

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-197319

(P2019-197319A)

(43) 公開日 令和1年11月14日(2019.11.14)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 9/50 (2006.01)	G06F 9/50 150D	5B042
G06F 11/34 (2006.01)	G06F 11/34 133	
G06F 13/14 (2006.01)	G06F 11/34 119	
	G06F 13/14 310F	

審査請求 未請求 請求項の数 8 O L (全 15 頁)

(21) 出願番号 特願2018-90015 (P2018-90015)
 (22) 出願日 平成30年5月8日 (2018.5.8)

(71) 出願人 00005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番1号
 (74) 代理人 110002147
 特許業務法人酒井国際特許事務所
 (72) 発明者 三吉 貴史
 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
 Fターム(参考) 5B042 JJ20 MA14 MC33

(54) 【発明の名称】 情報処理装置、情報処理方法および情報処理プログラム

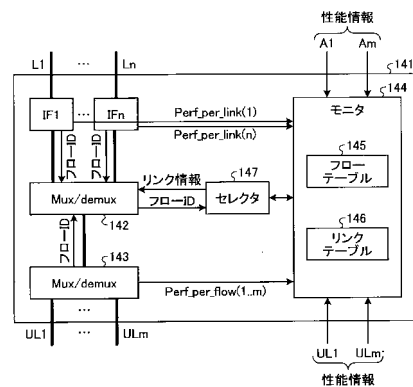
(57) 【要約】

【課題】システム全体の性能を最大化するようにフローを振り分けることができる情報処理装置、情報処理方法および情報処理プログラムを提供する。

【解決手段】情報処理装置は、複数のリンクを経由してCPUと接続されるオフロード回路を備える。情報処理装置のオフロード回路は、論理回路と、収集部と、選択部とを備える。論理回路は、アプリケーションの処理を演算する。収集部は、アプリケーションの処理に対応するフローごとのリンクの性能情報を示す値と、リンクごとの使用可能な性能情報の最大値とを収集する。選択部は、フローごとのリンクの性能情報を示す値に基づいて、要する性能情報を満たしていないフローを判定する。また、選択部は、リンクごとの使用可能な性能情報の最大値と、リンクごとの現在使用されている性能情報の値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。

【選択図】 図4

制御回路の機能構成の一例を示すブロック図



【特許請求の範囲】**【請求項 1】**

複数のリンクを経由してCPUと接続されるオフロード回路において、アプリケーションの処理を演算する論理回路と、前記アプリケーションの処理に対応するフローごとの前記リンクの性能情報を示す値と、前記リンクごとの使用可能な性能情報の最大値とを収集する収集部と、前記フローごとの前記リンクの性能情報を示す値に基づいて、要する性能情報を満たしていないフローを判定し、前記リンクごとの使用可能な性能情報の最大値と、前記リンクごとの現在使用されている性能情報の値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける選択部と、

10

を備えることを特徴とする情報処理装置。

【請求項 2】

前記収集部は、前記アプリケーションの処理に対応するフローごとの前記リンクの使用帯域を示す値と、前記リンクごとの使用可能な帯域の最大値とを収集し、

前記選択部は、前記フローごとの前記リンクの使用帯域を示す値に基づいて、要する帯域を満たしていないフローを判定し、前記リンクごとの使用可能な帯域の最大値と、前記リンクごとの現在使用されている帯域の値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける、

ことを特徴とする請求項 1 に記載の情報処理装置。

20

【請求項 3】

前記収集部は、前記アプリケーションの処理、または、前記論理回路の性能が性能要件を満たしていない場合、対応する前記フローの優先度を上げる、

ことを特徴とする請求項 2 に記載の情報処理装置。

【請求項 4】

前記選択部は、前記リンクごとの使用可能な帯域の最大値と、前記リンクごとの現在使用されている帯域の値とに基づいて、前記フローの振り分け先のリンクを選択し、選択した前記振り分け先のリンクに該フローを振り分けると、前記振り分け先のリンクの使用可能な帯域の最大値を超える場合、前記振り分け先のリンクを使用する最も優先度が低いフローを、振り分け元のリンクに振り分ける、

ことを特徴とする請求項 2 または 3 に記載の情報処理装置。

30

【請求項 5】

前記収集部は、所定時間ごとに、前記リンクの使用帯域を示す値と、前記リンクごとの使用可能な帯域の最大値とを収集し、

前記選択部は、前記所定時間ごとに、前記要する帯域を満たしていないフローを判定する、

ことを特徴とする請求項 2 ~ 4 のいずれか 1 つに記載の情報処理装置。

【請求項 6】

前記収集部は、前記アプリケーションの処理に対応するフローごとの前記リンクのレイテンシを示す値と、前記リンクごとのレイテンシの最大値とを収集し、

前記選択部は、前記フローごとの前記リンクのレイテンシを示す値に基づいて、要するレイテンシを満たしていないフローを判定し、前記リンクごとのレイテンシの最大値と、前記リンクごとの現在のレイテンシの値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける、

40

ことを特徴とする請求項 1 に記載の情報処理装置。

【請求項 7】

複数のリンクを経由してCPUと接続され、アプリケーションの処理を演算する論理回路を備えるオフロード回路において、

前記アプリケーションの処理に対応するフローごとの前記リンクの性能情報を示す値と、前記リンクごとの使用可能な性能情報の最大値とを収集し、

前記フローごとの前記リンクの性能情報を示す値に基づいて、要する性能情報を満たし

50

ていないフローを判定し、前記リンクごとの使用可能な性能情報の最大値と、前記リンクごとの現在使用されている性能情報の値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける、

ことを特徴とする情報処理方法。

【請求項 8】

複数のリンクを経由して CPU と接続され、アプリケーションの処理を演算する論理回路を備えるオフロード回路において、

前記アプリケーションの処理に対応するフローごとの前記リンクの性能情報を示す値と、前記リンクごとの使用可能な性能情報の最大値とを収集し、

前記フローごとの前記リンクの性能情報を示す値に基づいて、要する性能情報を満たしていないフローを判定し、前記リンクごとの使用可能な性能情報の最大値と、前記リンクごとの現在使用されている性能情報の値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける、

処理をコンピュータに実行させることを特徴とする情報処理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、情報処理方法および情報処理プログラムに関する。

【背景技術】

【0002】

近年、CPU (Central Processing Unit) で動作するアプリケーションを、高速化のために FPGA (Field Programmable Gate Array) 等のハードウェアにオフロードして動作させることが提案されている。オフロードしたアプリケーションの一部は、ユーザロジック (オフロード回路) と呼ばれ、1つの FPGA に複数のユーザロジックを搭載する場合がある。複数のユーザロジックは、それぞれ CPU で動作するアプリケーションに応じて、任意の組み合わせで動作することが想定されている。CPU と FPGA との間は、例えば、PCI Express (以下、PCIe ともいう。) 等の I/O バス経由での接続と、OpenCAPI (Open Coherent Accelerator Processor Interface) 等のメモリコヒーレントバス経由での接続とがある。

【0003】

I/O バス経由の接続は、アクセス遅延が大きく、CPU でキャッシュできないが、帯域は中から広帯域であり、FPGA 側からの DMA (Direct Memory Access) 転送に向く。一方、メモリコヒーレントバス経由の接続は、アクセス遅延が小さく、CPU でキャッシュできるが、帯域は小から中帯域であり、FPGA 側からのメモリアクセスが可能である。CPU と FPGA との間の接続 (リンク) は、I/O バス経由の接続と、メモリコヒーレントバス経由の接続との両方が用意される場合がある。この場合、ユーザロジックは、どのリンクを使用するかを事前に固定的に選択する。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2014 - 170363 号公報

【特許文献 2】特表 2008 - 546072 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、ユーザロジックは、他のユーザロジックがどのリンクを使用するのが運用時まで判らない。このため、例えば、ユーザロジック、コマンド種別およびアドレス範囲の組み合わせを表すフローを単位として各リンクに振り分ける場合、フローが要求する帯域に対してリンクの帯域が不足する等により、アプリケーションの性能要求を満たさない場合がある。

10

20

30

40

50

【 0 0 0 6 】

一つの側面では、システム全体の性能を最大化するようにフローを振り分けることができる情報処理装置、情報処理方法および情報処理プログラムを提供することにある。

【課題を解決するための手段】

【 0 0 0 7 】

一つの態様では、情報処理装置は、複数のリンクを経由してCPUと接続されるオフロード回路を備える。情報処理装置のオフロード回路は、論理回路と、収集部と、選択部とを備える。論理回路は、アプリケーションの処理を演算する。収集部は、前記アプリケーションの処理に対応するフローごとの前記リンクの性能情報を示す値と、前記リンクごとの使用可能な性能情報の最大値とを収集する。選択部は、前記フローごとの前記リンクの性能情報を示す値に基づいて、要する性能情報を満たしていないフローを判定する。また、選択部は、前記リンクごとの使用可能な性能情報の最大値と、前記リンクごとの現在使用されている性能情報の値とに基づいて、前記フローの振り分け先のリンクを選択して振り分ける。

10

【発明の効果】

【 0 0 0 8 】

システム全体の性能を最大化するようにフローを振り分けることができる。

【図面の簡単な説明】

【 0 0 0 9 】

【図 1】図 1 は、複数のリンクでCPUとFPGAを接続する場合の一例を示す図である。

20

【図 2】図 2 は、実施例の情報処理装置のハードウェア構成の一例を示すブロック図である。

【図 3】図 3 は、実施例の情報処理装置の機能構成の一例を示すブロック図である。

【図 4】図 4 は、制御回路の機能構成の一例を示すブロック図である。

【図 5】図 5 は、フローテーブルの一例を示す図である。

【図 6】図 6 は、リンクテーブルの一例を示す図である。

【図 7】図 7 は、実施例の経路制御処理の一例を示すフローチャートである。

【図 8】図 8 は、実施例の切替処理の一例を示すフローチャートである。

【発明を実施するための形態】

30

【 0 0 1 0 】

以下、図面に基づいて、本願の開示する情報処理装置、情報処理方法および情報処理プログラムの実施例を詳細に説明する。なお、本実施例により、開示技術が限定されるものではない。また、以下の実施例は、矛盾しない範囲で適宜組みあわせてもよい。

【実施例】

【 0 0 1 1 】

まず、図 1 を用いてリンクの帯域が不足する場合について説明する。図 1 は、複数のリンクでCPUとFPGAを接続する場合の一例を示す図である。図 1 の例では、CPU 10 にはメモリ 11 が接続され、CPU 10 とFPGA 12 との間のリンクとして、PCIe 経由の接続と、OpenCAPI 経由の接続との両方が用意される。ユーザロジックUL 201 ~ UL 203 は、Mux / Demux 13 を介して、PCIe およびOpenCAPI と接続される。

40

【 0 0 1 2 】

ここで、PCIe は、Gen 4 x 16 lane x 2 slot であり、帯域が 51.2 GB / s、アクセス遅延が 1 μ s ~ であるとする。また、OpenCAPI は、帯域が 25.6 GB / s、アクセス遅延が 100 ns ~ であるとする。ユーザロジックの要求性能は、ユーザロジックUL 201 は、必要帯域 10 GB / s、許容アクセス遅延 200 ns である。ユーザロジックUL 202 は、必要帯域 40 GB / s、許容アクセス遅延 2 μ s である。ユーザロジックUL 203 は、必要帯域 20 GB / s、許容アクセス遅延 1 μ s である。

50

【 0 0 1 3 】

ユーザロジックUL201～UL203は、事前にどのリンクを使用するかは固定選択であるとし、運用時まで他のユーザロジックがどれだけリンクを使用するか不明であるとする。この場合、図1の例では、ユーザロジックUL201は、OpenCAPIに割り振り、ユーザロジックUL202は、PCIeに割り振る。ところが、ユーザロジックUL203は、OpenCAPIとPCIeとのうち、どちらに割り振っても、必要帯域がリンクの帯域を上回ることになる。つまり、図1は、フローが要求する帯域に対してリンクの帯域が不足し、アプリケーションの性能要求を満たさない場合の例である。これに対し、動的にユーザロジックUL203の通信を2つの性質の異なるリンクに割り振ることができれば、システム全体としての性能を向上させることができる。

10

【 0 0 1 4 】

図2は、実施例の情報処理装置のハードウェア構成の一例を示すブロック図である。図2に示すように、情報処理装置100は、通信部110と、表示部111と、操作部112と、HDD(Hard Disk Drive)113と、メモリ120と、CPU130と、FPGA140とを有する。なお、通信部110、表示部111、操作部112、HDD113、CPU130およびFPGA140は、バス114を介して相互に接続される。また、CPU130は、メモリ120およびFPGA140と接続される。なお、情報処理装置100は、図2に示す機能部以外にも既知のコンピュータが有する各種の機能部、例えば各種の入力デバイスや音声出力デバイス等の機能部を有することとしてもかまわない。

20

【 0 0 1 5 】

通信部110は、例えば、NIC(Network Interface Card)等によって実現される。通信部110は、図示しないネットワークを介して他の情報処理装置と有線または無線で接続され、他の情報処理装置との間で情報の通信を司る通信インタフェースである。

【 0 0 1 6 】

表示部111は、各種情報を表示するための表示デバイスである。表示部111は、例えば、表示デバイスとして液晶ディスプレイ等によって実現される。表示部111は、CPU130から図示しない表示制御部を介して入力された表示画面等の各種画面を表示する。

【 0 0 1 7 】

操作部112は、情報処理装置100のユーザから各種操作を受け付ける入力デバイスである。操作部112は、例えば、入力デバイスとして、キーボードやマウス等によって実現される。操作部112は、ユーザによって入力された操作を操作情報としてCPU130に出力する。なお、操作部112は、入力デバイスとして、タッチパネル等によって実現されるようにしてもよく、表示部111の表示デバイスと、操作部112の入力デバイスとは、一体化されるようにしてもよい。

30

【 0 0 1 8 】

HDD113は、補助記憶装置であり、CPU130で動作するOS(Operating System)や各種データを記憶する。なお、HDD113は、ハードディスクドライブの他にもフラッシュメモリ等の半導体メモリ素子を用いたSSD(Solid State Drive)や光ディスク等の記憶装置によって実現されてもよい。

40

【 0 0 1 9 】

メモリ120は、主記憶装置であり、例えば、各種のSDRAM(Synchronous Dynamic Random Access Memory)等のようなRAM(Random Access Memory)等の半導体メモリ素子等の記憶装置によって実現される。メモリ120は、CPU130での処理に用いる情報を記憶する。なお、メモリ120は、バス114に接続してもよいし、FPGA140と直接接続してもよい。

【 0 0 2 0 】

CPU130は、HDD113等の記憶部に記憶されているプログラムに従って、メモリ120等のRAMを作業領域として各種処理を実行する。すなわち、CPU130は、OSやVM(Virtual Machine)によって制御され、各種処理を実行する。

50

【0021】

FPGA140は、CPU130で動作するアプリケーションをオフロードするユーザロジックを動作させる。FPGA140は、CPU130とバス114（例えば、PCIe。）経由、および、メモリコヒーレントバス（例えば、OpenCAPI。）経由で接続される。すなわち、FPGA140は、複数のリンクを経由してCPU130と接続されるオフロード回路を構成する。なお、複数のリンクは、異なる種類のバスでなく、同じ種類のバスであってもよい。

【0022】

図3は、実施例の情報処理装置の機能構成の一例を示すブロック図である。図3に示すように、情報処理装置100では、CPU130でOS/VM131が動作し、さらに、OS/VM131上でアプリケーションA1~Amが動作している。また、FPGA140は、制御回路141と複数のユーザロジックUL1~ULmとを有する。制御回路141は、リンクL1~LnでCPU130と接続される。リンクL1~Lnは、バス114およびメモリコヒーレントバスに対応する。ユーザロジックUL1~ULmは、アプリケーションA1~Amと対応付けられており、対応するアプリケーションA1~Amとの間で、制御回路141およびリンクL1~Lnを介して通信を行う。つまり、ユーザロジックUL1~ULmは、アプリケーションA1~Amの処理を演算する論理回路である。また、アプリケーションA1~Am、および、ユーザロジックUL1~ULmは、使用するリンクL1~Lnのいずれかを用いて、それぞれ性能情報を制御回路141に出力する。なお、性能情報は、例えば、性能要件を満たす場合を「1」とし、性能要件を満たさない場合を「0」とする。

10

20

【0023】

図4は、制御回路の機能構成の一例を示すブロック図である。図4に示すように、制御回路141は、インタフェースIF1~IFnと、Mux/demux142と、Mux/demux143と、モニタ144と、セクタ147とを有する。また、モニタ144は、フローテーブル145と、リンクテーブル146とを有する。

【0024】

インタフェースIF1~IFnは、リンクL1~Lnに対応するインタフェースである。インタフェースIF1~IFnは、リンクごとの使用帯域を表すデータであるPerf_per_link(1)~(n)を、それぞれモニタ144に出力する。また、インタフェースIF1~IFnは、フローを識別するフローID(Identifier)をMux/demux142に出力する。なお、フローは、ユーザロジックUL1~ULm、コマンド種別(Read/Write)およびアドレス範囲の組み合わせに基づくデータの流れを表し、フローIDを用いて各フローを識別する。

30

【0025】

Mux/demux142は、マルチプレクサおよびデマルチプレクサである。Mux/demux142は、インタフェースIF1~IFn側から入力されるフローを多重化し、Mux/demux143に出力する。また、Mux/demux142は、Mux/demux143から入力される多重化されたフローを複数のフローに戻し、対応するインタフェースIF1~IFnに出力する。Mux/demux142は、インタフェースIF1~IFn側、および、Mux/demux143側から入力された各フローのフローIDをセクタ147に出力する。Mux/demux142には、セクタ147からフローIDに応じたリンク情報が入力される。Mux/demux142は、リンク情報に応じたインタフェースIF1~IFnに、対応するフローを出力するようにフローの経路を切り替える。

40

【0026】

Mux/demux143は、マルチプレクサおよびデマルチプレクサである。Mux/demux143は、ユーザロジックUL1~ULm側から入力されるパケットからフローを抽出して多重化し、Mux/demux142に出力する。また、Mux/demux143は、Mux/demux142から入力される多重化されたフローを複数のフ

50

ローに戻し、対応するユーザロジックUL1～ULmに出力する。Mux/demux143は、各フローに対応するフローIDをMux/demux142に出力する。また、Mux/demux143は、抽出したフローに関するフロー情報と、フローごとの使用帯域とを含むデータであるPerf_per_flow(1...m)をモニタ144に出力する。

【0027】

モニタ144は、リンクL1～Lnの使用帯域と、各フローの使用帯域と、性能情報とを収集し、フローテーブル145に記憶する。ここで、図5および図6を用いて、フローテーブル145およびリンクテーブル146について説明する。

【0028】

フローテーブル145は、フローごとにリンク情報（経路情報）、優先度、性能情報等を対応付けて記憶する。図5は、フローテーブルの一例を示す図である。図5に示すように、フローテーブル145は、「フローID」、「フロー情報」、「リンク情報」、「優先度」、「性能情報（現在値）」といった項目を有する。また、「フロー情報」は、「ユーザロジック」、「コマンド種別」、「アドレス範囲」といった項目を有する。また、「性能情報（現在値）」は、「アプリ性能」、「UL性能」、「使用帯域」といった項目を有する。

10

【0029】

「フローID」は、フローを識別する識別子である。「ユーザロジック」は、当該フローに対応するユーザロジックを識別する識別子である。「コマンド種別」は、Read/Writeの別を示す情報である。「コマンド種別」は、Read/Writeのどちら

20

【0030】

「優先度」は、各フローの振り分けの優先度を示し、例えば、「0～255」の256段階で表すことができる。この場合、数値が大きいほど優先度が高いとしている。「アプリ性能」は、アプリケーションA1～Amから収集した性能情報である。「アプリ性能」は、OKであれば性能要件を満たし、NGであれば性能要件を満たさないことを示す。ここで、「アプリ性能」が性能要件を満たすとは、例えば、アプリケーションからユーザロジックに対してリクエストを送信してレスポンスが返ってくるまでの応答時間やスループット等が条件を満たす場合である。

30

【0031】

「UL性能」は、ユーザロジックUL1～ULmから収集した性能情報である。「UL性能」は、OKであれば性能要件を満たし、NGであれば性能要件を満たさないことを示す。ここで、「UL性能」が性能要件を満たすとは、例えば、ユーザロジックからアプリケーションにレスポンスを送信してACKが返ってくるまでの応答時間や、アプリケーションからの単位時間あたりのリクエスト回数等が条件を満たす場合である。

【0032】

「アプリ性能」および「UL性能」では、OKは、モニタ144が収集した性能情報「1」に対応し、NGは、モニタ144が収集した性能情報「0」に対応する。性能情報は、例えば、モニタ144内に、アプリケーションA1～Am、および、ユーザロジックUL1～ULmに対応するレジスタを設け、それぞれがレジスタに「0」または「1」を書き込み、モニタ144が定期的にレジスタを参照することで収集できる。「使用帯域」は、当該フローが使用しているリンクL1～Lnの帯域を示す情報である。

40

【0033】

リンクテーブル146は、リンクL1～Lnの性能情報を記憶する。図6は、リンクテーブルの一例を示す図である。図6に示すように、リンクテーブル146は、「リンクID」、「性能情報」といった項目を有する。また、「性能情報」は、「最大帯域」、「最小遅延」といった項目を有する。

【0034】

50

「リンクID」は、リンクL1～Lnを識別する識別子である。「最大帯域」は、当該リンクが収容可能な最大の帯域を示す情報である。「最小遅延」は、当該リンクにおける最小の遅延時間を示す情報である。

【0035】

図4の説明に戻って、モニタ144は、使用帯域や性能情報の収集とともに、各フローの優先度を制御する。モニタ144は、電源が投入されると、初期状態を設定する。モニタ144は、初期状態として、フローテーブル145の優先度欄を全て「1」に設定する。モニタ144は、アプリケーションA1～Am、および、ユーザロジックUL1～ULmの動作が開始されると、フローテーブル145のフロー情報、リンク情報および性能情報を更新する。なお、フローIDは、例えば、アプリケーションの起動時に設定されたフローIDをアプリケーションから取得する。また、フローIDは、例えば、アプリケーションの終了時にアプリケーションからの指示により削除される。

10

【0036】

つまり、モニタ144は、Mux/demux143から入力されるフロー情報に基づいて、フローテーブル145のフロー情報を更新する。モニタ144は、リンクテーブル146を参照し、インタフェースIF1～IFnから入力されるリンクごとの使用帯域に基づいて、使用可能な帯域が多いリンクから順にフローを割り振って、フローテーブル145のリンク情報を更新する。モニタ144は、Mux/demux143から入力されるフローごとの使用帯域に基づいて、フローテーブル145のフローごとの使用帯域を更新する。モニタ144は、アプリケーションA1～Am、および、ユーザロジックUL1～ULmに対応するレジスタを参照し、フローテーブル145のアプリ性能およびUL性能を更新する。

20

【0037】

モニタ144は、フローテーブル145のアプリ性能およびUL性能がNGであるフローがある場合、当該フローの優先度をインクリメントする。次に、モニタ144は、セクタ147に対して、切替処理の実行を指示する。

【0038】

モニタ144は、切替処理の終了後、フローテーブル145を参照し、全フローにおいて、アプリ性能またはUL性能がNGであれば、全フローの優先度をデクリメントする。また、モニタ144は、フローテーブル145を参照し、全フローにおいて、アプリ性能およびUL性能がOKであれば、全フローの優先度をデクリメントし、一定時間の待機後、再びアプリ性能およびUL性能を更新して優先度の制御を繰り返す。

30

【0039】

言い換えると、モニタ144は、アプリケーションの処理に対応するフローごとのリンクの性能情報を示す値と、リンクごとの使用可能な性能情報の最大値とを収集する収集部の一例である。また、モニタ144は、アプリケーションの処理に対応するフローごとのリンクの使用帯域を示す値と、リンクごとの使用可能な帯域の最大値とを収集する。また、モニタ144は、アプリケーションの処理、または、論理回路の性能が性能要件を満たしていない場合、対応するフローの優先度を上げる。また、モニタ144は、所定時間ごとに、リンクの使用帯域を示す値と、リンクごとの使用可能な帯域の最大値とを収集する。

40

【0040】

セクタ147は、Mux/demux142からフローIDが入力されると、フローテーブル145を参照し、入力されたフローIDに対応するリンク情報をMux/demux142に出力する。すなわち、セクタ147は、フローテーブル145を参照し、リンクL1～Lnを流れるパケットを、適切なユーザロジックUL1～ULmに接続する。

【0041】

また、セクタ147は、モニタ144から切替処理の実行を指示されると、フローとリンクとの対応を切り替える切替処理を実行する。セクタ147は、フローテーブル1

50

45を参照し、未判定のフローがあるか否かを判定する。セクタ147は、未判定のフローがないと判定した場合には、切替処理を終了する。

【0042】

セクタ147は、未判定のフローがあると判定した場合には、未判定のフローのうち、優先度が最も高いフローを対象フローとして選択する。セクタ147は、選択した対象フローについて、アプリ性能およびUL性能がOKであるか否かを判定する。セクタ147は、アプリ性能およびUL性能がOKであると判定した場合には、次のフローの判定に進む。

【0043】

セクタ147は、アプリ性能およびUL性能がOKでないと判定した場合には、フローテーブル145およびリンクテーブル146を参照し、最も使用可能な帯域が大きいリンクを選択する。セクタ147は、選択したリンクが対象フローのリンクと同じであるか否かを判定する。セクタ147は、選択したリンクが対象フローのリンクと同じであると判定した場合には、次のフローの判定に進む。

10

【0044】

セクタ147は、選択したリンクが対象フローのリンクと同じでないと判定した場合には、対象フローに選択したリンクを設定する。セクタ147は、選択したリンクの使用帯域が最大帯域に収まるか否かを判定する。セクタ147は、選択したリンクの使用帯域が最大帯域に収まると判定した場合には、切替処理を終了する。

【0045】

セクタ147は、選択したリンクの使用帯域が最大帯域に収まらないと判定した場合には、フローテーブル145を参照し、選択したリンクを使用する最も優先度が低いフローを選択する。セクタ147は、選択した最も優先度が低いフローのリンクに、対象フローの元のリンクを設定し、切替処理を終了する。すなわち、セクタ147は、対象フローの元のリンクと、最も優先度が低いフローのリンクとを入れ替える。

20

【0046】

言い換えると、セクタ147は、フローごとのリンクの性能情報を示す値に基づいて、要する性能情報を満たしていないフローを判定する。セクタ147は、リンクごとの使用可能な性能情報の最大値と、リンクごとの現在使用されている性能情報の値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。つまり、セクタ147は、選択部の一例である。すなわち、セクタ147は、フローごとのリンクの使用帯域を示す値に基づいて、要する帯域を満たしていないフローを判定する。セクタ147は、リンクごとの使用可能な帯域の最大値と、リンクごとの現在使用されている帯域の値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。

30

【0047】

また、セクタ147は、リンクごとの使用可能な帯域の最大値と、リンクごとの現在使用されている帯域の値とに基づいて、フローの振り分け先のリンクを選択する。セクタ147は、選択した振り分け先のリンクに該フローを振り分けると、振り分け先のリンクの使用可能な帯域の最大値を超える場合、振り分け先のリンクを使用する最も優先度が低いフローを、振り分け元のリンクに振り分ける。また、セクタ147は、所定時間ごとに、要する帯域を満たしていないフローを判定する。

40

【0048】

次に、実施例の情報処理装置100の動作について説明する。図7は、実施例の経路制御処理の一例を示すフローチャートである。

【0049】

制御回路141のモニタ144は、電源が投入されると、初期状態を設定する(ステップS1)。モニタ144は、アプリケーションA1~Am、および、ユーザロジックUL1~ULmの動作が開始されると、フローテーブル145のフロー情報、リンク情報および性能情報を更新する。モニタ144は、動作開始後、一定時間待機する(ステップS2)。

50

【 0 0 5 0 】

その後、モニタ 1 4 4 は、アプリケーション A 1 ~ A m、および、ユーザロジック U L 1 ~ U L m に対応するレジスタを参照し、フローテーブル 1 4 5 のアプリ性能および U L 性能を更新する（ステップ S 3）。

【 0 0 5 1 】

モニタ 1 4 4 は、フローテーブル 1 4 5 のアプリ性能および U L 性能が N G であるフローの優先度をインクリメントする（ステップ S 4）。次に、モニタ 1 4 4 は、セクタ 1 4 7 に対して、切替処理の実行を指示する（ステップ S 5）。

【 0 0 5 2 】

ここで、図 8 を用いて切替処理について説明する。図 8 は、実施例の切替処理の一例を示すフローチャートである。

10

【 0 0 5 3 】

セクタ 1 4 7 は、モニタ 1 4 4 から切替処理の実行を指示されると、フローとリンクとの対応を切り替える切替処理を実行する。セクタ 1 4 7 は、フローテーブル 1 4 5 を参照し、未判定のフローがあるか否かを判定する（ステップ S 5 1）。セクタ 1 4 7 は、未判定のフローがないと判定した場合には（ステップ S 5 1：否定）、切替処理を終了し、経路制御処理に戻る。

【 0 0 5 4 】

セクタ 1 4 7 は、未判定のフローがあると判定した場合には（ステップ S 5 1：肯定）、未判定のフローのうち、優先度が最も高いフローを対象フローとして選択する（ステップ S 5 2）。セクタ 1 4 7 は、選択した対象フローについて、アプリ性能および U L 性能が O K であるか否かを判定する（ステップ S 5 3）。セクタ 1 4 7 は、アプリ性能および U L 性能が O K であると判定した場合には（ステップ S 5 3：肯定）、次のフローの判定を行うため、ステップ S 5 1 に戻る。

20

【 0 0 5 5 】

セクタ 1 4 7 は、アプリ性能および U L 性能が O K でないと判定した場合には（ステップ S 5 3：否定）、フローテーブル 1 4 5 およびリンクテーブル 1 4 6 を参照し、最も使用可能な帯域が大きいリンクを選択する（ステップ S 5 4）。セクタ 1 4 7 は、選択したリンクが対象フローのリンクと同じであるか否かを判定する（ステップ S 5 5）。セクタ 1 4 7 は、選択したリンクが対象フローのリンクと同じであると判定した場合には（ステップ S 5 5：肯定）、次のフローの判定を行うため、ステップ S 5 1 に戻る。

30

【 0 0 5 6 】

セクタ 1 4 7 は、選択したリンクが対象フローのリンクと同じでないと判定した場合には（ステップ S 5 5：否定）、対象フローに選択したリンクを設定する（ステップ S 5 6）。セクタ 1 4 7 は、選択したリンクの使用帯域が最大帯域に収まるか否かを判定する（ステップ S 5 7）。セクタ 1 4 7 は、選択したリンクの使用帯域が最大帯域に収まると判定した場合には（ステップ S 5 7：肯定）、切替処理を終了し、経路制御処理に戻る。

【 0 0 5 7 】

セクタ 1 4 7 は、選択したリンクの使用帯域が最大帯域に収まらなると判定した場合には（ステップ S 5 7：否定）、フローテーブル 1 4 5 を参照し、選択したリンクを使用する最も優先度が低いフローを選択する（ステップ S 5 8）。セクタ 1 4 7 は、選択した最も優先度が低いフローのリンクに、対象フローの元のリンクを設定し（ステップ S 5 9）、切替処理を終了して経路制御処理に戻る。これにより、セクタ 1 4 7 は、優先度の高いフローからリンクに振り分けることができる。

40

【 0 0 5 8 】

図 7 の説明に戻って、モニタ 1 4 4 は、切替処理の終了後、フローテーブル 1 4 5 を参照し、全フローにおいて、アプリ性能または U L 性能が N G であれば、全フローの優先度をデクリメントする（ステップ S 6）。モニタ 1 4 4 は、フローテーブル 1 4 5 を参照し、全フローにおいて、アプリ性能および U L 性能が O K であれば、全フローの優先度をデ

50

クリメントし（ステップS7）、ステップS2に戻る。これにより、制御回路141は、システム全体の性能を最大化するようにフローを振り分けることができる。また、制御回路141は、ユーザロジックに依存することなく、システム全体の性能を最大化する通信パターンを自動的に決定できる。また、制御回路141は、ユーザロジックのFPGA接続方式に依存する部分を削減することができる。すなわち、情報処理装置100では、他のシステムへのユーザロジックの移植が容易となる。また、情報処理装置100では、インタフェース設計を共通化できるので、開発工数を削減できる。

【0059】

なお、上記実施例では、リンクの性能情報として使用帯域を用いたが、これに限定されない。例えば、リンクの性能情報としてレイテンシを用いてもよい。この場合、制御回路141は、Mux/demux142にレイテンシチェッカを接続し、定期的に計測用パケットを各リンクに送信して取得したレイテンシ情報に基づいて、各フローを各リンクに振り分ける。すなわち、モニタ144は、アプリケーションの処理に対応するフローごとのリンクのレイテンシを示す値と、リンクごとのレイテンシの最大値とを収集する。セクタ147は、フローごとのリンクのレイテンシを示す値に基づいて、要するレイテンシを満たしていないフローを判定する。セクタ147は、リンクごとのレイテンシの最大値と、リンクごとの現在のレイテンシの値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。

10

【0060】

このように、情報処理装置100は、複数のリンクを経由してCPUと接続されるオフロード回路を備える。オフロード回路は、FPGA140として、制御回路141と、ユーザロジックである論理回路とを備える。論理回路は、アプリケーションの処理を演算する。制御回路141は、アプリケーションの処理に対応するフローごとのリンクの性能情報（レイテンシ）を示す値と、リンクごとの使用可能な性能情報の最大値とを収集する。また、制御回路141は、フローごとのリンクの性能情報（レイテンシ）を示す値に基づいて、要する性能情報（レイテンシ）を満たしていないフローを判定する。また、制御回路141は、リンクごとの使用可能な性能情報の最大値と、リンクごとの現在使用されている性能情報の値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。その結果、制御回路141は、システム全体の性能を最大化するようにフローを振り分けることができる。

20

【0061】

また、制御回路141は、アプリケーションの処理に対応するフローごとのリンクの使用帯域を示す値と、リンクごとの使用可能な帯域の最大値とを収集する。また、制御回路141は、フローごとのリンクの使用帯域を示す値に基づいて、要する帯域を満たしていないフローを判定する。また、制御回路141は、リンクごとの使用可能な帯域の最大値と、リンクごとの現在使用されている帯域の値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。その結果、制御回路141は、リンクの使用帯域に基づいてシステム全体の性能を最大化するようにフローを振り分けることができる。

30

【0062】

また、制御回路141は、アプリケーションの処理、または、論理回路の性能が性能要件を満たしていない場合、対応するフローの優先度を上げる。その結果、制御回路141は、優先度に応じてフローを振り分けることができる。

40

【0063】

また、制御回路141は、リンクごとの使用可能な帯域の最大値と、リンクごとの現在使用されている帯域の値とに基づいて、フローの振り分け先のリンクを選択する。制御回路141は、選択した振り分け先のリンクに該フローを振り分けると、振り分け先のリンクの使用可能な帯域の最大値を超える場合、振り分け先のリンクを使用する最も優先度が低いフローを、振り分け元のリンクに振り分ける。その結果、制御回路141は、優先度の高いフローと優先度の低いフローとが使用するリンクを入れ替えることができる。

【0064】

また、制御回路141は、所定時間ごとに、リンクの使用帯域を示す値と、リンクごと

50

の使用可能な帯域の最大値とを収集する。また、制御回路141は、所定時間ごとに、要する帯域を満たしていないフローを判定する。その結果、制御回路141は、動的にフローを振り分けることができる。

【0065】

また、制御回路141は、アプリケーションの処理に対応するフローごとのリンクのレイテンシを示す値と、リンクごとのレイテンシの最大値とを収集する。また、制御回路141は、フローごとのリンクのレイテンシを示す値に基づいて、要するレイテンシを満たしていないフローを判定する。制御回路141は、リンクごとのレイテンシの最大値と、リンクごとの現在のレイテンシの値とに基づいて、フローの振り分け先のリンクを選択して振り分ける。その結果、制御回路141は、リンクのレイテンシに基づいてシステム全体の性能を最大化するようにフローを振り分けることができる。

10

【0066】

また、図示した各部の各構成要素は、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各部の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況等に応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。例えば、Mux/demux142とMux/demux143とを統合してクロスバースイッチとしてもよい。また、図示した各処理は、上記の順番に限定されるものでなく、処理内容を矛盾させない範囲において、同時に実施してもよく、順序を入れ替えて実施してもよい。

【0067】

さらに、制御回路141で行われる各種処理機能は、CPU（またはMPU、MCU（Micro Controller Unit）等のマイクロ・コンピュータ）上で、その全部または任意の一部を実行するようにしてもよい。また、各種処理機能は、CPU（またはMPU、MCU等のマイクロ・コンピュータ）で解析実行されるプログラム上、またはワイヤードロジックによるハードウェア上で、その全部または任意の一部を実行するようにしてもよいことは言うまでもない。

20

【0068】

なお、上記実施例で説明した制御回路141は、プログラムを読み込んで実行することで、図4、図7、図8等で説明した処理と同様の機能を実行することができる。例えば、制御回路141は、モニタ144、セクタ147と同様の処理を実行するプロセスを実行することで、上記実施例と同様の処理を実行することができる。

30

【0069】

これらのプログラムは、インターネットなどのネットワークを介して配布することができる。また、これらのプログラムは、ハードディスク、フレキシブルディスク（FD）、CD-ROM、MO、DVDなどのコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行することができる。

【符号の説明】

【0070】

- 100 情報処理装置
- 110 通信部
- 111 表示部
- 112 操作部
- 113 HDD
- 114 バス
- 120 メモリ
- 130 CPU
- 131 OS/VM
- 140 FPGA
- 141 制御回路
- 142, 143 Mux/demux

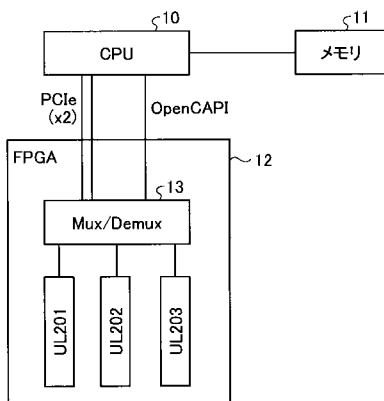
40

50

- 1 4 4 モニタ
- 1 4 5 フローテーブル
- 1 4 6 リンクテーブル
- 1 4 7 セレクタ
- A 1 ~ A m アプリケーション
- I F 1 ~ I F n インタフェース
- L 1 ~ L n リンク
- U L 1 ~ U L m ユーザロジック

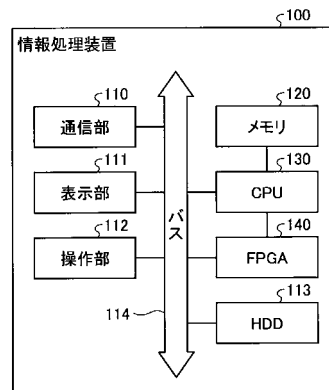
【 図 1 】

複数のリンクでCPUとFPGAを接続する場合の一例を示す図



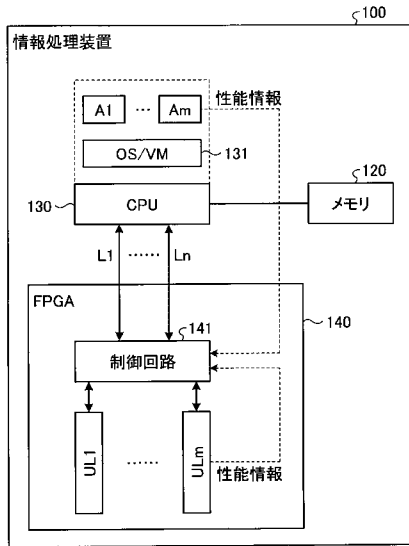
【 図 2 】

実施例の情報処理装置のハードウェア構成の一例を示すブロック図



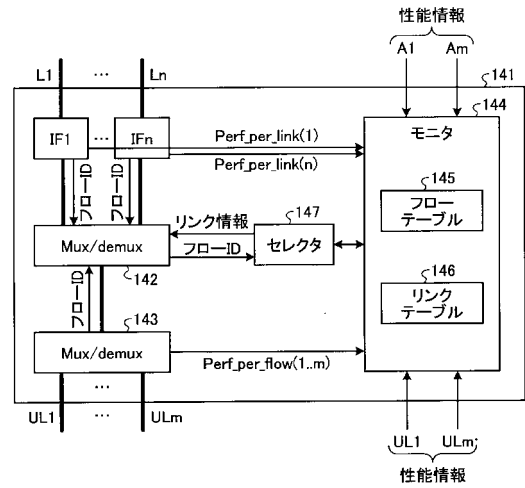
【 図 3 】

実施例の情報処理装置の機能構成の一例を示すブロック図



【 図 4 】

制御回路の機能構成の一例を示すブロック図



【 図 5 】

フローテーブルの一例を示す図

フローID	フロー情報			リンク情報	優先度 0:低~255:高	性能情報(現在値)		
	ユーザロジック	コマンド種別	アドレス範囲			アプリ性能	UL性能	使用帯域
1	UL1	Read	All	L1	10	OK	NG	10GB/s
2	UL1	Write	All	L1	1	OK	OK	1GB/s
3	UL2	R/W	All	L1	11	NG	OK	20GB/s
4	UL3	R/W	0x0000-0xFFFF	L2	12	OK	OK	30GB/s
5	UL3	R/W	0x10000-0x1FFF	L2	10	OK	OK	10GB/s
...

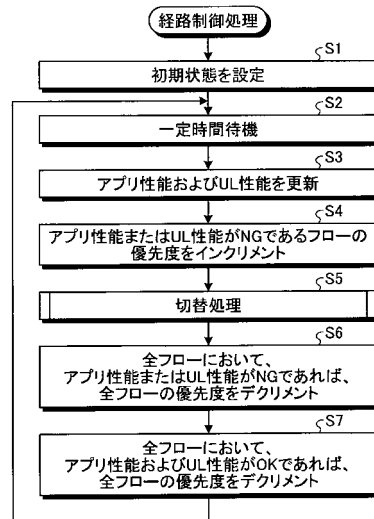
【 図 6 】

リンクテーブルの一例を示す図

リンクID	性能情報		
	最大帯域	最小遅延	...
L1	50GB/s	1000ns	...
L2	20GB/s	300ns	...
L3	10GB/s	100ns	...
...

【 図 7 】

実施例の経路制御処理の一例を示すフローチャート



【 図 8 】

