

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7502338号
(P7502338)

(45)発行日 令和6年6月18日(2024.6.18)

(24)登録日 令和6年6月10日(2024.6.10)

(51)国際特許分類 F I
 G 0 6 F 9/50 (2006.01) G 0 6 F 9/50 1 5 0 C
 G 0 6 F 9/455(2018.01) G 0 6 F 9/455 1 5 0

請求項の数 20 (全14頁)

(21)出願番号	特願2021-569031(P2021-569031)	(73)特許権者	314015767
(86)(22)出願日	令和2年4月16日(2020.4.16)		マイクロソフト テクノロジー ライセン
(65)公表番号	特表2022-533997(P2022-533997		シング,エルエルシー
	A)		アメリカ合衆国 ワシントン州 9 8 0 5
(43)公表日	令和4年7月27日(2022.7.27)		2 レッドモンド ワン マイクロソフト
(86)国際出願番号	PCT/US2020/028603		ウェイ
(87)国際公開番号	WO2020/236363	(74)代理人	100118902
(87)国際公開日	令和2年11月26日(2020.11.26)		弁理士 山本 修
審査請求日	令和5年4月14日(2023.4.14)	(74)代理人	100106208
(31)優先権主張番号	62/850,421		弁理士 宮前 徹
(32)優先日	令和1年5月20日(2019.5.20)	(74)代理人	100196508
(33)優先権主張国・地域又は機関	米国(US)		弁理士 松尾 淳一
(31)優先権主張番号	16/808,286	(72)発明者	チオウ, デレク
(32)優先日	令和2年3月3日(2020.3.3)		アメリカ合衆国 ワシントン州 9 8 0 5
	最終頁に続く		2 - 6 3 9 9 レッドモンド ワン マイ
			最終頁に続く

(54)【発明の名称】 S O CおよびF P G Aを有するサーバオフロードカード

(57)【特許請求の範囲】

【請求項1】

サーバであって、

C P U (中央処理装置) コンプレックスと、

オフロードカードであって、

S o C (システムオンチップ) と、

前記 S o C の外部にあり、且つ、前記 S o C に連結された F P G A (フィールドプログラマブルゲートアレイ) と

を含むオフロードカードと、

を備え、

前記 C P U コンプレックスが、1つまたは複数の仮想マシン (V M) を実行するように構成され、

前記 S o C が、前記1つまたは複数の V M に関連付けられたハイパーバイザの1つまたは複数の第1の機能をソフトウェアで実行するように構成され、

前記 F P G A が、前記1つまたは複数の V M に関連付けられた前記ハイパーバイザの1つまたは複数の第2の機能をハードウェアで実行するように構成される、

サーバ。

【請求項2】

請求項1に記載のサーバであって、前記 S o C および前記 F P G A が、前記オフロードカードの内部にある P C I e (ペリフェラルコンポーネントインターコネクトエクスプレス

) インターフェースを介して、および、前記オフロードカードの内部にあるイーサネットインターフェースを介して、互いに通信可能に連結される、サーバ。

【請求項 3】

請求項 2 に記載のサーバであって、前記 S o C および前記 F P G A が、前記オフロードカードの内部にある J T A G (ジョイントテストアクショングループ) インターフェースを介して、互いに通信可能に連結される、サーバ。

【請求項 4】

請求項 1 に記載のサーバであって、前記オフロードカードが、前記サーバのメインボードに P C I e エッジコネクタインターフェースを介して挿入される、サーバ。

【請求項 5】

請求項 4 に記載のサーバであって、前記 S o C が、前記サーバのベースボードマネジメントコントローラ (B M C) に、前記 P C I e エッジコネクタインターフェースを通じて通信可能に連結される、サーバ。

【請求項 6】

請求項 4 に記載のサーバであって、前記 F P G A が、前記 C P U コンプレックスと前記 P C I e エッジコネクタインターフェースを通じて通信可能に連結される、サーバ。

【請求項 7】

請求項 1 に記載のサーバであって、前記 S o C が、前記オフロードカード上にある 1 つまたは複数の揮発性メモリモジュールと通信可能に連結され、前記 1 つまたは複数の揮発性メモリモジュールが、前記 S o C が前記 1 つまたは複数の第 1 の機能をそこから実行できる作業メモリとして機能する、サーバ。

【請求項 8】

請求項 1 に記載のサーバであって、前記 S o C が、前記オフロードカード上にあるフラッシュメモリモジュールと通信可能に連結され、前記フラッシュメモリモジュールが、前記 1 つまたは複数の第 1 の機能のためのプログラムコードを格納する、サーバ。

【請求項 9】

請求項 1 に記載のサーバであって、前記 F P G A が、前記オフロードカード上にある 1 つまたは複数の揮発性メモリモジュールと通信可能に連結され、前記 1 つまたは複数の揮発性メモリモジュールが、前記 1 つまたは複数の第 2 の機能を実行する時に前記 F P G A のための作業メモリとして機能する、サーバ。

【請求項 10】

請求項 1 に記載のサーバであって、前記 F P G A が、前記オフロードカード上にあるフラッシュメモリモジュールと通信可能に連結され、前記フラッシュメモリモジュールが、1 つまたは複数の第 2 の機能を実行するように前記 F P G A を構成するための少なくとも 1 つの構成イメージを格納する、サーバ。

【請求項 11】

請求項 10 に記載のサーバであって、前記フラッシュメモリモジュールが、前記 F P G A のための正常動作構成に対応する第 1 の構成イメージ、および、前記 F P G A のためのフェールセーフ動作構成に対応する第 2 の構成イメージを格納する、サーバ。

【請求項 12】

請求項 11 に記載のサーバであって、前記第 1 の構成イメージが、前記オフロードカードの電源投入時にデフォルトで前記 F P G A に適用される、サーバ。

【請求項 13】

請求項 12 に記載のサーバであって、前記第 1 の構成イメージを適用している間にエラーが発生した場合、前記第 2 の構成イメージが、前記 F P G A に適用される、サーバ。

【請求項 14】

請求項 1 に記載のサーバであって、前記 F P G A が、T O R (トップオブラック) ネットワークスイッチと通信可能に連結された第 1 の外部ネットワークインターフェース、および、前記サーバの N I C (ネットワークインターフェースカード) と通信可能に連結された第 2 の外部ネットワークインターフェースを含む、サーバ。

10

20

30

40

50

【請求項 15】

請求項 1 に記載のサーバであって、前記 S o C が、前記オフロードカード上にある B I O S (ベーシックインプット/アウトプット)フラッシュ構成要素に、セキュリティチップを介して通信可能に連結され、前記セキュリティチップが、前記 B I O S フラッシュ構成要素に格納されたファームウェアの完全性を検証するように構成される、サーバ。

【請求項 16】

請求項 1 に記載のサーバであって、前記 1 つまたは複数の第 1 の機能が、ネットワーク制御プレーン機能またはストレージ制御プレーン機能を含む、サーバ。

【請求項 17】

請求項 1 に記載のサーバであって、前記 1 つまたは複数の第 2 の機能が、ネットワークデータプレーン機能またはストレージデータプレーン機能を含む、サーバ。

10

【請求項 18】

サーバであって、

1 つまたは複数の仮想マシン (V M) を実行するように構成された C P U (中央処理装置)コンプレックスと、

オフロードカードであって、

前記 1 つまたは複数の V M に関連付けられたハイパーバイザの 1 つまたは複数の第 1 の機能をソフトウェアで実行するための手段と、

前記 1 つまたは複数の V M に関連付けられた前記ハイパーバイザの 1 つまたは複数の第 2 の機能をハードウェアで実行するための手段と、

20

を含むオフロードカードと

を備えるサーバ。

【請求項 19】

方法であって、

サーバのオフロードカード上にある F P G A (フィールドプログラマブルゲートアレイ)が、前記サーバの N I C (ネットワークインターフェースカード)からネットワークパケットを受け取るステップであって、前記ネットワークパケットが、前記 F P G A と前記 N I C を相互接続するイーサネットインターフェースを介して受け取られる、ステップと、

前記 F P G A が、前記ネットワークパケットのヘッダに基づいてフローテーブルへのルックアップをハードウェアで実施するステップと、

30

マッチするエントリが前記ヘッダのための前記フローテーブル内で見つからなかったと判定すると、前記 F P G A が、前記オフロードカード上にある S o C (システムオンチップ)に前記ネットワークパケットを転送するステップであって、前記ネットワークパケットが、前記 F P G A と前記 S o C を相互接続するイーサネットインターフェースを介して転送される、ステップと、

前記 S o C が、前記ネットワークパケットのためのネクストホップの宛先をソフトウェアで計算するステップと、

前記 S o C が、前記ネクストホップの宛先を含む新しいフローエントリで前記フローテーブルをソフトウェアで更新するステップと

を含む方法。

40

【請求項 20】

請求項 19 に記載の方法であって、マッチするエントリが前記フローテーブル内で見つかったと判定すると、

前記 F P G A が、前記マッチするエントリに基づいて前記ネットワークパケットを更新するステップと、

前記 F P G A が、前記 F P G A の外部ネットワークインターフェースを介して外部ネットワークに前記ネットワークパケットを伝送するステップと、

をさらに含む方法。

【発明の詳細な説明】

【背景技術】

50

【0001】

[0001]Microsoft AzureおよびAmazon AWSなどのクラウドプラットフォームは、地理的に拡散されたデータセンタの全体にわたって分散された物理サーバの大きい集団（本明細書ではクラウドサーバと呼ばれる）の上で動く。これらのクラウドサーバのかなりの部分が、仮想マシン（VM）のホスティングを可能にするハイパーバイザとして知られる仮想化ソフトウェアレイヤを実装する。数ある中でも、仮想化ソフトウェアレイヤは、クラウドプラットフォームの顧客が購入し、VMを使用して顧客のアプリケーションワークロードを実行するIaaS（サービスとしてのインフラストラクチャ）シナリオを可能にする。

【発明の概要】

10

【発明が解決しようとする課題】

【0002】

[0002]伝統的に、ハイパーバイザを実装する各クラウドサーバでは、クラウドサーバのCPU（中央処理装置）コアの一定の割合が、ハイパーバイザ用に確保される。この確保は、ハイパーバイザがその機能を実行するのに十分な計算リソースを有することを保証するが、例えば顧客VMによる使用のために利用可能なCPUコアの数も低減させる。規模の点で、これは、顧客とじかに接するクラウドプラットフォームの全般的な計算容量の大幅な低下になる恐れがある。

【課題を解決するための手段】

【0003】

20

[0003]SoC（システムオンチップ（system-on-chip））およびFPGA（フィールドプログラマブルゲートアレイ（field programmable gate array））を含むオフロードシステムを有する物理サーバが開示される。オフロードシステムの1つの可能な実施形態はカード上にある。実施形態の1つのセットによれば、SoCは、ソフトウェアでの実行に適したサーバのCPUコンプレックスから1つまたは複数のハイパーバイザ機能をオフロードするように構成することができ、FPGAは、ハードウェアでの実行に適したCPUコンプレックスから1つまたは複数のハイパーバイザ機能をオフロードするように構成することができる。

【図面の簡単な説明】

【0004】

30

【図1】[0004]いくつかの実施形態による、SoCおよびFPGAを有するオフロードカードを含む物理サーバトポロジを描写する図である。

【図2】[0005]いくつかの実施形態による、図1のオフロードカードのためのアーキテクチャを描写する図である。

【図3】[0006]いくつかの実施形態による、JTAG（ジョイントテストアクショングループ（Joint Test Action Group））多重化装置の実装形態を描写する図である。

【図4】[0007]いくつかの実施形態による、实例のネットワーク処理フロー図である。

【発明を実施するための形態】

【0005】

40

[0008]以下の説明では、説明のために、非常に多くの例および詳細が、様々な実施形態を理解するために示されている。それでも、いくつかの実施形態を、これらの詳細の一部がなくても実践できること、または、変更形態もしくはその同等物で実践できることが当業者には明らかであろう。

【0006】

1. 全体像

[0009]本開示の実施形態は、SoC（システムオンチップ）およびFPGA（フィールドプログラマブルゲートアレイ）を備えるオフロードカードを用いる物理サーバ設計を対象とする。様々な実施形態では、SoCおよびFPGAは、サーバのCPUコンプレックスによって伝統的に実行されるハイパーバイザ機能を動かすことができ、これにより、こ

50

これらの機能の処理負担をCPUコンプレックスからオフロードする。例えば、オフロードカードのSoCは、汎用プロセッサの柔軟性（例えばネットワーキングおよびストレージ制御プレーン機能）を必要とするか、汎用プロセッサの柔軟性から利益を得るハイパーバイザ機能を動かすことができるが、オフロードカードのFPGAは、ハードウェア（例えばネットワーキングおよびストレージデータプレーン機能）での実装/アクセラレーションに適したハイパーバイザ機能を動かすことができる。

【0007】

[0010]この全体的なアーキテクチャにより、全てではないがほとんどのハイパーバイザ処理をサーバのCPUコンプレックスからオフロードカードに移すことができ、オフロードカードは、都合のよいことに、CPUコンプレックスがテナント（例えば顧客）VMワークロードを実行することに集中することを可能にする。ハイパーバイザがCPUコンプレックスから完全に立ち退かされるケースでは、テナントコードは、CPUコンプレックス上で「ベアメタル」のように（すなわち介在するいかなるハイパーバイザ仮想化レイヤもない状態で）動かせる可能性がある。

【0008】

[0011]なおさらに、オフロードカード上でのハイパーバイザコード/ロジックの実行は、CPUコンプレックス上のテナントコードの実行から物理的に切り離されるので、このソリューションは、テナントコードを攻撃ベクトルとして使用しようとし得るサイドチャネル攻撃からハイパーバイザを保護する。

【0009】

[0012]さらに、ハードウェア実装に対応可能な一定のハイパーバイザ機能を加速させるためにFPGAを用いることによって、オフロードカードは、構造的柔軟性を同時に維持しつつ、サーバの効率性を改善することができる。例えば、必要ならFPGAは、1つのタイプ/クラスの機能（例えばネットワーキング）の加速から、別のタイプ/クラスの機能（例えばストレージ）の加速にプログラムし直すことができる。これは、ASIC（特定用途向け集積回路）などのハードロジックベースのアクセラレータでは可能ではない。

【0010】

[0013]本開示の前述および他の態様が、以下のセクションでさらに詳細に説明される。

2. サーバトポロジ

[0014]図1は、本開示の一定の実施形態による、物理サーバ100の高レベルトポロジを示す簡易ブロック図である。実施形態の1つのセットでは、物理サーバ100は、クラウドプラットフォームのインフラストラクチャの一部として導入されたクラウドサーバでよい。これらの実施形態では、物理サーバ100は、クラウドプラットフォームプロバイダによって運用されるデータセンタ内のサーバラックにマウントされてよい。他の実施形態では、物理サーバ100は、例えばスタンドアロンサーバの形で、敷地内の企業IT環境など、他の状況でおよび/または他のフォームファクタを介して、導入されてよい。

【0011】

[0015]背景技術セクションで述べたように、クラウドサーバは仮想化のためにハイパーバイザを実装することが多く、これによりクラウドプラットフォームが、IaaS（サービスとしてのインフラストラクチャ）などのサービスを提供することができる。それでも、ハイパーバイザ（「ホスト」としても知られる）用に、CPUコアを含むこれらのプラットフォームリソースの一部を使用することにより、従来のクラウドサーバは、これらのCPU容量の全てをVMに露出することができず、これによりプラットフォームの効率性を低下させる。

【0012】

[0016]この問題および他の問題に対処するために、物理サーバ100は、SoC104およびFPGA106を備える斬新なオフロードカード102を含む。図示の実施形態では、オフロードカード102は、PCIe（ペリフェラルコンポーネントインターフェースエクスプレス（Peripheral Component Interface Express））ベースの拡張カードとして実装され、したがって、標準PCIe x16

10

20

30

40

50

3.0 エッジコネクタインターフェース108を介して物理サーバ100のメインボードとインターフェースする。他の実施形態では、オフロードカード102は、周辺インターフェースの他のいずれかのタイプを使用して実装されてもよい。

【0013】

[0017]図示のように、SoC104は、独自のRAM(ランダムアクセスメモリ(random access memory))110およびフラッシュメモリ112を有し、オフロードカード102の内部にある少なくともインターフェース(すなわちPCIeインターフェース114およびイーサネットインターフェース116)を介してFPGA106と通信可能に連結される。さらにSoC104は、I2Cインターフェース108およびいくつかの他のチャネル(例えばUSBおよびCOM)を通じて、物理サーバ100のベースボードマネージメントコントローラ(BMC: base board management controller)118と通信可能に連結される。

10

【0014】

[0018]FPGA106も、独自のRAM120およびフラッシュメモリ122を有し、PCIeエッジコネクタインターフェース108を通じて物理サーバ100のCPUコンプレックス124と通信可能に連結される。このCPUコンプレックスは、物理サーバ100のメインCPUコアおよび関連付けられたRAMモジュールを備える。さらにFPGA106は、2つの外部イーサネットインターフェースを含み、これらの一方は、(例えばTOR(トップオブラック)スイッチまたは他のいくつかのネットワークデバイスを介して)外部ネットワーク126に接続し、他方は、物理サーバ100内のNIC(ネットワークインターフェースカード/コントローラ(network interface card/controller))128に接続する。

20

【0015】

[0019]一般的に言えば、図1に示されたトポロジは、物理サーバ100のCPUコンプレックス124上で伝統的に動かされるハイパーバイザ機能の一部または全てをオフロードカード102のSoC104およびFPGA106上で代わりに動かし、したがってオフロードカード102のSoC104およびFPGA106にオフロードすることを可能にする。例えば、汎用プロセッサの柔軟性から利益を得る(またはハードウェアで実装するには単に複雑/強すぎる)ハイパーバイザ機能をSoC104上で動かすことができ、SoC104には1つまたは複数の汎用処理コアが組み込まれている。このような機能の例は、SDN(ソフトウェア定義ネットワークング(software-defined networking))制御プレーン機能を含み、SDNは複雑なルーティング計算を必要とし、新しいプロトコルおよび特徴をサポートするために比較的頻繁に更新される必要がある。

30

【0016】

[0020]その一方で、ハードウェアアクセラレーションに適したハイパーバイザ機能は、FPGA106上の論理ブロックを介して実装することができる。このような機能の例は、SDNデータプレーン機能および(データ複製、重複排除などの)ストレージデータプレーン機能を含み、SDNデータプレーン機能は、制御プレーンの決定に従ってネットワークデータトラフィックを転送することを伴う。

40

【0017】

[0021]このソリューションにより、従来のサーバ設計を超えるいくつかの長所が達成される。第1に、一定のホスト処理の役目をCPUコンプレックス124から取り除くことによって、ハイパーバイザによって使用される(CPUコンプレックス124内のCPUコアを含む)プラットフォームリソースの量を減少させることができ、(「ゲスト」としても知られる)VMが利用可能なプラットフォーム容量を増加させる。これは、サーバ能力のあらゆる漸進的向上が、規模の点で重大なインパクトを与える恐れがあるパブリッククラウドプラットフォームにおいて特に有益である。いくつかの実施形態では、ハイパーバイザは、CPUコンプレックス124から全面的に立ち退かされ、オフロードカード102に移されてもよく、このケースでは、CPUコンプレックス124は最低限のハイパ

50

ーバイザを動かすことができるか（CPUコンプレックス124は一定のレジスタにアクセスすることなど、CPUコンプレックス自体でしか動かせない問題を処理する）、またはハイパーバイザが全くなく、CPUコンプレックス124の計算能力の残りをゲストワークロードに充てることできる。

【0018】

[0022]第2に、SoC104（非ハードウェアアクセラレーション機能を扱う）と、FPGA106（ハードウェアアクセラレーション機能を扱う）との両方をオフロードカード102上に実装し、これら2つをしっかりと連結することによって、SoC104上で動くハイパーバイザコードがFPGA106に実装されたロジックと相互作用しやすくなり、逆もまた同様である。オフロードカード102上にハードウェアアクセラレータを単

10

【0019】

[0023]第3に、オフロードカード102上で動くホストコードは、CPUコンプレックス124上で動くゲストコードから物理的に切り離されるので、悪意のエンティティがVMを介してハイパーバイザを攻撃するのが困難になる。これは、現代のCPUアーキテクチャにおける一定のサイドチャネル脆弱性の最近の発見の観点から特に適切である。これらの既知の脆弱性にはパッチを当てることができるが、他の同様の脆弱性が将来見つかる

20

【0020】

[0024]第4に、ハードウェアアクセラレーションのためにASICではなくFPGAを使用することによって、FPGAをプログラムし直すことによって改善される種々のユースケースまたは同じユースケースのためにオフロードカード102を簡単に再利用することができ、必要ならFPGAのロジックを更新することができる。これは、既に現場にある多くのカードを引き抜き、置き替えることが望ましくないことがある大規模な導入に有利である。

【0021】

[0025]図1の物理サーバ100について示された特定のトポロジは例証であり、様々な変更形態が可能であることを理解されたい。例えば、SoC104およびFPGA106は、周辺機器（例えばPCIe）インターフェースを介して物理サーバのメインボードとインターフェースする拡張カード（すなわちオフロードカード102）上に実装されるものとして示されているが、いくつかの実施形態では、代替のオフロードアーキテクチャが使用されてもよい。特定の实施形態では、SoC104および/またはFPGA106の1つまたは複数は、サーバのメインボードに直接実装されてもよい。

30

【0022】

[0026]別の例として、NIC128がスタンドアロンの構成要素として描写されているが、いくつかの実施形態では、NIC128の機能は、FPGA106内など、図1に示された1つまたは複数の他の構成要素に組み込まれてもよい。当業者は、他の変形形態、変更形態、および代替形態を認識するはずである。

40

3. オフロードカードアーキテクチャ

[0027]図2は、いくつかの実施形態による図1のオフロードカード102のアーキテクチャに関するさらなる詳細を提示する概略図200である。次に本アーキテクチャの様々な態様が下記で論じられる。

3.1 SoC

[0028]SoC104は、1つまたは複数の汎用処理コア、（メモリ、ストレージ、および周辺機器のための）インターフェース、ならびにNICを含む、いくつかの既存のシステムオンチップ設計のいずれか1つを使用して実装することができる。特定の实施形態では、SoC104は、ARMマイクロプロセッサアーキテクチャに基づく汎用処理コアを

50

組み込むことができる。}

[0029]図示のように、SoC 104は、3つの別個のインターフェースを介してFPGA 106と通信可能に連結され、このインターフェースは、下記のセクション3.2で論じられる。さらにSoCは、(1)メモリアンターフェース204を介して図1のRAM 110に対応する1つまたは複数のDRAM(ダイナミックRAM)モジュール202に、(2)ストレージインターフェース208を介して図1のフラッシュメモリ112に対応するeMMC(組込型マルチメディアカード: embedded multimedia card)デバイス206に、(3)SPI(シリアルペリフェラルインターフェース: Serial Peripheral Interface)インターフェース212および介在するセキュリティチップ214を介してBIOSフラッシュメモリ構成要素210に、ならびに、(4)(FPGA 106およびPCIeエッジコネクタインターフェース108にも接続する)I2Cバス222を介して(EEPROM 216、ホットスワップコントローラ218、および温度センサ220などの)いくつかのI2C(集積デバイス間)デバイスに、接続される。

10

【0023】

[0030](1)に関して、SoC 104は、物理サーバ100のCPUコンプレックス124からオフロードされたハイパーバイザコードを含むプログラムコードを動かすためのSoC 104の作業メモリとして、DRAMモジュール202を使用することができる。DRAMモジュール202の固有の数および容量、ならびにメモリアンターフェース204の仕様は、実装形態に応じて変化させることができる。特定の形態では、DRAMモジュール202は、シングル1024M(メガビット)×64ビット+ECC(エラー訂正コード)メモリバンクとして編成された8GB(ギガバイト)のDDR4 DRAMを備えることができ、メモリアンターフェース204は、シングルDDR4-2400メモリチャネルとして構成することができる。

20

【0024】

[0031](2)に関して、SoC 104は、CPUコンプレックス124からオフロードされたハイパーバイザコードを含む、SoC上で実行されることになるプログラムコードを格納およびブートするため、ならびに、FPGA 106に適用されることになるFPGA構成イメージを格納するための非一時的ストレージ媒体として、eMMCデバイス204を使用することができる。

30

【0025】

[0032](3)に関して、BIOSフラッシュメモリ構成要素210は、SoC 104のためのシステムファームウェアを保持することができ、セキュリティチップ214は、数ある中でも、このシステムファームウェアが、攻撃者によって意図的または偶然に修正されることも破損させられることもないことを保証することができる。

【0026】

[0033](4)に関して、I2Cデバイス216、218、および220は、オフロードカード102に関する様々な管理情報をBMC 118に提供することができる。これらの情報は、動作温度データ、製造情報、および電力消費量データなどの情報を含むことができる。

40

【0027】

[0034]上記に加えて、SoC 104は、外部ヘッダ228、230、および232へのUSB(ユニバーサルシリアルバス)、COM、およびJTAG(ジョイントテストアクセシビリティグループ)インターフェース224、225、および226をそれぞれ含み、これらは、デバッグまたは管理のためにSoC 104をBMC 118または外部デバイスと接続するために使用することができる。BMC 118によってPCIeエッジコネクタインターフェース108を通じてSoC 104に送ることができるパワースロットル信号234もある。

3.2 SoCとFPGAとの間のインターフェース

[0035]既に前で述べたように、SoC 104は、図2の3つの内部のチップ間(chi

50

p - t o - c h i p) インターフェース (P C I e インターフェース 2 3 6、イーサネットインターフェース 2 3 8、および J T A G インターフェース 2 4 0) を介して F P G A 1 0 6 と通信可能に連結される。様々な実施形態では、P C I e インターフェース 2 3 6 は、制御能力とデータ転送 / 交換能力の両方を提供する。制御能力については、S o C 1 0 4 は、P C I e インターフェース 2 3 6 (または代替として J T A G インターフェース) を使用して、F P G A 1 0 6 を管理および更新することができる。例えば、S o C 1 0 4 は、R A M 1 1 0 から F P G A 1 0 6 に転送された F P G A 構成イメージが正しいことを確認することができ、このインターフェースを使用して F P G A 上または F P G A のフラッシュメモリ 1 2 2 内のイメージを更新することができる。データ能力については、P C I e インターフェース 2 3 6 は、S o C 1 0 4 上で動くプログラムコードが F P G A 1 0 4 にデータを送ること、および F P G A 1 0 4 からデータを受け取ることができる。これは、例えば P C I e を介してデータを交換するように既に書かれているハイパーバイザコードに有用であり、なぜなら、S o C 1 0 4 上での実行 (または F P G A 1 0 6 での実装) のために、このようなコードを比較的少ない変更で移植することができるからである。特定の実施形態では、P C I e インターフェース 2 3 6 は、8 つの P C I 3 . 0 レーンを有することができる (すなわち P C I 3 . 0 8 x インターフェースに対応させることができる)。他の実施形態では、4、1 2、1 6 など、他の任意の数の P C I レーンがサポートされ得る。

10

【 0 0 2 8 】

[0 0 3 6] イーサネットインターフェース 2 3 8 は、ネットワークパケットの形のデータを S o C 1 0 4 および F P G A 1 0 6 が交換することを可能にする。これは、例えばネットワークパケットを介してデータを交換するように既に書かれているハイパーバイザコードに有用であり、なぜなら、S o C 1 0 4 上での実行 (または F P G A 1 0 6 上での実装) のために、このようなコードを比較的少ない変更で移植することができるからである。例えば、ネットワークフローベースの転送が F P G A 1 0 6 上のハードウェアで実装され、ネットワークフローのためのルートを決定するためのネットワーク制御プレーンが S o C 1 0 4 上のソフトウェアで実装されるシナリオを考察する。このケースでは、フローテーブルの例外およびルールをネットワークパケットの形で、F P G A 1 0 6 と S o C 1 0 4 との間で通信することができる。特定の実施形態では、イーサネットインターフェース 2 3 8 は、2 5 G (ギガビット) イーサネットをサポートすることができる。

20

30

【 0 0 2 9 】

[0 0 3 7] J T A G インターフェース 2 4 0 は、S o C 1 0 4 が、低レベルテスト (例えばデバッグ) およびプログラミングのために F P G A 1 0 6 と通信するための通路を提供する。いくつかの実施形態では、外部ヘッダ 2 3 2 を介して接続された外部プログラマデバイスがインターフェース 2 4 0 を駆動させることを可能にする J T A G 多重化装置を、S o C 1 0 4 と F P G A 1 0 6 との間の J T A G 経路に挿入することができる。これらの実施形態では、外部プログラマデバイスからの「現在」の信号が、J T A G インターフェース 2 4 0 の信号経路を S o C 1 0 4 からデバイスにスイッチする。これは、最初のビットストリームをロードするときの最初のオフロードカードの立ち上げ (b r i n g - u p) に有用であり、S o C から F P G A への J T A G 経路が用意されていないときの F P G A アプリケーション開発に有用である。図 3 は、いくつかの実施形態による J T A G 多重化装置 3 0 2 を含む本アーキテクチャの実例の図 3 0 0 を描写する。

40

3 . 3 F P G A

[0 0 3 8] F P G A 1 0 6 は、いくつかの既存の F P G A チップのいずれか 1 つを使用して実装することができる。特定の実施形態では、F P G A 1 0 6 は、一定の最低限の数のプログラム可能論理素子 (例えば 1 0 0 0 K 個の素子) および一定の最低限のトランシーバ / F P G A ファブリックスピードグレード (例えばグレード 2) をサポートする既存の F P G A チップを使用して実装することができる。図 2 に示されているように、F P G A 1 0 6 は、I 2 C バス 2 2 2 と通信可能に連結され、上で論じたインターフェース 2 3 6 ~ 2 4 0 を介して S o C 1 0 4 と通信可能に連結される。さらに F P G A 1 0 6 は、(1)

50

内部 P C I e インターフェース 2 4 2 を介して P C I e エッジコネクタインターフェース 1 0 8 に、(2) メモリインターフェース 2 4 6 を介して図 1 の R A M 1 2 0 に対応する 1 つまたは複数の D R A M モジュール 2 4 4 に、(3) ストレージインターフェース 2 4 9 を介して図 1 のフラッシュメモリ 1 2 2 に対応する Q S P I (クアッドシリアルペリフェラルインターフェース (Q u a d S e r i a l P e r i p h e r a l I n t e r f a c e)) フラッシュメモリモジュール 2 4 8 に、ならびに(4) イーサネットインターフェース 2 5 4 および 2 5 6 を介して 2 つのネットワークランシーバモジュール 2 5 0 および 2 5 2 に、それぞれ接続される。

【 0 0 3 0 】

[0 0 3 9] (1) に関して、内部 P C I e インターフェース 2 4 2 は、C P U コンプレックス 1 2 4、および(例えば N I C 1 2 8 を含む)物理サーバ 1 0 0 にインストールされた他の P C I e デバイスと、F P G A 1 0 6 が通信することを可能にする。特定の実施形態では、P C I e インターフェース 2 4 2 は、P C I e 3 . 0 x 1 6 インターフェースでよい。

10

【 0 0 3 1 】

[0 0 4 0] (2) に関して、F P G A 1 0 6 は、C P U コンプレックス 1 2 4 からオフロードされたハイパーバイザロジックを含む、デバイスにプログラムされたロジックを実行するときの、F P G A 1 0 6 の作業メモリとして D R A M モジュール 2 4 4 を使用することができる。D R A M モジュール 2 4 4 の固有の数および容量、ならびにメモリインターフェース 2 4 6 の仕様は、実装形態に応じて変化させることができる。特定の実施形態では、D R A M モジュール 2 0 2 は、2 つの 4 G B バンクの 5 1 2 M x 6 4 ビット + E C C として編成された 8 G B (ギガバイト)の D D R 4 D R A M を備えることができ、メモリインターフェース 2 4 6 は、デュアル D D R 4 - 2 4 0 0 メモリチャネルとして構成することができる。

20

【 0 0 3 2 】

[0 0 4 1] (3) に関して、Q S P I フラッシュメモリモジュール 2 4 8 は、F P G A 1 0 6 の指定の機能を実施するように F P G A 1 0 6 自体を構成するために、電源投入時に F P G A 1 0 6 がロードできる 1 つまたは複数の F P G A 構成イメージを保持することができる。いくつかの実施形態では、Q S P I フラッシュメモリモジュール 2 4 8 は、下記のセクション 3 . 4 に記述された少なくとも 3 つの別個のイメージを保持することができる。フラッシュメモリからの構成の他に、F P G A 1 0 6 は、外部 J T A G プログラムデバイスを介した構成、J T A G インターフェース 2 4 0 を介して S o C 1 0 4 によって送られた J T A G コマンド、P C I e を介した C v P (プロトコルを介した構成 (C o n f i g u r a t i o n v i a P r o t o c o l))、および P C I e を介した部分再構成もサポートすることができる。

30

【 0 0 3 3 】

[0 0 4 2] (4) に関して、ネットワークランシーバモジュール 2 5 0 は、F P G A 1 0 6 が、外部ネットワーク 1 2 6 から入って来るネットワークトラフィックを受け取ること、および、外部ネットワーク 1 2 6 に出るネットワークトラフィックを伝送することを可能にする。さらに、ネットワークランシーバモジュール 2 5 2 は、F P G A 1 0 6 がネットワークトラフィックを N I C 1 2 8 と交換することを可能にする。これは、F P G A 1 0 6 がモジュール 2 5 2 を介して N I C 1 2 8 から出て行くネットワークパケットを受け取り、ネットワークパケットを適切に処理 / 変換し、モジュール 2 5 0 を介して外部ネットワーク 1 2 6 にネットワークパケットを送出することができるので、F P G A 1 0 6 がネットワークプレーン機能を実装するシナリオに有用である。逆に、F P G A 1 0 6 は、モジュール 2 5 0 を介して外部ネットワーク 1 2 6 から入って来るネットワークパケットを受け取り、ネットワークパケットを適切に処理 / 変換し、モジュール 2 5 2 を介して N I C 1 2 8 にネットワークパケットを送ることができる(この時点でネットワークパケットを正しい宛先 V M に通信することができる)。このようにネットワークデータプレーンアクセラレーションのために F P G A 1 0 6 を活用する実例のネットワークデータ

40

50

フローが、下記のセクション4で論じられる。特定の実施形態では、ネットワークトランシーバモジュール250および252はQSPF28光モジュールであることが可能であり、イーサネットインターフェース254および256は100Gイーサネットをサポートすることができる。

3.4 FPGAフラッシュ構成の詳細

[0043]実施形態の1つのセットでは、QSPIフラッシュメモリモジュール248は、FPGA106のための最低限の3つの別個の構成イメージ（ゴールデンイメージ、フェールセーフイメージ、およびユーザアプリケーションイメージ）を格納することができる。ゴールデンイメージは、最初の製造時に工場試験され、FPGA106のための正常な意図する機能を備える。フェールセーフイメージは工場プログラムされ、製造後は決して上書きされない。様々な実施形態では、このフェールセーフイメージは、電源投入時にオフロードカード102によって要求される機能の最小限のセットを収め、FPGA106のネットワークインターフェースは、FPGAによるいかなる中間処理もないインターフェース間を全てのトラフィックが直接通過する迂回モードを強制される。最後に、ユーザアプリケーションイメージは、ユーザ/顧客によって定義されたイメージである。

【0034】

[0044]オフロードカード102が電源投入された時、デフォルトでゴールデンイメージは、QSPIフラッシュメモリモジュール248からロードされ、FPGA106の構造を構成するためにFPGA106に適用されることになる。この電源投入処理に伴う何らかのエラーがある場合（またはサーバ実行時間中に問題が見つかった場合）、ゴールデンイメージの代わりにフェールセーフイメージをロードするためにカードをリブートすることができる。

4 実例のネットワーク処理ワークフロー

[0045]前述のオフロードカードアーキテクチャを考慮に入れて、図4は、いくつかの実施形態による物理サーバ100によって実装され得る実例のネットワーク処理ワークフローのフローチャート400を描写する。フローチャート400は、オフロードカード106のFPGA106が、SOC104上で動くネットワーク制御プレーンによって決定されたネットワークフローを備えるフローテーブルを維持すること、およびフローテーブルに従ってデータパケットを転送することを行うように構成されると仮定する。

【0035】

[0046]ブロック402でスタートし、物理サーバ100のNIC128は、SR-IOV（シングルルートIO仮想化）インターフェースを、サーバ100上で動くVMに提示することができる。このSR-IOVインターフェース（仮想機能と呼ばれる）は、ハイパーバイザを伴わずにVMがNIC128と直接通信することを可能にする。

【0036】

[0047]ブロック404において、VMは、リモートの宛先に伝送されることになるネットワークパケットのためのデータペイロードを作り出すことができ、これをNIC128に通知することができる。これに回答して、NIC128は、VMのゲストメモリ空間からデータペイロードを読み込むこと（ブロック406）、数ある中でも、VMのIPアドレス、および意図される宛先のIPアドレスを識別するヘッダを伴う1つまたは複数のネットワークパケットにデータペイロードを組み立てること（ブロック408）、ならびに、FPGA106のネットワークトランシーバモジュール252に接続されたNIC128の出口ポートからネットワークパケットを出力すること（ブロック410）を行うことができる。

【0037】

[0048]ブロック412および414において、FPGA106はネットワークパケットを受け取り、フローテーブルへのネットワークパケットの5-タプル（ソースIPアドレス、ソースポート、宛先IPアドレス、宛先ポート、プロトコル）のルックアップを実施するために、FPGA106のネットワークデータプレーンロジックを適用することができる。マッチするエントリがテーブル内で見つかった場合（ブロック416）、FPGA

106は、エントリ内でネットワークパケットのためのネクストホップの宛先を識別すること(ブロック418)、パケットのヘッダを更新すること(ブロック420)、ネットワークランシバモジュール250から外部ネットワーク126にパケットを送ること(ブロック422)を行うことができ、これによりワークフローを終える。

【0038】

[0049]その一方で、ブロック416において、マッチするエントリがテーブル内で見つからなかった場合(これがフローにおける第1のパケットであることを示す)、FPGA106は、内部イーサネットインターフェース238を介してネットワークパケットをSOC104に送ることができる(ブロック424)。SOC104上で動くネットワーク制御プレーン構成要素は、次に、パケットのためのネクストホップの宛先を計算し、インターフェース238を介してパケットのネットワークフローのための新しいエントリをFPGAのフローテーブルに追加することができる(ブロック426)。この新しいエントリを使用して、FPGA106はブロック420および422を実行することができ、ワークフローを終えることができる。

10

【0039】

[0050]上記の説明は、本開示の様々な実施形態を、これらの実施形態の態様がどのように実装され得るかについての例と共に示す。上記の例および実施形態は唯一の実施形態であるとは見なされるべきではなく、以下の特許請求の範囲によって定義されるような本開示の柔軟性および長所を示すために提示される。例えば、いくつかの実施形態が特定の処理フローおよびステップについて説明されてきたが、本開示の範囲が、説明されたフローおよびステップに厳密に限定されるわけではないことが当業者には明らかなはずである。連続的なものと説明されたステップは同時に実行されてもよく、ステップの順序は変化してもよく、ステップは修正、結合、追加、または省略されてもよい。別の例として、いくつかの実施形態がハードウェアとソフトウェアの特定の組合せを使用して説明されてきたが、ハードウェアとソフトウェアの他の組合せが可能であること、および、ソフトウェアで実装されるものとして説明された特定の動作をハードウェアでも実装することができ、逆もまた同様であることを認識するはずである。

20

【0040】

[0051]したがって、本明細書および図面は制限的な意味ではなく、例証的なものとみなされるべきである。他の配置、実施形態、実装形態、および同等物が当業者には明らかなはずであり、以下の特許請求の範囲において示されるような本開示の趣旨および範囲から逸脱することなく採用されてもよい。

30

40

50

【 図面 】

【 図 1 】

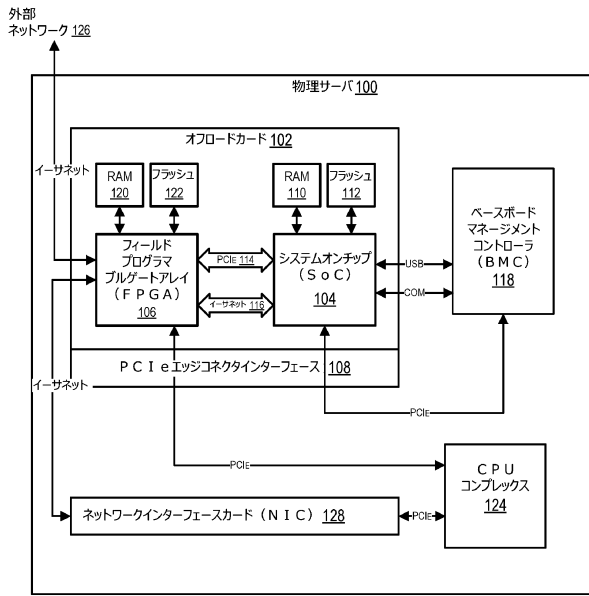


FIG. 1

【 図 2 】

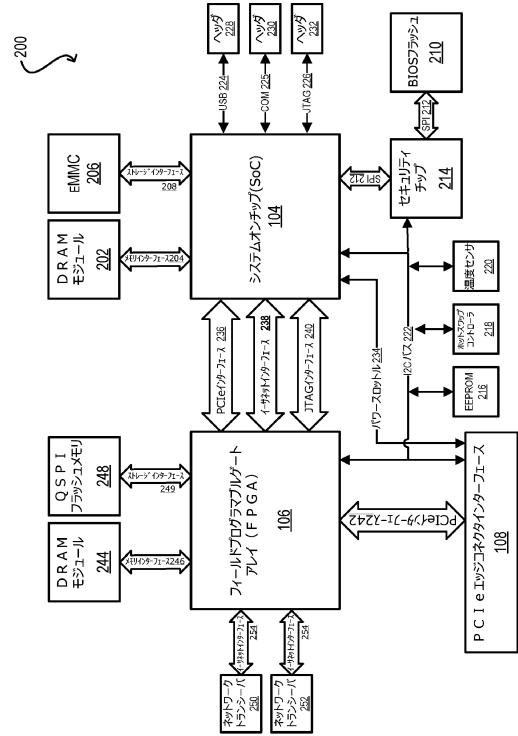


FIG. 2

【 図 3 】

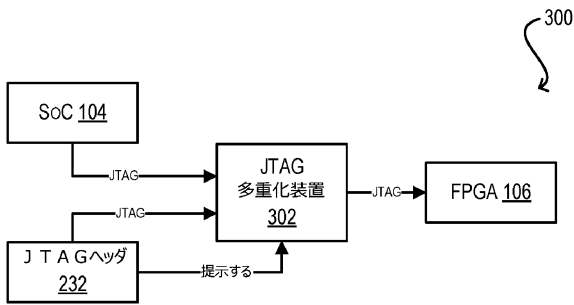


FIG. 3

【 図 4 】

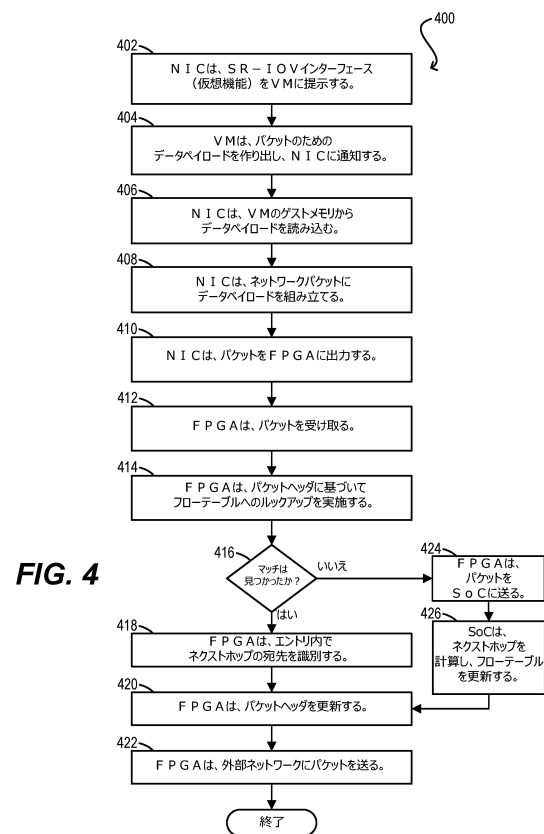


FIG. 4

10

20

30

40

50

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

マイクロソフト ウェイ, マイクロソフト テクノロジー ライセンシング, エルエルシー

(72)発明者 パトナム, アンドリュー

アメリカ合衆国 ワシントン州 98052-6399 レッドモンド ワン マイクロソフト ウェイ, マイクロソフト テクノロジー ライセンシング, エルエルシー

(72)発明者 ファイアストーン, ダニエル

アメリカ合衆国 ワシントン州 98052-6399 レッドモンド ワン マイクロソフト ウェイ, マイクロソフト テクノロジー ライセンシング, エルエルシー

(72)発明者 ラビアー, ジャック

アメリカ合衆国 ワシントン州 98052-6399 レッドモンド ワン マイクロソフト ウェイ, マイクロソフト テクノロジー ライセンシング, エルエルシー

審査官 坂東 博司

(56)参考文献

国際公開第2019/083977(WO, A1)

特表2021-501407(JP, A)

国際公開第2018/064415(WO, A1)

特表2019-535092(JP, A)

特開2012-156796(JP, A)

特表2015-515798(JP, A)

特開2018-149755(JP, A)

特開2017-174301(JP, A)

米国特許出願公開第2018/0004539(US, A1)

DANIEL FIRESTONE et al, Azure Accelerated Networking: SmartNICs in Public Cloud, Proceedings of the 15th USENIX Symposium on Networked Systems Design and Implementation, USENIX Association, 2018年04月09日, 51-64, [online], 発行日2018年4月9日, [令和5年10月17日検索], インターネット <URL: https://www.usenix.org/sites/default/files/nsdi18_full_proceedings_interior.pdf>

(58)調査した分野 (Int.Cl., DB名)

G06F 9/50

G06F 9/455