシュをするかどうかを設定される。

コマンド・イネーブルは、コマンドバスのコマンドが有効かどうかを表す。1の時に有効となる。

Maestro1 では、中央制御回路で全ての制御を行っていたが、この方法は制御線が多くなると、制御手順を頭で考えるのが困難になる。それで、Maestro2 では、各モジュールがコマンドによって、制御される方式にした。こうすることで、各モジュールごとにシミュレーションで動作を確認するのも都合が良い。

4.2 FPGA 内部モジュール間のデータバスと制御バス

FPGA 内部モジュール間のデータバスは64ビット幅のバスを2チャンネル実装してある。現在は、チャンネルごと

にデータを流す方向が決まっている。これは、主にデータ・セクタを節約して、ゲート規模を節約するためである。データバスはSystem Arbiter に実装しており、データバスの内部には、FIFO やバイト・レーン変換をするための回路が実装してある。バイト・レーン変換機能は、8バイトの境界の途中からDMA 転送する場合に、バイト・レーンの入替を行う回路である。(セクション5.8で詳しく述べる)図5にFPGA 内のデータの流れを示す。図5でDMA 転送はバースト転送可能なCH0、CH1のFIFOを通るデータバスを使用する。

制御バスは、データバスに対応して2チャンネル実装されている。制御バスは次に示す3種類の信号がある。

- レディ
- データ転送終了
- バイト・イネーブル

レディは、”データ出力や入力する準備が出来た”ということを表す信号である。各チャンネルごとにソース・レディとディスティネーション・レディがある。

データ転送終了はコマンドで指定されたバースト長のデータ転送が終了したことを表す。これもレディと同様に各チャンネルごとにソースとディスティネーションの信号がある。

バイト・イネーブルは、どのバイト・レーンが有効化を示す信号である。データバスは64ビット、8バイト幅なので、バイト・イネーブルは8ビット幅である。ソースから出力され、ディスティネーションでデバイスに書き込む時にどのバイト・レーンが有効であるかを示す。

5. FPGA 内部モジュール

FPGA 内部モジュールを順に詳しく説明する。

5.1 MLX32

MLX32 は、セクション4.2で記述したデータバスからデータをケーブルが接続されている先のNI やSB に送り出すためにカプセル化する。データは32バイト単位の1つ又は複数のパケットにまとめられ、ヘッダを付けられて、送り出される。ハードウェアによるエラー訂正機能も実装されていて、受信側でエラーが検出されると、送信側は同じパケットを再送信し、正常に受信できるまで再送する。MLX32 はOSI 参照モデルで言うところ

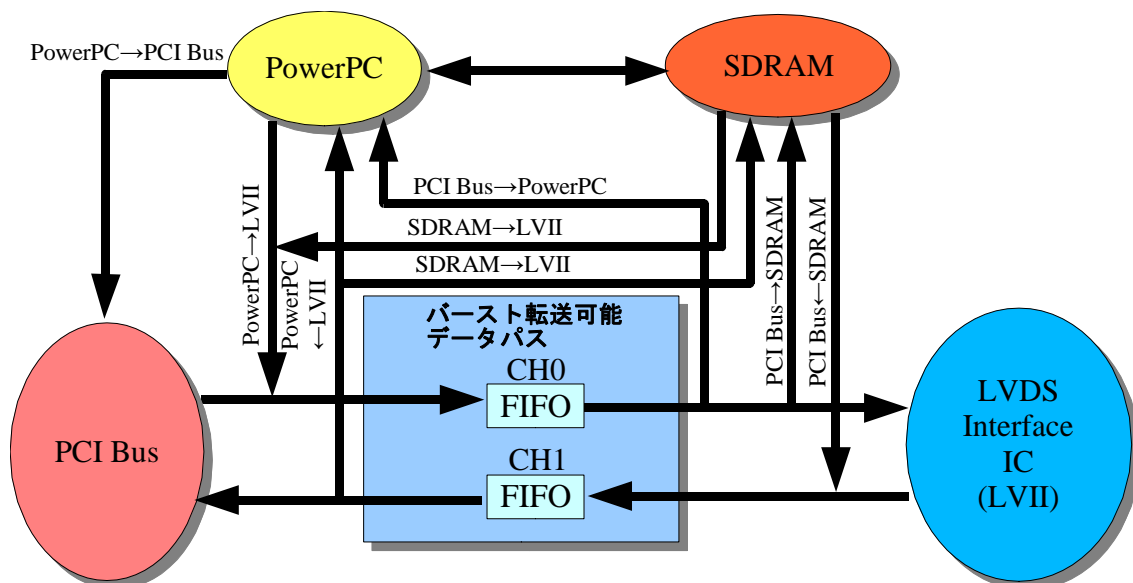


図5: FPGA 内のデータバス

タリンク層に相当する。MLX32の入力は64ビット66MHzで、出力は32ビット100MHzであるので、出力側の最大ビットレートは3.2Gbpsである。MLX32は山際氏が作製した。

5.2 PCI

PCIは64ビット、66MHzのPCIバスとのインターフェースを行う。PCIは、マスタ、スレーブ、コンフィギュレーションの各モジュールに分けられる。

マスタは自分から主体的にデータ転送をするモジュールで、PCIの機能で一番使われる主要なモジュールである。ホストPCメモリからの読み出し、書き込み双方向とも独立にデータ転送することが出来て、66MHz、64ビットでバースト転送することが出来る。

PCI仕様2.2では、64ビット幅のPCIデバイスは、32ビットのターゲットにアクセスする場合に、64ビットアクセスとしてデータ転送を始めていても、32ビットアクセスに変更しなければならないという規定がある。しかし、32ビットアクセスへの変更は、回路の複雑化、動作スピードの低下、およびスループットの低下を招くため、本マスタでは全て64ビット幅でアクセスする。このように、性能向上の観点から、あえてPCI規格を逸脱した実装になっている部分が数ヶ所ある。

スレーブは、ホストPCからのPCIデータ・アクセスを受け取って、NIのリソースへ書き込んだり、リソースからホストPCへ読み出ししたりする。スレーブは64ビットと32ビットの両方のアクセスに対応する。ペンティアム・プロセッサはデータバスが64ビット幅であるが、基本的には、32ビット・プロセッサであり、キャッシュが禁止されているメモリ領域には、32ビット・アクセスが行われる。よって、32ビット・アクセスに対応しないと、データ転送が出来ない。その代わりといっは何だが、スレーブはバースト転送には対応していない。連続的にデータが来ても1つのデータ転送でPCIスレーブ・アクセスを切断する。ペンティアム・プロセッサでは、キャッシュ禁止メモリ領域への読み書き(NIへの読み書きに相当する)は、単一のPCIスレーブ・アクセスとなる。しかし、サーバーワークスのチップセットでは、連続したアドレスへの読み書きではバースト転送となる場合もあるようだ。

コンフィギュレーションは、主にコンフィギュレーション・アクセスの処理とコンフィギュレーション・レジスタの実装を行っている。

ホストPCの電源ON時には、アドレスマップがホストPCのBIOSによって、各PCIデバイスに設定される。このように、PCIのコンフィギュレーション・アクセスは、初期設定をホストPCから各PCIデバイスに設定する作業である。コンフィギュレーションの中には、コンフィギュレーション・アクセスでアクセスするレジスタ(コンフィギュレーション・レジスタ)を実装してある。コンフィギュレーション・レジスタには、ホストPCからPCIスレーブ・アクセスでアクセスする時に必要なベース・アドレスを設定するレジスタなどがある。

5.3 PCI Target Controller

PCI Target Controller(PCITC)は、PCIのスレーブに直結されていて、PCIスレーブ・アクセスのFPGA内部処理を行う。PCIスレーブ・アクセスがあると、PCITCは起動し、FPGA内部バス経由で使用するデバイス用のコントローラ(NBL Controller、SDRAM Controller)にアクセ

スする場合に、前述したようにマスタ・デバイスとしてFPGAの内部バスを要求する。Internal RegistersとFLASH ROM Controllerへアクセスする場合には、PowerPC Interfaceと共用の専用バスがあるので、PowerPC Interfaceとバスの調停を行うだけで、それらを使用することが出来る。

なお、NBL Controller、SDRAM Controllerにアクセスする場合には、PowerPCのPCIマスタ・アクセスとの競合を避けるために、リソースをロックしなければならない。その機能は、Internal Registersに実装してある。

5.4 PowerPC Interface

PowerPC Interfaceは、PowerPCプロセッサ(MPC603e,300MHz,バス動作周波数66MHz)とのインターフェースを実装する。インターフェースの基本構造はMaestro1と同一なので、大部分を流用した。ただし、使用言語をAHDLからVHDLに変更したので、変換を行った。

PowerPCのメモリ・アクセス手順はアドレス転送とデータ転送に分けられる。PowerPC Interfaceはアドレス転送を処理して、デバイスへのアクセス要求を受け取る。ここから、PowerPCのデータ転送が始まる。もし、要求がLVDSインターフェース用ICかPCIバスへのアクセスの場合には、PowerPC Interfaceは、内部バスを使用するので、System Arbiterにデバイスを要求する。その後、デバイスが使用可能になった時に、許可を受け取る。次に、PowerPC Interfaceは、そのデバイスの管理を受け持つコントローラとデータをやり取りする。データのやり取りが終了したら、PowerPCのデータ転送が終了し、完了を送出して終了する。

PowerPCの要求がInternal RegistersかFLASH ROMへのアクセスの場合は、専用のアービタでPCIスレーブとの調停をおこなった後に、データ転送を行う。この場合のデータ転送は、全て単一転送となる。

5.5 Internal Registers

Internal Registersは、FPGA内部の制御レジスタや状態レジスタ、PCIバス用のコンフィギュレーション・レジスタなどが実装されている。Internal Registersは、PowerPCとPCIスレーブ・アクセスから専用バスで書き込み、読み出しができるが、バースト転送は出来ずに、単一転送のみとなる。

アドレス・オフセットが00h(16進数)から3Fhまでは、PCIバスのコンフィギュレーション領域となっている。デバイスID、ベンダID、ベースアドレス・レジスタなどの、各種レジスタが実装されている。

アドレス・オフセットが40hからFFhまでは、FPGA内部の22個の制御、状態レジスタが実装されている。DMA設定用レジスタもここにある。面白い機能としては、FPGAを初期化する設定ビットもあり、ここを設定すると、もう一度FLASH ROMからコンフィギュレーション・データをFPGAに読み込んで、FPGAが完全に初期化される。

アドレス・オフセットが100hからFFFhまでは、内部Block RAM領域である。FPGA内部のBlock RAMをSRAM(Static Random Access Memory)としてマップしてある。Internal Registersとしてのアクセスなので、単一転送のみであるが、ホストPCからアクセスする時に、SDRAMのようにリソースの競合を防ぐためのリソース・

ロックが必要ないため、結果としてSDRAMよりも高速にアクセスが出来る。

5.6 NBL Controller

NBL Controllerは、MLX32とのインターフェースを行う。NBL Controllerは、入出力共に3エントリのFIFOを持っていて、MLX32と内部データバスのデータのプロトコル変換を行う。特にMLX32からのデータ読み出しの場合は、イネーブルを入れてから、1クロック後にデータが出てくるので、内部バスのプロトコルと合わない。よって、NBL Controllerでイネーブルと同時にデータが出力されるように変換を行う。

NBL Controllerは、MLX32のバッファに読み出すデータがあるかどうか、また、MLX32のバッファにデータが書き込めるかどうかを常時監視している。その結果、データ転送が行えない場合には、イネーブルを落として、FPGA内の伝送相手のデバイスを停止する。

5.7 SDRAM Controller

SDRAM Controllerは、66MHz動作、64ビット幅、64MbyteのSDRAMを制御する。SDRAMとPowerPCのデータバスは共通となっていて、PowerPCとのやり取りでは、PowerPCとSDRAM間で直接データをやり取りし、FPGA内とのデータのやり取りはない。(図2参照)

SDRAMの初期化時には、モードレジスタにバースト長やアクセス方法などを設定する。それにより、バースト長が決定されるが、本SDRAM Controllerではバースト長を1に設定してある。それは、PowerPCの非キャッシュメモリ領域のアクセスの時に、単一転送とすることが出来るからである。バースト転送時の動作は以下で説明する。

SDRAMは、行アドレスと列アドレスの2つのアドレスを入力する。始めに行アドレスを入れておいて、その後、READ又はWRITEコマンドと共に列アドレスを入力する。ここで転送を終了させれば、バースト転送長は1にセットされているため単一転送となる。バースト転送時にはREAD又はWRITEコマンドと列アドレスを連続的に入れることで、バースト転送を実現している。バースト転送は列アドレスが9ビットなので、512個までしか行えない。それ以上は、一度バースト転送を中断して、もう一度、行アドレスを入れなおさなければならぬ。

SDRAMは、ダイナミック・メモリなので、そのままではデータが消えてしまう。データを消さないためには、15.6usecごとにリフレッシュを行って、データの再書き込みを行う必要がある。SDRAMとPCIバス間のバースト・データ転送の際にこの点が問題となった。定期的なリフレッシュを行う必要がSDRAMにはあるが、PCIバスは、どの位の時間でデータ転送が終了するかは、ホストPCの都合による。つまり、バースト転送中にリフレッシュを必要とするタイミングになってしまうかもしれない。そこで、その時には、PCIバスの方を待たせておき、SDRAMの方はいったんデータ転送を打ち切って、リフレッシュを行うこととした。リフレッシュ後に、SDRAM Controllerは残りのアドレスを出力し、続きを行う。この辺が一番苦勞した箇所である。リフレッシュによる中断そして再開が15.6usecごとにあるために、なかなかバグ

を取りきれなかった。バグが出る条件の特定に時間がかかってしまった。

5.8 DMA Controller

DMA ControllerはDMAの制御を受け持っている。DMA Controllerは制御だけを受け持ち、実際のデータはSystem Arbiter内にあるFIFOを経由する。図5に示した”バースト転送可能データバス”がそれである。そのFIFOはバイト・レーン変換機能を持っている。

DMA Controllerは、2つのスレーブ・デバイスにコマンドを送って、相互間でデータ転送を行わせる。最初に、ソース・デバイス、ディスティネーション・デバイス、バースト長などをInternal Registers内のDMA設定レジスタに書き込むことによって、DMA転送は、起動する。

DMA転送のチャンネルは前述してきたように2チャンネルあるので、設定用のレジスタも2チャンネル分用意されている。各チャンネルは、転送方向が決まっている。チャンネル0はPCIから読み出してLVDSへ出力、PCIから読み出してSDRAMへ書き込み、SDRAMから読み出してLVDSへ出力の3種類のDMA転送をサポートする。(図5参照)

チャンネル1はLVDS入力からPCIへ出力、LVDS入力からSDRAMへ書き込み、SDRAMから読み出してPCIへ出力の3種類である。(図5参照)

このように、転送可能デバイスを分けることによって、データ入力用マルチプレクサの回路規模が小さくなるため、FPGAのリソースを節約することが出来る。

DMA転送のバースト長の最大値は、64Mbyte-1である。これは、主にSDRAMの容量が64Mbyteであるためである。設定用のレジスタに書かれたDMA転送は、DMA Controllerで、内部的に最大256バーストの転送単位に分けられて実行される。特に、SDRAMは列アドレスが9ビットで512のアドレスであるため、512の境界では、DMA転送を終了して、再起動しなくてはならないためでもある。

その他に、DMA転送は便利な機能を持っている。それはバイト・レーンの変換機能である。

ペンティアム・プロセッサ、PowerPCプロセッサやNIの内部バスなどは64ビット幅のデータバスを持っている。通常のCPU(Central Processing Unit)は、8ビットCPUからの慣習で8ビット単位のデータの塊(バイト)を1単位として、アドレスを付加している。64ビットのデータバス幅を持つCPUで、バイト単位にアクセスすると

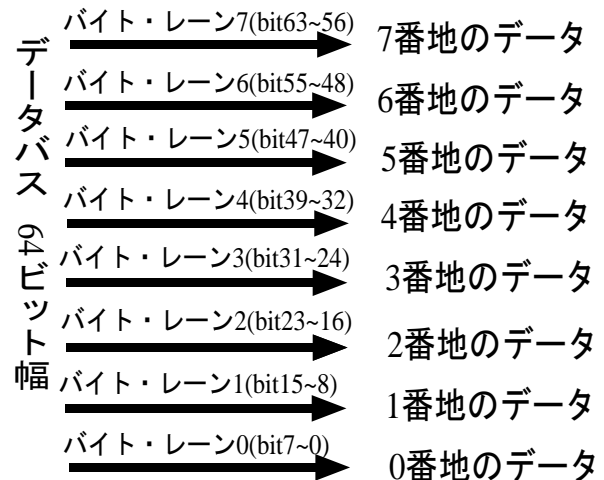
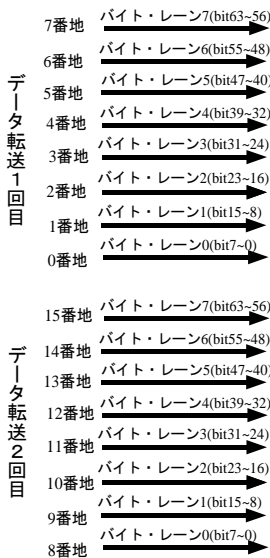
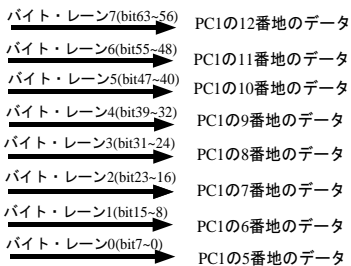


図 6: x86 プロセッサのバイト・レーン

PC1のバイト・レーン配置



PC1とPC2間のバイト・レーン配置



バイト・レーン変換機能付きFIFOでバイト・レーンを変換

PC2のバイト・レーン配置

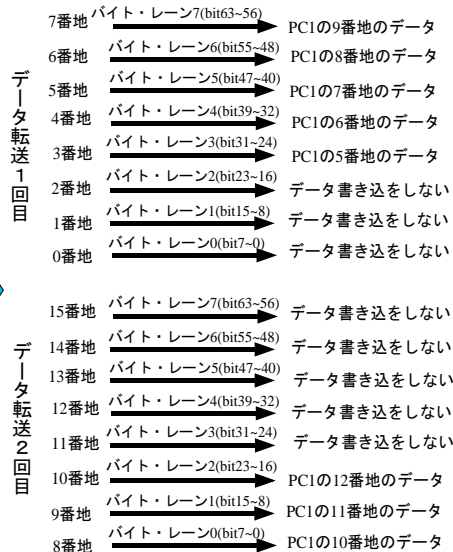


図 7: バイト・レーン変換機能を使用したデータ転送

64 ビット幅のうちの8 ビットを使用することになる。例えば0 番地から7 番地までのバイト・アクセスでは、64 ビット幅のうち、違った場所のデータバスを使用することになる。これをバイト・レーンと呼ぶ。バイト・レーンを図 6 に示す。

以上のことを踏まえて、バイト単位でデータを読んで、それを違った番地には書き込むことを考えてみよう。まず、8 バイトの境界から、8 バイト単位でアクセスしていれば、問題はない。しかし、そうでないと、単純なデータ転送では問題の生ずる場合がある。例えば5 番地のデータをその他のPC に送ろうとした時に、送り先のPC では3 番地にデータを入れたいとしよう。その場合には、5 番地のバイト・レーンと3 番地のバイト・レーンが異なるため、単純に64 ビット幅のデータを転送していたので、正しくデータを転送することが出来ない。

図 6 で説明すると、5 番地はバイト・レーン5 だが、3 番地はバイト・レーン3 になり、データを64 ビット幅でそのまま転送すると、自分のPC の5 番地からのデータを相手方のPC の3 番地からに転送することは無理である。その問題を解決するために、DMA 転送時にバイト・レーン変換機能を使用する。バイト・レーン変換機能は

前述した System Arbiter 内にある FIFO に実装されている。これを BLTFIFO と呼ぶことにする。

バイト・レーン変換機能を説明するために、自分のPC(PC1 とする)の5 番地から12 番地のデータを他のPC(PC2 とする)の3 番地から10 番地に送るとしよう。まずは、PC1 で0 番地から15 番地までのデータをDMA で読み出す。次に BLTFIFO は0~4 番地、13~15 番地までのデータは捨てて、5~12 番地までのデータをバイト・レーン0 から詰めなおし(この転送を Packed 転送と呼ぶことにする)、詰めなおしたデータをPC2 に送る。PC2 の BLTFIFO は、3 番地のバイト・レーン(バイト・レーン3)にもう一度、詰めなおして、PC2 のメモリにDMA で書き込む。その際に、データは8 バイト分だけなので、余計な分(0~2 番地、11~15 番地)は、書き込まない(この転送を Unpacked 転送と呼ぶことにする)。その様子を図 7 に示す。

BLTFIFO はバイト・レーン変換機能と FIFO から構成されている。バイト・レーン変換機能は図 8 の DFF0~7、MUX1、TR0~6、MUX2-0~7 から構成される。

前述の Packed 転送の場合は、まずは必要な5~7 番地のデータを MUX1 でバイト・レーンを変換して

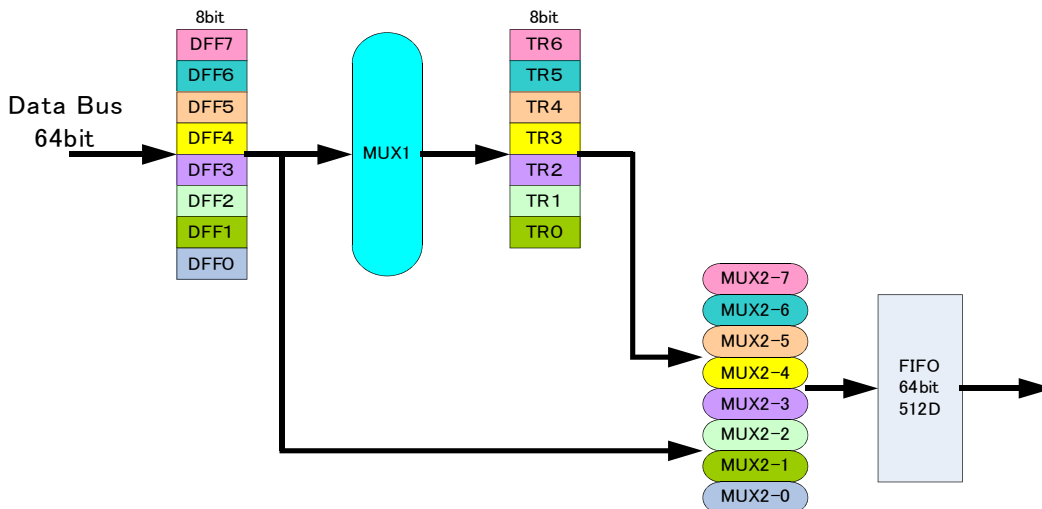


図 8: BLTFIFO の構成

TR0~2にラッチする。次のクロックでDFFから出力されている8~15番地のデータから8~12番地のデータをMUX2-3~7で選択してFIFOに入力する。MUX2-0~2はTR0~2のラッチされたデータ(5~7番地のデータ)を選択してFIFOに入力する。こうして、FIFOには5~12番地の64ビット幅1個に詰め替えられたデータが入力され、他のPCに転送される。

Unpacked転送の場合は、PC1の5~12番地のデータのうち5~9番地のデータだけが(図8のDFF0~4)はMUX2-0~4を通してFIFOに入力される。PC1の10~12番地のデータ(図8のDFF5~7)はTR0~2にラッチされる。次のクロック(つまりFIFOの次のエントリ)で、TR0~2にラッチされたデータは、MUX2-0~2を通してFIFOに入力される。入力されなかったデータは不定となる。

6. まとめ

本報告では、クラスタ向けのネットワークインターフェース(NI)の開発について述べた。本NIは、PCメモリ上のデータを0コピーで転送する機能により、高速にデータを他のPCへ伝送できる。NIのピンポン転送時の実効最大スループットは2.4Gpbs(Giga Bit Per Second)である。

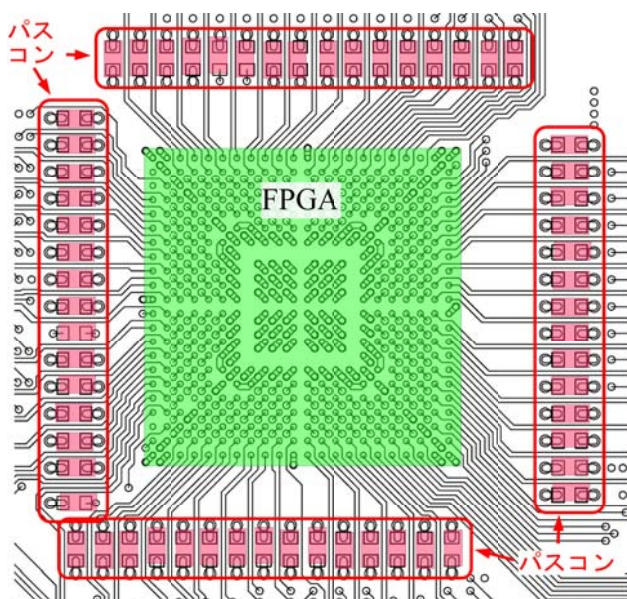


図 9: 部品実装面のプリント配線パターン図面

NI制作上の反省点としては、FPGAの電源用バイパス・コンデンサ(パスコン)の位置が適切でなかったことが挙げられる。パスコンは当初あまり重要だとは考えてなかったため、パターン設計会社の実装位置を決定してもらった。図9に部品実装面のプリント配線パターン図面の一部を示す。図9の緑色の四角形がFPGA、赤色の長方形がパスコンである。パスコンはFPGAの周りに実装されているのがわかる。このパスコンの位置では、FPGAの電源からの距離が長くなり、電源のインピーダンスが低くならないという欠点があることがわかった。その結果として、100MHz動作のLVDSインターフェース用ICへのデータパスにグラント・バウンズによるデータの誤りが生じてしまった。現在は、IOの出力電流値を12mAから4mAに変更したところ、データの誤りはなくなった。本来ならば、パスコンの大きさを小さくして、裏

面のFPGAの直下に実装し、FPGAの電源からの距離をより小さくしなければならない。

7. 謝辞

Maestro2システムの開発補助をさせていただいたシステム情報工学研究科コンピュータ・サイエンス専攻の和田耕一教授に深く感謝いたします。また、共同制作者であるINESC-IDの山際伸一氏とシステム情報工学研究科の青木圭一氏に深く感謝いたします。

参考文献

- [1] 小野雅晃, Xilinx社製FPGAを搭載したPCIボードのシミュレーション, 平成15年度高エネルギー加速器研究機構技術研究会報告集ポスターセッション(2004)
- [2] Shinichi Yamagiwa, Keiichi Aoki, Masaaki Ono, Tetsuya Sakurai, Koichi Wada, and Luis Miguel Campos. Maestro2: A new challenge for high performance cluster network. In The 6th World Multiconference on Systemics, Cybernetics and Informatics, volume XI, Computer Science II, pp.382-387, 2002.
- [3] Keiichi Aoki, Shinichi Yamagiwa, Masaaki Ono, Koichi Wada, Luis Miguel Campos. An architecture of high performance cluster network: Maestro2. In 2003 IEEE Pacific Rim Conference on Communications, Computers and Signal Processing(PacRim03), 2003.
- [4] Shinichi Yamagiwa, Kevin Ferreira, Luis Miguel Campos, Keiichi Aoki, Masaaki Ono, Koichi Wada, Munehiro Fukuda, Leonel Sousa. On the Performance of Maestro2 High Performance Network Equipment, Using New Improvement Techniques. In 23rd IEEE International Performance Computing and Communications Conference(IPCCC 2004), 2004.
- [5] Keiichi Aoki, Shinichi Yamagiwa, Kevin Ferreira, Luis Miguel Campos, Masaaki Ono, Koichi Wada, Leonel Sousa. Maestro2: High Speed Network Technology for High Performance Computing. In 2004 IEEE International Conference on Communications (ICC 2004), 2004.