

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第766632号
(P766632)

(45)発行日 令和7年4月22日(2025.4.22)

(24)登録日 令和7年4月14日(2025.4.14)

(51)国際特許分類	F I			
G 0 6 F 11/34 (2006.01)	G 0 6 F	11/34	1 8 0	
G 0 6 F 15/82 (2006.01)	G 0 6 F	11/34	1 7 6	
	G 0 6 F	15/82	6 1 0 N	

請求項の数 9 (全25頁)

(21)出願番号	特願2023-559356(P2023-559356)	(73)特許権者	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(86)(22)出願日	令和3年11月12日(2021.11.12)	(74)代理人	100098394 弁理士 山川 茂樹
(86)国際出願番号	PCT/JP2021/041778	(74)代理人	100153006 弁理士 小池 勇三
(87)国際公開番号	WO2023/084749	(74)代理人	100064621 弁理士 山川 政樹
(87)国際公開日	令和5年5月19日(2023.5.19)	(74)代理人	100121669 弁理士 本山 泰
審査請求日	令和6年5月7日(2024.5.7)	(72)発明者	有川 勇輝 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
		(72)発明者	三浦 直樹

最終頁に続く

(54)【発明の名称】 コンピュータシステムおよびその制御方法

(57)【特許請求の範囲】

【請求項1】

複数の演算部とホスト部とを備え、前記複数の演算部がそれぞれトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記入力データの種別ごとに、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、

前記ホスト部と前記演算部とのいずれかが、前記トレースバッファ内に記録された、前記イベントの検出時刻が属する前記入力データの種別を判定するステップと、

前記ホスト部と前記演算部とのいずれかが、判定の結果に応じて、同一の前記入力データの種別に属する前記イベントの検出時刻のうち、最新の前記イベントの検出時刻を記録し、前記最新の前記イベントの検出時刻以外の前記イベントの検出時刻を消去するステップと

を備えるコンピュータシステムの制御方法。

【請求項2】

前記ホスト部が、前記トレースバッファを定期的に監視するステップとを備え、

前記ホスト部が、複数の前記イベントの検出時刻が同一の前記入力データの種別に属する場合に、前記イベントの検出時刻に基づき、前記トレースバッファ内に記録された前記イベントの検出時刻を消去する

ことを特徴とする請求項1に記載のコンピュータシステムの制御方法。

【請求項3】

前記演算部が、前記記録の前に、前記判定を実行し、前記イベントの検出時刻が同一の前記入力データの種別に属する場合に、前記イベントの検出時刻に基づき、前記トレースバッファ内に記録された前記イベントの検出時刻を上書きする

ことを特徴とする請求項 1 に記載のコンピュータシステムの制御方法。

【請求項 4】

複数の演算部がそれぞれ演算器とトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、

前記複数の演算部のうち、一の演算部の演算器が、前記入力データを処理し、他の演算部の演算器に転送するステップと、

前記他の演算部が、受信完了通知を発信するステップと、

前記一の演算部が、前記受信完了通知を受信すると、前記一の演算部の前記トレースバッファに記録される前記イベントの検出時刻を消去するステップと

を備えるコンピュータシステムの制御方法。

【請求項 5】

複数の演算部がそれぞれ演算器とトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、

所定の時間の経過を契機として、前記トレースバッファに記録された前記イベントの検出時刻を消去することを特徴とするコンピュータシステムの制御方法。

【請求項 6】

入力データを処理するコンピュータシステムであって、

複数の演算部と、

前記複数の演算部と接続し、前記複数の演算部を制御するホスト部と

を備え、

前記複数の演算部の間で前記処理されたデータが転送され、

前記演算部が、前記入力データからの所定のイベントの検出を契機として、トレースデータを記録するトレースバッファ

を備え、

前記トレースデータが、前記演算部の動作周波数を基にする前記イベントの検出時刻であるタイムスタンプ値を有し、

前記ホスト部と前記演算部とのいずれかが、所定の条件で、前記トレースバッファに記録される前記トレースデータを消去する

ことを特徴とするコンピュータシステム。

【請求項 7】

前記トレースデータが、さらに、前記入力データの種別と、前記イベントを検出する箇所を示す情報と、前記イベントの内容を区別する情報と、任意のデータとの少なくともいずれかを有する

ことを特徴とする請求項 6 に記載のコンピュータシステム。

【請求項 8】

前記ホスト部と前記演算部とのいずれかが、前記トレースバッファ内に記録された、複数の前記トレースデータが同一の前記入力データの種別に属する場合には、前記複数の前記トレースデータのうち、最新の前記トレースデータを保持し、前記最新の前記トレースデータ以外の前記トレースデータを消去する

ことを特徴とする請求項 6 又は請求項 7 に記載のコンピュータシステム。

【請求項 9】

前記複数の演算部のうち、一の演算部の演算器が、前記入力データを処理し、他の演算部の演算器に転送し、

前記他の演算部が、前記転送されたデータを受信し、受信完了通知を発信し、

前記一の演算部が、前記受信完了通知を受信し、前記一の演算部の前記トレースバッファ

10

20

30

40

50

アに記録される前記イベントの検出時刻を消去する

ことを特徴とする請求項 6 又は請求項 7 に記載のコンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数の演算部を有するコンピュータシステムおよびその制御方法に関する。

【背景技術】

【0002】

機械学習や人工知能 (AI) や IoT (Internet of Things) など多くの分野で技術革新が進み、様々な情報やデータを活用することで、サービスの高度化・付加価値の提供が盛んに行われている。このような処理では、大量の計算をする必要があり、そのための情報処理基盤が必須である。

10

【0003】

例えば、非特許文献 1 では、既存の情報処理基盤をアップデートする試みが展開されているが、急増するデータに対して現状のコンピュータが対応しきれていない。今後進展するためには、ムーアの法則を越える「ポストムーア技術」が確立されなければならないと指摘している。

【0004】

ポストムーア技術として、例えば、非特許文献 2 では、フローセントリックコンピューティングという技術が開示されている。フローセントリックコンピューティングでは、データのある場所で処理を行うという従来のコンピューティングの考えではなく、計算機能が存在する場所にデータを移動して処理を行うという新たな概念が導入された。

20

【0005】

上記のフローセントリックコンピューティングを実現するためには、データ移動に必要な広帯域な通信ネットワークが必要になるだけでなく、所望の演算性能が得るためには計算リソースを効率よく制御する必要がある。

【0006】

フローセントリックコンピューティング (例えば、非特許文献 2) では、複数の演算機能を連動させる手法が開示されている。

【先行技術文献】

30

【非特許文献】

【0007】

【文献】“NTT Technology Report for Smart World 2020,” 日本電信電話株式会社, 2020年. https://www.rd.ntt/_assets/pdf/techreport/NTT_TRFSW_2020_EN_W.pdf.

【文献】R. Takano and T. Kudoh, “Flow-centric computing leveraged by photonic circuit switching for the post-moore era,” Tenth IEEE/ACM International Symposium on Networks-on-Chip (NOCS), Nara, 2016, pp. 1-3. <https://ieeexplore.ieee.org/abstract/document/7579339>.

【発明の概要】

40

【発明が解決しようとする課題】

【0008】

しかしながら、複数の演算部が連動するコンピュータシステムにおいて、ホスト部を経由せずに演算部同士が主体的にデータを移動させるため、演算部内で生じた障害を特定することが困難であった。

【0009】

また、入力データがある時刻に通過した演算部を特定するなどコンピュータシステムの内部状態を把握することが困難であった。

【課題を解決するための手段】

【0010】

50

上述したような課題を解決するために、本発明に係るコンピュータシステムの制御方法は、複数の演算部とホスト部とを備え、前記複数の演算部がそれぞれトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記入力データの種別ごとに、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、前記ホスト部と前記演算部とのいずれかが、前記トレースバッファ内に記録された、前記イベントの検出時刻が属する前記入力データの種別を判定するステップと、前記ホスト部と前記演算部とのいずれかが、判定の結果に応じて、同一の前記入力データの種別に属する前記イベントの検出時刻のうち、最新の前記イベントの検出時刻を記録し、前記最新の前記イベントの検出時刻以外の前記イベントの検出時刻を消去するステップとを備える。

10

【0011】

また、本発明に係るコンピュータシステムの制御方法は、複数の演算部がそれぞれ演算器とトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、前記複数の演算部のうち、一の演算部の演算器が、前記入力データを処理し、他の演算部の演算器に転送するステップと、前記他の演算部が、受信完了通知を発信するステップと、前記一の演算部が、前記受信完了通知を受信すると、前記一の演算部の前記トレースバッファに記録される前記イベントの検出時刻を消去するステップとを備える。

【0012】

また、本発明に係るコンピュータシステムの制御方法は、複数の演算部とホスト部とを備え、複数の演算部がそれぞれトレースバッファを備え、入力データからの所定のイベントの検出を契機として、前記演算部の動作周波数を基に取得される前記イベントの検出時刻を前記トレースバッファに記録するコンピュータシステムの制御方法であって、前記ホスト部と前記演算部とのいずれかが、所定の時間の経過を契機として、前記トレースバッファに記録された前記イベントの検出時刻を消去することを特徴とする。

20

【0013】

また、本発明に係るコンピュータシステムは、入力データを処理するコンピュータシステムであって、複数の演算部と、前記複数の演算部と接続し、前記複数の演算部を制御するホスト部を備え、前記複数の演算部の間で前記処理されたデータが転送され、前記演算部が、前記入力データからの所定のイベントの検出を契機として、トレースデータを記録するトレースバッファを備え、前記トレースデータが、前記演算部の動作周波数を基にする前記イベントの検出時刻であるタイムスタンプ値を有し、前記ホスト部と前記演算部とのいずれかが、所定の条件で、前記トレースバッファに記録される前記トレースデータを消去することを特徴とする。

30

【発明の効果】

【0014】

本発明によれば、障害が生じた場合に、容易に、障害が生じた箇所を特定でき、障害発生時の内部でのデータ状況を把握できるコンピュータシステムおよびその制御方法を提供できる。

40

【図面の簡単な説明】

【0015】

【図1】図1は、本発明の第1の実施の形態に係るコンピュータシステムの構成を示すブロック図である。

【図2】図2は、本発明の第1の実施の形態に係るコンピュータシステムにおける演算部の構成を示すブロック図である。

【図3A】図3Aは、本発明の第1の実施の形態に係るコンピュータシステムにおける演算部の拡張および縮小を説明するための図である。

【図3B】図3Bは、本発明の第1の実施の形態に係るコンピュータシステムにおける演算部の拡張および縮小を説明するための図である。

50

【図 4 A】図 4 A は、本発明の第 1 の実施例に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【図 4 B】図 4 B は、本発明の第 1 の実施例に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【図 4 C】図 4 C は、本発明の第 1 の実施例に係るコンピュータシステムの制御方法を説明するための図である。

【図 5】図 5 は、本発明の第 2 の実施例に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【図 6】図 6 は、本発明の第 3 の実施例に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【図 7 A】図 7 A は、本発明の第 1 の実施の形態に係るコンピュータシステムの制御方法の一例を説明するための図である。

【図 7 B】図 7 B は、本発明の第 1 の実施の形態に係るコンピュータシステムの制御方法の一例を説明するための図である。

【図 8】図 8 は、本発明の第 2 の実施の形態に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【図 9】図 9 は、本発明の第 2 の実施の形態に係るコンピュータシステムの制御方法を説明するためのフローチャート図である。

【発明を実施するための形態】

【0016】

< 第 1 の実施の形態 >

本発明の第 1 の実施の形態に係るコンピュータシステムとその制御方法について、図 1 ~ 図 3 を参照して説明する。

【0017】

< コンピュータシステムの構成 >

本実施の形態に係るコンピュータシステム 10 は、図 1 に示すように、N 個の演算部 11_1 ~ 11_N (N は 1 以上の整数) と、演算部 11_1 ~ 11_N を接続する内部通信部 13 と、演算部 11_1 ~ 11_N に対して動作パラメータを設定・管理するホスト部 12 とを備える。

【0018】

演算部 11_1 ~ 11_N は、プロセッサやアクセラレータなどにより構成され、トレース部 14_1 ~ 14_N を備える。

【0019】

トレース部 14_1 ~ 14_N は、演算部 11_1 ~ 11_N が連動する際に、各演算部 11_1 ~ 11_N の任意の観測ポイントにおいて、所定のイベントの検出を契機として、各演算部 11_1 ~ 11_N の動作周波数を基にしたイベント検出時刻を記録する。

【0020】

ここで、トレース部 14_1 ~ 14_N は、データ種別またはイベント種別ごとに、イベント検出時刻を記録できる。さらに、任意のデータを記録してもよい。

【0021】

演算部 11_1 ~ 11_N の連動において、演算部 11_1 で処理されたデータが、内部通信部 13 を介して、演算部 11_2 に転送される。引き続き、データの転送が繰り返され、データが演算部 11_N に転送される。

【0022】

なお、演算部 11_1 ~ 11_N の連動の方法として、複数の演算部を直列に接続する処理方法や、複数の演算部を並列に接続する処理方法や、両者を組み合わせた処理方法などが挙げられる。複数の演算部が連動することで、所望のサービスを提供し、アプリケーションを処理する。

【0023】

演算部 11_1 ~ 11_N (N は 1 以上の整数) は、コンピュータシステム 10 の外部

10

20

30

40

50

から入力される入力データに対して所定の演算処理を実行する機能を有する。演算処理とは、例えば、画像データが入力された際に画像サイズを縮小・拡大する処理や、画像データから特定の物体を検出する処理、画像データを復号・暗号化する処理など、入力データに対する加工、集計、結合などの一般的な演算処理である。

【0024】

また、演算部11__1～11__Nは、当該システムの停止中・稼働中を問わず、追加・削除してもよい。例えば、演算器の一部分のみを動的再構成が可能なデバイスであるFPGAを用いることで実現することができる。また、演算部11__1～11__Nの実装方法として、特定の演算に特化した専用回路を具備するアクセラレータカードを追加しても良い。また、演算部の中に演算機能を提供する演算器を複数具備することもできる。

10

【0025】

ホスト部12は、演算部11__1～11__Nに対して動作パラメータを設定・管理する機能を有し、詳細には演算部11__1～11__Nを制御する機能や、データを記憶する機能を有する。動作パラメータは、例えば、画像処理において複数のアルゴリズムを切り替えて使う場合、アルゴリズムを特定するための情報であり、演算処理における係数や閾値などである。

【0026】

また、当該システムの稼働開始後であっても演算部の追加・削除ができる場合、ホスト部12は、その演算部に対して、所望の処理内容を実行するための回路情報を演算部に対して設定するなど、コンピュータシステム10全体の管理を行う。

20

【0027】

内部通信部13は、演算部11__1～11__Nを接続するとともに、演算部11__1～11__N間でデータの授受を行うための通信機能を有する。具体的には、PCIeやイーサネットなど市中の通信規格と前記通信規格を満足する物理構成、すなわちPCIeスイッチやイーサネットスイッチが挙げられる。

【0028】

また、演算部11__1～11__Nの中に演算機能を提供する演算器を複数具備する場合、前記観測ポイントを演算部11__1～11__Nの中に複数設けてもよい。

【0029】

<演算部の構成>

コンピュータシステム10における演算部11__1は、図2に示すように、複数(N台)の演算器15__1(1)～15__N(1)と、トレース部14__1とを備える。ここで、演算器は1台であってもよい。

30

【0030】

トレース部14__1は、イベントジェネレータ16__1__1～16__2__Nと、タイムスタンプ部17と、トレースバッファ18とを備える。

【0031】

イベントジェネレータ16__1__1～16__2__Nは、演算器15__1(1)～15__N(1)の入力側と出力側それぞれに接続される。タイムスタンプ部17の出力が、イベントジェネレータ16__1__1～16__2__Nに接続される。イベントジェネレータ16__1__1～16__2__Nの出力は、トレースバッファ18に接続される。

40

【0032】

ここで、イベントジェネレータ16__1__1～16__2__Nは、演算器15__1(1)～15__N(1)の入力側と出力側いずれかに配置されればよく、任意の箇所に配置されればよく、少なくとも1台配置されればよい。また、トレースバッファ18は複数台配置されてもよく、少なくとも1台配置されればよい。

【0033】

イベントジェネレータ16__1__1～16__2__Nは、演算部11__1～11__Nの任意の箇所に挿入されて、データの種別(ユーザID、セッションID、ストリームID、サービスID)ごとにイベント(ストリームの先頭、末尾)を検出し、検出時刻(以下、

50

「タイムスタンプ値」という。)を含むトレースデータを後述するトレースバッファ18に記録する契機を発生させる。

【0034】

また、データの種別は上記に限らず、データを整理するために用いるパケットのヘッダ情報や、データと並走する信号が有する情報など、データの整理に利用できる情報であれば適用することができる。

【0035】

タイムスタンプ部17は、少なくとも1台のクロックカウンタを備え、複数のイベントジェネレータ16_1_1~16_2_N(観測ポイント)間を同期させるとともに、演算部11_1~11_Nの動作周波数の精度で時刻を取得する。ここで、演算部11_1~11_Nの動作周波数(クロック周波数)は、FPGA(field-programmable gate array)を用いて当該機能を実現する場合、通常、数ナノ秒程度である。

10

【0036】

トレースバッファ18は、イベントジェネレータ16_1_1~16_2_Nによるイベント検出を契機として、トレースデータを記録する。ここで、トレースデータは、タイムスタンプ部17から取得した各検出時刻(タイムスタンプ値)と、インスタンスIDと、イベント種別(イベントID)と、データ種別(TID)と、任意のデータとを有する。ここで、トレースデータは、少なくともタイムスタンプ値を有すればよい。

【0037】

また、トレースバッファ18は、イベントジェネレータの数に依存することなく、一定のバッファ量を提供する。

20

【0038】

ここで、タイムスタンプ値は、演算部(FPGA)内で統一された値である。

【0039】

また、インスタンスIDは、イベントジェネレータ・インスタンスを区別し、前記イベントを検出する箇所(観測ポイント)を示すIDである。

【0040】

また、イベント種別(イベントID)は、イベント内容を区別するIDである。例えば、ストリームの先頭の通過またはストリームの終端の通過により区別する。また、データの任意の箇所にイベント検出用のフラグ等を用意して、当該フラグが通過したことを検出する。

30

【0041】

また、任意のデータは、画像データ、数値データ、文章データなど通常コンピュータシステムで処理されるデータである。

【0042】

また、データ種別は、例えば、入力データの属性などを識別・分類したりするために用いられ、ユーザID、セッションID、ストリームID、サービスIDなど、データ本体に付随する情報である。また、データ種別を識別するための情報は必ずしもパケットのヘッダに付与されなくともよく、例えば、パケットのペイロードに独自に定義してもよい。また、演算部の内部において、データと並走する信号を用いる場合、データ種別を取得するのに並走信号を用いてもよい。

40

【0043】

<コンピュータシステムの動作>

本実施の形態に係るコンピュータシステム10の動作を、以下に説明する。

【0044】

コンピュータシステム10の演算部11_1において、演算器15_1(1)~15_N(1)にデータが入力される。入力されるデータは、様々な要素から構成され、イベント種別(イベントID)と、データ種別(TID)と、任意のデータを含む。

【0045】

50

ここで、演算部 1 1 __ 1 の演算器 1 5 __ 1 (1) で処理されたデータは送信され、演算部 1 1 __ 2 の演算器 1 5 __ 1 (2) に入力する。

【 0 0 4 6 】

イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N において、まず、演算部 1 1 __ 1 ~ 1 1 __ N に入力するデータの制御信号 (動作) を観測する。

【 0 0 4 7 】

イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N がイベントを検出すると、イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N が、入力データから、イベント種別と、データ種別と、任意のデータを取得する。ここで、イベント発生は、例えば、ストリームの先頭が通過した時、またはストリームの先頭が通過した時である。

【 0 0 4 8 】

イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N で、取得されたイベント種別とデータ種別と任意のデータに、インスタンス ID と、タイムスタンプ部 1 7 から送信されたタイムスタンプ値とが加えられる。その結果、トレースデータは、タイムスタンプ値と、インスタンス ID と、イベント種別と、データ種別と、任意のデータとで構成される。

【 0 0 4 9 】

ここで、トレースデータは、少なくともタイムスタンプ値を含めばよく、タイムスタンプ値を基に処理時間、データ流量などのコンピュータシステム内の情報を把握できる。また、不具合が生じた時刻なども把握できる。

【 0 0 5 0 】

さらに、トレースデータは、インスタンス ID を有することにより、不具合が生じた箇所を把握できる。

【 0 0 5 1 】

また、トレースデータは、イベント種別を有することにより、イベント発生時を把握できる。

【 0 0 5 2 】

また、トレースデータは、データ種別を有することにより、データ種別ごとの稼働状況を把握でき、データ消去時の判定 (後述) に用いることができる。

【 0 0 5 3 】

また、トレースデータは、任意のデータを有することにより、処理の再開 (後述) 時に用いることができる。

【 0 0 5 4 】

また、トレースデータは、サービスの優先度情報を有することにより、優先度に基づくトレースデータの管理に用いることができる。

【 0 0 5 5 】

最後に、イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N が、トレースデータをトレースバッファ 1 8 に送信する。

【 0 0 5 6 】

また、タイムスタンプ部 1 7 において、ホスト部 1 2 から送信されたカウン트의開始又は停止設定が受信 (書き込み) される。

【 0 0 5 7 】

これを契機に、クロックカウンタにカウン트의開始又は停止が設定される。

【 0 0 5 8 】

カウン트의開始設定によりカウン트는開始され、カウン트는各演算部 1 1 __ 1 ~ 1 1 __ N の動作周波数を基にして実行される。一方、カウン트의停止設定によりカウン트는停止される。

【 0 0 5 9 】

イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N はイベントを検出すると、クロックカウンタがカウントする時刻を、タイムスタンプ値として読み出し、タイムスタンプ値がイベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N に送信される。

10

20

30

40

50

【 0 0 6 0 】

また、イベントの検出は例えば、データと並走する信号のうち、有効なデータであるか否かを示す信号のON/OFF用いて判定したり、データの特定の領域にイベント検出用のフィールドを用意して、当該フィールドのビット列を用いたりすることで、イベントを検出する。

【 0 0 6 1 】

ここで、必要に応じて、演算部（FPGA）間で同期させる。なお、演算部間で同期させる方法として、ホスト部から同期をする演算部に対して同期させるための信号や、リセット信号、を入力することで、同期を図る。

【 0 0 6 2 】

また、ホスト部 1 2 は、トレースバッファ 1 8 からトレースデータを読み出す時に、タイムスタンプ部 1 7 にリセット信号を送信し、クロックカウンタの値をリセットする。

【 0 0 6 3 】

トレースバッファ 1 8 において、イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N より受信されたトレースデータが、トレースバッファ 1 8 に記録され蓄積される。

【 0 0 6 4 】

ここで、複数のイベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N から送信されたトレースデータが記録される。また、トレースデータの書き込み、読み出しは全 T I D で共通の F I F O (F i r s t - I n F i r s t - O u t) で実行される。

【 0 0 6 5 】

次に、ホスト部 1 2 が、トレースデータを、トレースバッファ 1 8 から読み出す。

【 0 0 6 6 】

最後に、ホスト部 1 2 で、読み出した（回収）データの後処理を実行する。詳細には、T I D ごとに検索（G R E P）する。引き続き、タイムスタンプ部 1 7 でソートした後に可視化する。

【 0 0 6 7 】

ここでは、イベントジェネレータが、イベント検出時に入力データから、イベント種別と、データ種別と、任意のデータを取得する例を示したが、イベント種別と、データ種別と、任意のデータを取得しなくとも、上述の通り、少なくともタイムスタンプ値を取得すればコンピュータシステム 1 0 を動作できる。

【 0 0 6 8 】

< 演算部および演算器の追加・削除 >

コンピュータシステム 1 0 において、演算部 1 1 __ 1 ~ 1 1 __ N を、コンピュータシステム 1 0 の停止中・稼働中を問わず、追加または削除できる。

【 0 0 6 9 】

また、例えば、図 3 A に示す演算部 1 1 __ 1 において、図 3 B に示すように、イベントジェネレータ 1 6 __ 1 __ 2 ~ 1 6 __ 2 __ N を新たに配置（追加）してタイムスタンプ部 1 7 とトレースバッファ 1 8 と接続することにより、演算器 1 5 __ 1 (1) ~ 1 5 __ N (1) を追加できる。

【 0 0 7 0 】

また、演算部 1 1 __ 1 において、イベントジェネレータ 1 6 __ 1 __ 2 ~ 1 6 __ 2 __ N を削除することにより、演算器 1 5 __ 1 (1) ~ 1 5 __ N (1) を削除できる。

【 0 0 7 1 】

演算器 1 5 __ 1 (1) ~ 1 5 __ N (1) を追加するとき、トレース部 1 4 __ 1 ~ 1 4 __ N において、複数のイベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N を任意の箇所に配置することができる。

【 0 0 7 2 】

ここで、イベントジェネレータ 1 6 __ 1 __ 1 ~ 1 6 __ 2 __ N を演算器 1 5 __ 1 (1) ~ 1 5 __ N (1) の入力の前段と出力の後段との両方に配置してもよく、入力の前段と出力の後段とのいずれか一方に配置してもよい。

10

20

30

40

50

【0073】

演算器15__1(1)~15__N(1)を削除するとき、トレース部14__1~14__Nにおいて、演算器15__1(1)~15__N(1)の前後のイベントジェネレータ16__1__1~16__2__Nを削除すればよい。

【0074】

ここで、削除する演算器15__1(1)~15__N(1)の入力の前段と出力の後段との両方のイベントジェネレータを削除してもよく、入力の前段と出力の後段とのいずれか一方のイベントジェネレータを削除してもよい。

【0075】

本実施の形態によれば、複数のイベントジェネレータを任意の箇所を追加できるため、イベントジェネレータのみを新たに追加することにより、容易に演算部の内部に新たに演算器を追加でき、演算器における処理時間の計測、トレースデータの収集が可能となる。

10

【0076】

また、複数のイベントジェネレータを任意の箇所から削除できるため、イベントジェネレータのみを削除すれば、容易に演算部の内部で演算器を削除できる。

【0077】

また、演算器を追加する場合にイベントジェネレータとタイムスタンプ部とトレースバッファとを全て追加する場合に比べて、回路規模を縮小でき、消費電力を抑制できる。

【0078】

また、演算部から演算器を削除する場合、演算器の前後に配置されているイベントジェネレータを削除してもよい。このとき、イベントジェネレータに関するトレースデータを削除してもよい。また、削除されるイベントジェネレータで検出したイベントに付随するトレースデータは、必ずしも削除しなくともよい。

20

【0079】

<第1の実施例>

本発明の第1の実施例に係るコンピュータシステム10の制御方法について、図4A、Bを参照して説明する。

【0080】

本実施例では、コンピュータシステム10において、演算部11__1~11__Nのトレース部14__1~14__Nが記録するトレースデータを用いて、不具合時に効率的に処理を再開する。ここで、不具合とは、通常発生し得るパケットロスや、演算部11__1~11__Nの内部の機能ブロックにおける処理のスタックなどである。

30

【0081】

コンピュータシステム10において、例えば、演算部11__1のトレース部14__1は、トレースバッファ18で任意のトレースデータとして、タイムスタンプ値と、インスタンスID(イベントを検出する箇所を示す情報)と、任意のデータとを記録する。

【0082】

本実施例に係るコンピュータシステム10の制御方法の一例として、ホスト部12が制御(管理)する場合を、図4Aを参照して説明する。

【0083】

初めに、ホスト部12が所定の周期でコンピュータシステムのトレースバッファ18をモニタして、演算部内の処理時間を計測する(ステップS11A)。

40

【0084】

ここで、処理時間の計測において、まず、ホスト部12が、演算器15__1(1)の入力側のイベントジェネレータ16__1__1(第1のイベントジェネレータ)と、演算器15__1(1)の出力側のイベントジェネレータ16__2__1(第2のイベントジェネレータ)それぞれの所定のイベント(第1のイベント、第2のイベント)の検出を契機とする第1のタイムスタンプ値と第2のスタンプ値を、トレースバッファ18から読み出す(取得する)。

【0085】

50

引き続き、ホスト部 1 2 が、第 1 のタイムスタンプ値と第 2 のスタンプ値との差分より、処理時間を算出する。

【 0 0 8 6 】

このように、任意の箇所（区間、例えば、演算器 1 5 __ 1 (1) ）の処理時間は、任意の箇所（区間、例えば、演算器 1 5 __ 1 (1) の前後）に配置されたイベントジェネレータそれぞれを入力データが通過する時刻（第 1 のタイムスタンプ値と第 2 のスタンプ値）の差分により取得される。

【 0 0 8 7 】

次に、計測された処理時間を、予め設定された所定の閾値と比較する（ステップ S 1 2 A ）。その結果、処理時間が所定の閾値より長い場合に、不具合が発生したと判定する。

10

【 0 0 8 8 】

最後に、不具合が発生した場合には、インスタンス ID より処理検出箇所を把握して、この処理検出箇所よりも前段のいずれかのトレースバッファ 1 8 に記録される任意のデータを用いて、処理を再開する（ステップ S 1 3 A ）。

【 0 0 8 9 】

また、本実施例に係るコンピュータシステム 1 0 の制御方法の一例として、演算部 1 1 __ 1 ~ 1 1 __ N が制御（管理）する場合を、図 4 B、C を参照して説明する。

【 0 0 9 0 】

初めに、演算部 1 1 __ 1 ~ 1 1 __ N は、トレースデータを記録するとともに、所定の周期でトレースバッファ 1 8 を監視し、任意の演算部（例えば演算部 1 1 __ 1 ）のトレースバッファ 1 8 にタイムスタンプ値が記録されてから、次段の演算部（例えば演算部 1 1 __ 2 ）の演算器 1 5 __ 1 (2) にデータが入力されるまでの時間（以下、「演算部間の処理時間」という。）を計測する（ステップ S 1 1 B ）。

20

【 0 0 9 1 】

この演算部間の処理時間の計測では、例えば、図 4 C に示すように、まず、演算部 1 1 __ 1 のトレースバッファ 1 8 の前段のイベントジェネレータ 1 6 __ 2 __ 1 がタイムスタンプ値（第 1 のスタンプ値）を取得し、演算部 1 1 __ 1 のトレースバッファ 1 8 で記録される。

【 0 0 9 2 】

引き続き、演算部 1 1 __ 1 のトレースバッファ 1 8 から演算部 1 1 __ 2 の演算器 1 5 __ 1 (2) に通知信号が送信され（図中、点線矢印）、この信号を契機に、演算部 1 1 __ 1 の演算器 1 5 __ 1 (1) から演算部 1 1 __ 2 の演算器 1 5 __ 1 (2) にデータが転送される。

30

【 0 0 9 3 】

引き続き、演算部 1 1 __ 2 の演算器 1 5 __ 1 (2) に入力されるデータのイベントの検出を契機に、演算部 1 1 __ 2 の演算器 1 5 __ 1 (2) の前段のイベントジェネレータ 1 6 __ 1 __ 1 が、タイムスタンプ値（第 2 のスタンプ値）を取得し、演算部 1 1 __ 2 のトレースバッファ 1 8 で記録される。

【 0 0 9 4 】

これらの第 1 のスタンプ値と第 2 のスタンプ値の差分より、演算部間の処理時間が計測される。

40

【 0 0 9 5 】

次に、計測された処理時間を、予め設定された所定の閾値と比較する（ステップ S 1 2 B ）。その結果、計測された時間が所定の閾値より長い場合に、不具合が発生したと判定する。

【 0 0 9 6 】

最後に、不具合が発生した場合には、トレースバッファ 1 8 に記録されるインスタンス ID により把握される処理検出箇所よりも前段のいずれかのトレースバッファ 1 8 に記録される任意のデータを用いて、処理を再開する（ステップ S 1 3 B ）。

【 0 0 9 7 】

50

ここで、トレースバッファに記録したタイムスタンプ値を用いて処理時間を計測する例を示したが、イベントジェネレータが取得したタイムスタンプ値を直接用いて処理時間を計測してもよい。

【0098】

このように、本実施例に係るコンピュータシステムの制御方法では、任意の位置に配置されるイベントジェネレータが、所定のイベントを検出し、これを契機として、一のタイムスタンプ値を取得し、他の位置に配置されるイベントジェネレータが、同様に他のタイムスタンプ値を取得し、一のタイムスタンプ値と、他のタイムスタンプ値との差分を算出することにより、不具合の発生を判定し、処理を再開する。

【0099】

本実施例に係るコンピュータシステムの制御方法によれば、トレース部14__1~14__Nが記録するトレースデータを用いることにより、不具合が発生した場合に、処理を最初から再開せずに、正常に動作していた箇所まで遡り処理を再開できる。

【0100】

また、トレースデータの処理を再開させる契機は、ホスト部12が管理する。または、演算部11__1~11__Nが管理してもよい。

【0101】

また、処理の再開は、必ずしも特定の機能ブロックに限定しなくともよく、例えば、演算部の入力まで遡って処理を再開してもよい。

【0102】

本実施例では、インスタンスIDにより処理検出箇所を把握する例を示したが、これに限らず、予め設定される演算器の処理速度と、計測されるタイムスタンプ値とを用いて処理検出箇所を導出してもよい。

【0103】

本実施例では、不具合の発生を把握するために、処理時間を用いる例を示したが、これに限らず、データ流量（後述）などを用いてもよい。

【0104】

本実施例では、トレースデータは、データ種別、イベント種別を有してもよい。また、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

【0105】

<第2の実施例>

本発明の第2の実施例に係るコンピュータシステムの制御方法について、図5を参照して説明する。本実施例では、コンピュータシステムにおいて、演算部11__1~11__Nのトレース部14__1~14__Nを用いて、システムの品質管理（ステート管理/ヘルスチェック）を実行する。

【0106】

コンピュータシステム10において、例えば、演算部11__1のトレース部14__1は、複数のイベントジェネレータ16__1__1~16__2__Nを備える。

【0107】

図5に、本実施例に係るコンピュータシステム10の制御方法のフローチャート図を示す。

【0108】

初めに、例えば、コンピュータシステム10の演算部11__1は、複数のイベントジェネレータ16__1__1~16__2__Nを用いて、データが各イベントジェネレータ16__1__1~16__2__Nを通過する時刻を収集する（ステップS21）。

【0109】

詳細には、演算部11__1が、異なる箇所に配置されるイベントジェネレータ（例えば、イベントジェネレータ16__1__1とイベントジェネレータ16__2__1）それぞれでの所定のイベントの検出に基づき、タイムスタンプ値（例えば、第1のタイムスタンプ値と第2のタイムスタンプ値）を収集する。

10

20

30

40

50

【0110】

次に、収集された第1のタイムスタンプ値と第2のタイムスタンプ値との差分を求めることにより、演算部11__1内の特定区間の通過に要する時間すなわち処理時間を計算する(ステップS22)。

【0111】

次に、処理時間を、予め設定された所定の閾値と比較する(ステップS23)。

【0112】

比較の結果、処理時間が閾値より大きい場合に、コンピュータシステム10の異常検知を通知する(ステップS24)。

【0113】

このように、本実施例に係るコンピュータシステムの管理方法によれば、コンピュータシステムが正常に稼働しているか否かをモニタできる。

【0114】

また、必ずしもすべてのデータに対して時間の解析処理を実行しなくともよい。例えば、予め指定した計測間隔で演算部11__1~11__N内の特定区間の通過に要した時間(処理時間)を観測し、所定の範囲内に収まるか、もしくは所定の閾値よりも処理時間が伸びているか否かを判定することにより、コンピュータシステムが正常に稼働しているか否かをモニタできる。

【0115】

また、テスト用のデータを入力して、このデータが演算部11__1~11__Nの内部を通過するのに要する時間を解析してもよい。

【0116】

本実施例では、イベントジェネレータが取得したタイムスタンプ値を直接用いて処理時間を計測する例を示したが、トレースバッファにタイムスタンプ値を記録した後に、トレースバッファに記録したタイムスタンプ値を用いて処理時間を計測してもよい。

【0117】

本実施例では、処理時間を用いてコンピュータシステムの状況を把握する例を示したが、これに限らず、データ流量などを用いてもよい。

【0118】

本実施例では、演算部がコンピュータシステムを制御する例を示したが、ホスト部がコンピュータシステムを制御してもよい。

【0119】

本実施例では、トレースデータは、タイムスタンプ値とともに、データ種別を有してもよい。さらに、インスタンスID、イベント種別、任意のデータを有してもよい。また、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

【0120】

<第3の実施例>

本発明の第3の実施例に係るコンピュータシステム10の制御方法について、図6を参照して説明する。本実施例では、コンピュータシステム10において、演算部11__1~11__Nのトレース部14__1~14__Nを用いて、コンピュータシステム10のフロー管理を実行する。

【0121】

コンピュータシステム10において、演算部11__1~11__Nのトレース部14__1~14__Nがトレースデータとして、タイムスタンプ値とインスタンスID(イベントを検出する箇所を示す情報)をトレースバッファ18に記録し、ホスト部12がトレースデータを読み出す。

【0122】

図6に、本実施例に係るコンピュータシステム10の制御方法のフローチャート図を示す。

【0123】

10

20

30

40

50

初めに、ホスト部 1 2 が、トレースバッファ 1 8 よりトレースデータとして、例えば、任意の箇所でイベントジェネレータ 1 6 _ 1 が異なるイベントで取得したタイムスタンプ値とインスタンス ID (イベントを検出する箇所を示す情報) とを収集する (ステップ S 3 1) 。

【 0 1 2 4 】

詳細には、イベントジェネレータ 1 6 _ 1 において、通過するデータのイベント、例えば先頭および終端それぞれの検出を契機として取得され、トレースバッファ 1 8 に記憶された先頭および終端のタイムスタンプ値を収集する。

【 0 1 2 5 】

次に、先頭のタイムスタンプ値と終端のタイムスタンプ値との差分をデータの通過時間として算出する。

10

【 0 1 2 6 】

次に、予め設定されている入力データ (または出力データ) のデータ量を、データの通過時間で除することにより、所定の箇所における単位時間当たりのデータ量 (データ流量) を計算する (ステップ S 3 2) 。

【 0 1 2 7 】

次に、データ流量を、予め設定された所定の閾値と比較する (ステップ S 3 3) 。

【 0 1 2 8 】

比較の結果、データ流量が閾値より大きい場合に、データの集中を回避するようにデータフロー (経路) を設定する。例えば、経路を割り当てる際に、インスタンス ID (イベントを検出する箇所を示す情報) により把握される所定の閾値を超えた経路を回避して、経路を設定する (ステップ S 3 4) 。

20

【 0 1 2 9 】

このように、本実施例に係るコンピュータシステムの管理方法によれば、データの集中を回避するようにデータフロー (経路) を設定できる。

【 0 1 3 0 】

本実施例で、トレースバッファに記録したタイムスタンプ値を用いて処理時間を計測する例を示したが、イベントジェネレータが取得したタイムスタンプ値を直接用いて処理時間を計測してもよい。

【 0 1 3 1 】

また、ホスト部がコンピュータシステムを制御する例を示したが、演算部がコンピュータシステムを制御してもよい。

30

【 0 1 3 2 】

また、トレースデータが、データ種別に関する情報を有すれば、ホスト部 1 2 は、データ種別ごとの稼働状況を把握できる。

【 0 1 3 3 】

これにより、コンピュータシステム 1 0 では、特定のフローにのみデータが集中している場合に、当該フローから他の負荷の低いフローにデータを移行させることにより、データの集中を回避でき、フロー管理を実行できる。

【 0 1 3 4 】

また、トレースデータは、インスタンス ID、イベント種別、任意のデータが記録されてもよい。また、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

40

【 0 1 3 5 】

また、特定のフローや特定の箇所に発生した障害を検出すれば、障害が発生した際に、当該フローを迂回する経路を設定するようにフローを管理できる。

【 0 1 3 6 】

また、データの経路が異なる複数のフローを有する演算部において不具合が発生した場合に、当該フロー以外のフローを他の演算部に回避させた後に、当該演算部の交換、リセット、解析などを実施できる。

50

【 0 1 3 7 】

< 演算部の計測例 >

本実施例に係るコンピュータシステム 10 の制御方法において、演算部における処理時間とデータ流量との計測の一例を、図 7 A、B を参照して説明する。

【 0 1 3 8 】

本計測例では、例えば、演算器 15 __ 1 (1) の前段のイベントジェネレータ 16 __ 1 __ 1 で入力データが観測され、後段のイベントジェネレータ 16 __ 2 __ 1 で入力データが観測される。

【 0 1 3 9 】

入力データの先頭がイベントジェネレータ 16 __ 1 __ 1 を通過する時をイベントの契機として、イベントジェネレータ 16 __ 1 __ 1 に入力データの先頭のタイムスタンプ値が取得される。

10

【 0 1 4 0 】

同様に、入力データの終端がイベントジェネレータ 16 __ 1 __ 1 を通過する時をイベントの契機として、イベントジェネレータ 16 __ 1 __ 1 に入力データの終端のタイムスタンプ値が取得される。

【 0 1 4 1 】

一方、出力データの先頭がイベントジェネレータ 16 __ 2 __ 1 を通過する時をイベントの契機として、イベントジェネレータ 16 __ 2 __ 1 に出力データの先頭のタイムスタンプ値が取得される。

20

【 0 1 4 2 】

同様に、出力データの終端がイベントジェネレータ 16 __ 2 __ 1 を通過する時をイベントの契機として、イベントジェネレータ 16 __ 2 __ 1 に出力データの終端のタイムスタンプ値が取得される。

【 0 1 4 3 】

図 7 A に、トレースデータの一例を示す。図 7 A において、タイムスタンプ値 (T i m e s t a m p : D e c 、 T i m e s t a m p : 0 x) 、 インスタンス ID (I n s) 、 イベント ID (E v t) 、 デコード化されたイベント ID (D e c) 、 T I D 、 イベントデータ (E v e n t D a t a) を示す。

【 0 1 4 4 】

デコード化されたイベント ID (D e c) において、H はデータの先頭、L はデータの終端を示す。

30

【 0 1 4 5 】

入力データと出力データのデータ量は 1 M B である。また、演算部 11 __ 1 ~ 11 __ N の動作周波数は 250 M H z 、 4 n s / サイクル) である。

【 0 1 4 6 】

入力データの先頭のタイムスタンプ値 (T i m e s t a m p : D e c) は、「 4 0 6 5 1 4 」であり、インスタンス ID (I n s) が「 1 0 」である (点線四角 4 1 内上段) 。また、イベント ID (D e c) の「 H - R - 」がイベント発生としてデータ先頭の通過を示す。

40

【 0 1 4 7 】

同様に、出力データの先頭のタイムスタンプ値 (T i m e s t a m p : D e c) は、「 4 0 1 6 5 6 」であり、インスタンス ID (I n s) が「 1 1 」である (点線四角 4 1 内下段) 。また、イベント ID (D e c) の「 H - R - 」がイベント発生としてデータ先頭の通過を示す。

【 0 1 4 8 】

また、入力データの終端のタイムスタンプ値 (T i m e s t a m p : D e c) は、「 5 4 7 7 9 1 」であり、インスタンス ID (I n s) が「 1 0 」である (点線四角 4 2 内上段) 。また、イベント ID (D e c) の「 - L R - 」がイベント発生としてデータ終端の通過を示す。

50

【 0 1 4 9 】

同様に、出力データの先頭のタイムスタンプ値 (T i m e s t a m p : D e c) は、「 5 4 7 7 9 4 」であり、インスタンスID (I n s) が「 1 1 」である (点線四角 4 2 内下段)。また、イベントID (D e c) の「 - L R - 」がイベント発生としてデータ終端の通過を示す。

【 0 1 5 0 】

図 7 B に、入力データ 4 3 と出力データ 4 4 との関係を模式的に示す。入力データ 4 3 がイベントジェネレータ 1 6 _ 1 _ 1 を通過する時間は、終端のタイムスタンプ値 (5 4 7 7 9 1 サイクル) と先頭のタイムスタンプ値 (4 0 1 5 6 4 サイクル) との差分 (矢印 4 5) より、 $1 4 6 2 2 7 \text{ サイクル} = 5 8 4 . 9 \mu \text{ s e c}$ と算出される。したがって、入力スループットすなわちデータ流量として、 $1 \text{ M B} / 5 8 4 . 9 \mu \text{ s e c} = \text{約 } 1 . 8 \text{ G B} / \text{s e c}$ が得られる。

10

【 0 1 5 1 】

同様に、出力データ 4 4 の終端のタイムスタンプ値と先頭のタイムスタンプ値との差分 (矢印 4 6) より、入力スループットすなわちデータ流量を算出できる。

【 0 1 5 2 】

このように、演算部において、データ流量は、データ量を、データのタイムスタンプ値の先頭と終端との差分で除することにより得られる。

【 0 1 5 3 】

また、処理時間 (レイテンシ) として、出力開始と入力開始のタイムスタンプ値の差分すなわち出力データ 4 4 の先頭のタイムスタンプ値 (4 0 1 6 5 6 サイクル) と入力データ 4 3 の先頭のタイムスタンプ値 (4 0 1 5 6 4 サイクル) との差分より、9 2 サイクルが得られる。

20

【 0 1 5 4 】

このように、演算部において、入力データの処理時間は、出力開始と入力開始のタイムスタンプ値の差分より得られる。

【 0 1 5 5 】

以上のように、本発明に係るコンピュータシステムの制御方法において、入力データと出力データのタイムスタンプ値を用いて、データの処理時間とデータ流量とを取得できる。

【 0 1 5 6 】

< 効果 >

本発明の実施の形態および実施例に係るコンピュータシステムおよびその管理方法によれば、コンピュータシステムにおける演算処理が、途中で停止する場合または正常に処理が完了しない場合に、複数の演算部のうち、処理が停止した演算部を容易に検出、特定できる。

30

【 0 1 5 7 】

また、トレースバッファ 1 8 がデータを記録しているため、処理の途中から再開することができる。その結果、すでに実行した処理を再度、最初から繰り返す必要がなくなる。また、正常に処理が完了しない場合に、処理時間を短縮できる。

【 0 1 5 8 】

また、ホスト部 1 2 で各演算部の状態を一元管理できるため、例えば、処理が停止した演算部を経由しないデータのフローを設定でき、正常に処理が完了しないデータの数を削減できる。

40

【 0 1 5 9 】

また、データ種別 (ユーザID、セッションID、ストリームID、サービスID) ごとにイベントを検出できトレースデータを蓄積できるため、容易に、特定のデータ種別に着目して品質管理や不具合解析等を実行できる。

【 0 1 6 0 】

また、演算部と独立してトレース部 1 4 _ 1 ~ 1 4 _ N を備えるため、演算部の異常状態を保持できる。

50

【0161】

また、ユーザ毎（セッション毎）の粒度でフローを管理できる。

【0162】

また、イベント検出するイベントジェネレータを任意の部分に挿入できるので、特定のフローのみに発生する不具合を検出できる。

【0163】

<第2の実施の形態>

本発明の第2の実施の形態に係るコンピュータシステムおよびその制御方法を、図8を参照して説明する。本実施の形態に係るコンピュータシステム10は、第1の実施の形態と同様の構成を有する。

【0164】

第1の実施の形態に係るコンピュータシステム10では、トレースバッファ18がオーバーフローして、トレースデータを記録することが困難となる。そこで、トレースデータを消去する必要がある。

【0165】

<コンピュータシステムの制御方法>

図8に、本実施の形態に係るコンピュータシステムの制御方法のフローチャート図を示す。

【0166】

本実施の形態に係るコンピュータシステムでは、トレースバッファ18に、トレースデータとして、タイムスタンプ値とともに、少なくともデータ種別が記録される。ここで、トレースデータとして、インスタンスID、イベント種別、任意のデータが記録されてもよい。または、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

【0167】

初めに、ホスト部12が、演算部11_1～11_Nのトレースバッファ18を所定の周期で監視する（ステップS51）。

【0168】

次に、トレースバッファ18内に同一のデータ種別に対するトレースデータが複数記録されているか否かを判定する（ステップS52）。

【0169】

判定の結果、トレースバッファ18内に同一のデータ種別に対するトレースデータが複数記録されている場合には、この複数のトレースデータのうち、タイムスタンプ値が最新のトレースデータを保持（記録）するとともに、過去に記録したトレースデータ（最新のトレースデータ以外のトレースデータ）を消去する（ステップS53）。このように、前記イベントの検出時刻に基づき、前記トレースバッファ内に記録された前記イベントの検出時刻を消去する。

【0170】

本実施の形態に係るコンピュータシステムおよびその制御方法によれば、トレースバッファ18の蓄積容量が有限であっても、バッファのオーバーフローを抑制でき、物理的なバッファ量を削減できる。

【0171】

また、バッファ量を過剰に搭載する必要がないので、演算部の消費電力削減でき、電力効率を向上できる。

【0172】

また、データ種別（ユーザID、セッションID、ストリームID、サービスID）ごとにトレースデータを記録することにより、特定のデータ種別（例えば最高優先のサービスなど）については、比較的長時間トレースデータを保持することで信頼性を高めるなど、柔軟性の高いコンピュータシステムを提供できる。

【0173】

10

20

30

40

50

また、本実施の形態は、当然、第 1 の実施の形態と同様の効果を奏する。

【 0 1 7 4 】

< 第 2 の実施の形態の変形例 1 >

本発明の第 2 の実施の形態の変形例 1 に係るコンピュータシステムおよびその制御方法を説明する。本変形例に係るコンピュータシステム 10 は、第 1 の実施の形態と同様の構成を有する。

【 0 1 7 5 】

本変形例に係るコンピュータシステムでは、トレースバッファ 18 に、トレースデータとして、タイムスタンプ値とともに、少なくともデータ種別が記録される。ここで、トレースデータとして、インスタンス ID、イベント種別、任意のデータが記録されてもよい。または、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

10

【 0 1 7 6 】

本変形例では、コンピュータシステム 10 の演算部 11__1 ~ 11__N が、イベントジェネレータ 16__1__1 ~ 16__2__N より受信されたトレースデータ（最新のトレースデータ）を、トレースバッファ 18 に記録する（書き込む）前に、トレースバッファ 18 が既に保持しているトレースデータの中に、最新のトレースデータと同一のデータ種別のトレースデータが記録されているか否かを判定する。

【 0 1 7 7 】

判定の結果、同一のデータ種別のトレースデータが記録されている場合には、既に保持しているデータ（最新のデータ以外のデータ）を消去し、引き続き、最新のトレースデータを記録する、すなわち最新のトレースデータを上書きする。このように、前記イベントの検出時刻に基づき、前記トレースバッファ内に記録された前記イベントの検出時刻を上書きする。

20

【 0 1 7 8 】

ここで、上書きを防ぐためにフラグ等を付与し、フラグの有無で上書きの可否を判定してもよい。

【 0 1 7 9 】

これにより、本実施の形態と同様の効果を奏する。

【 0 1 8 0 】

< 第 2 の実施の形態の変形例 2 >

本発明の第 2 の実施の形態の変形例 2 に係るコンピュータシステムおよびその制御方法を説明する。本変形例に係るコンピュータシステム 10 は、第 1 の実施の形態と同様の構成を有する。

30

【 0 1 8 1 】

本変形例に係るコンピュータシステムでは、トレースバッファ 18 に、トレースデータとして、少なくともタイムスタンプ値が記録される。ここで、トレースデータとして、データ種別、インスタンス ID、イベント種別、任意のデータが記録されてもよい。または、トレースデータが、データ種別またはイベント種別ごとに記録されてもよい。

【 0 1 8 2 】

本変形例では、初めに、前段の演算部 11__1 の演算器 15__1 (1) から後段の演算部 11__2 の演算器 15__1 (2) に、演算器 15__1 (1) で処理されたデータが送信される。

40

【 0 1 8 3 】

次に、後段の演算部 11__2 が、前段の演算部 11__1 に受信完了通知を送付する。

【 0 1 8 4 】

最後に、前段の演算部 11__1 が後段の演算部 11__2 から受信完了通知を受信すると、前段の演算部 11__1 のトレースバッファ 18 のデータを消去する。

【 0 1 8 5 】

または、後段の演算部 11__2 が、ホスト部 12 に受信完了通知を送付して、ホスト部 12 の指示により前段の演算部 11__1 のトレースバッファ 18 のデータを消去してもよ

50

い。

【0186】

これにより、本実施の形態と同様の効果を奏する。

【0187】

また、本実施の形態において、トレースバッファ18より大きい記憶容量を有するホスト部12が、トレースバッファ18のデータを回収し記録する場合には、ホスト部12によるデータの回収または回収の完了を契機として、トレースバッファ18のデータを消去してもよい。

【0188】

また、本実施の形態において、予め設定した時間の経過を契機として、トレースバッファ18のデータを消去してもよい。

10

【0189】

<第3の実施の形態>

本発明の第3の実施の形態に係るコンピュータシステムおよびその制御方法を説明する。本実施の形態に係るコンピュータシステム10は、第1の実施の形態と同様の構成を有する。

【0190】

第1の実施の形態に係るコンピュータシステム10では、複数の演算部が連動して動作するので、複数の演算部が各々計測するタイムスタンプの時刻を同期させることが困難である。

20

【0191】

とくに、各演算部の動作周波数が異なる場合、動作周波数をベースとしたクロックカウンタの値が異なるため、タイムスタンプの時刻を同期させることが困難である。

【0192】

<コンピュータシステムの制御方法>

本実施の形態に係るコンピュータシステムおよびその制御方法では、ホスト部12が、タイムスタンプ値を所定の値、例えば初期値に一元的に設定する。図9に、本実施の形態に係るコンピュータシステムの制御方法のフローチャート図を示す。

【0193】

初めに、ホスト部12がタイムスタンプ値の初期値として所定の値を設定する(ステップS61)。

30

【0194】

ここで、演算部11_1~11_Nが入力データの計測を開始する時のタイムスタンプ値を初期値としてもよい。

【0195】

または、ホスト部12が、タイムスタンプ部17にカウントの開始を書き込み、カウントが開始される時のタイムスタンプ値を初期値としてもよい。

【0196】

次に、各演算部11_1~11_Nが、タイムスタンプ値を取得する(ステップS62)。

40

【0197】

次に、ホスト部12が、取得されたタイムスタンプ値と、タイムスタンプ値に対する所定の基準値とを比較する(ステップS63)。

【0198】

ここで、タイムスタンプ値に対する所定の基準値は、予め設定される。所定の基準値はタイムスタンプ値の所定の値、例えば初期値に対するずれの許容範囲を示すものである。

【0199】

最後に、各演算部11_1~11_Nのタイムスタンプ値が所定の基準値の許容範囲を超えた場合に、ホスト部12が、各演算部11_1~11_Nのタイムスタンプ値を所定の値、例えば初期値に再設定する(ステップS64)。

50

【0200】

本実施の形態に係るコンピュータシステムおよびその制御方法によれば、動作周波数の異なる演算部が共存する場合（例えば、演算部の演算内容によって動作周波数は異なる場合）に、タイムスタンプ値を統一的に扱えるので、トレースデータ回収後に容易に解析できる。

【0201】

また、複数の演算部を容易に同期させることが可能となるため、スケールアウト性を向上できる。

【0202】

また、ホスト部12が、計測開始などを契機として一元的にタイムスタンプ値を所定の値、例えば初期値に設定するため、演算部間のクロックカウンタのずれを容易に補正できる。

10

【0203】

また、本実施の形態は、当然、第1の実施の形態と同様の効果を奏する。

【0204】

<第3の実施の形態の変形例1>

本発明の第3の実施の形態の変形例1に係るコンピュータシステムおよびその制御方法を説明する。本変形例に係るコンピュータシステム10は、第1の実施の形態と同様の構成を有する。

【0205】

本変形例では、各演算部11__1～11__Nで周波数が異なる場合に、ホスト部12が各演算部11__1～11__Nの動作周波数の差を調整する。

20

【0206】

例えば、演算部11__1の周波数が100MHz、1クロックサイクルが10ナノ秒であり、演算部11__2の周波数が200MHz、1クロックサイクルが5ナノ秒である場合に、ホスト部12が演算部11__1のタイムスタンプ値を2倍にすることで、演算部11__1と演算部11__2のタイムスタンプ値を同期させる。

【0207】

また、ホスト部12が演算部11__2のタイムスタンプ値を1/2倍することで、演算部11__1と演算部11__2のタイムスタンプ値を同期させてもよい。

30

【0208】

また、ホスト部12がトレースデータを読み出した後に、換算基準値（例：100MHz）を設け、全演算部のカウンタ値を基準値100MHzに合わせて換算してもよい。

【0209】

このように、本変形例に係るコンピュータシステムおよびその制御方法では、ホスト部が、一の演算部（例えば、演算部11__1）の動作周波数に対して、他の演算部（例えば、演算部11__2）の動作周波数が同一になるように他の演算部（例えば、演算部11__2）に設定される係数を、他の演算部（例えば、演算部11__2）のタイムスタンプ値に乗ずる。ここで、他の演算部は複数であってもよい。このように、演算部ごとに異なるカウンタ値の差を調整する。

40

【0210】

これにより、各演算部のカウンタ値が同一になるので、本実施の形態と同様の効果を奏する。

【0211】

<第3の実施の形態の変形例2>

本発明の第3の実施の形態の変形例2に係るコンピュータシステムおよびその制御方法を説明する。本変形例に係るコンピュータシステム10は、第1の実施の形態と同様の構成を有する。

【0212】

本変形例では、各演算部11__1～11__Nで周波数が異なる場合は、各演算部11__

50

1 ~ 1 1 __ N で動作周波数の差を調整する。

【0 2 1 3】

まず、予め、演算器ごとに周波数に応じて異なる換算値（係数）を設定する。例えば、演算部 1 1 __ 1 の周波数が 1 0 0 M H Z、演算部 1 1 __ 2 の周波数が 2 0 0 M H Z である場合に、基準値を 1 0 0 M H Z とすれば、演算部 1 1 __ 1 の換算値（係数）は「1」、演算部 1 1 __ 2 の換算値（係数）は「1 / 2」となる。

【0 2 1 4】

次に、トレースバッファ 1 8 にタイムスタンプ値を記録する前段で、タイムスタンプ値に換算値（係数）を乗じた値を記録する。

【0 2 1 5】

また、トレースバッファ 1 8 にタイムスタンプ値を記録した後に、クロックカウンタ値を読み出す際に同様に換算をしてもよい。

【0 2 1 6】

このように、本変形例に係るコンピュータシステムおよびその制御方法では、演算部のトレースバッファの前段で、一の演算部（例えば、演算部 1 1 __ 1）の動作周波数に対して、他の演算部（例えば、演算部 1 1 __ 2）の動作周波数が同一になるように他の演算部（例えば、演算部 1 1 __ 2）に設定される係数を、他の演算部（例えば、演算部 1 1 __ 2）のタイムスタンプ値に乘ずる。ここで、他の演算部は複数であってもよい。このように、演算部ごとに異なるカウンタ値の差を調整する。

【0 2 1 7】

これにより、各演算部のカウンタ値が同一になるので、本実施の形態と同様の効果を奏する。

【0 2 1 8】

本発明の実施の形態で、演算部が処理時間、データ流量等の計測、不具合などの判定などを実行する場合、演算部における演算器が実行してもよいし、演算部内に別途計測、判定などの処理機能を設けてもよい。

【0 2 1 9】

本発明の実施の形態では、第 1 の実施の形態と、第 2 の実施の形態と、第 3 の実施の形態とを組み合わせることにより、さらなる効果を奏することができる。

【0 2 2 0】

本発明の実施の形態では、コンピュータシステムの構成、管理方法などにおいて、各構成部の構造、寸法、材料等の一例を示したが、これに限らない。コンピュータシステムの機能を発揮し効果を奏するものであればよい。

【産業上の利用可能性】

【0 2 2 1】

本発明は、情報処理分野におけるコンピュータシステムに適用することができる。

【符号の説明】

【0 2 2 2】

1 0 コンピュータシステム

1 1 __ 1 ~ 1 1 __ N 演算部

1 2 ホスト部

1 3 内部通信部

1 4 __ 1 ~ 1 4 __ N トレース部

10

20

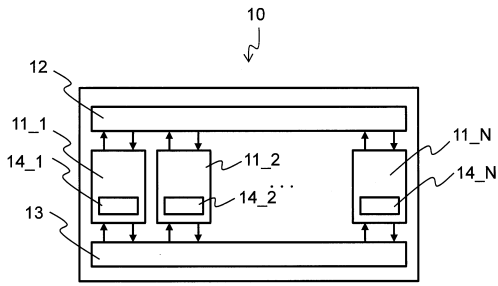
30

40

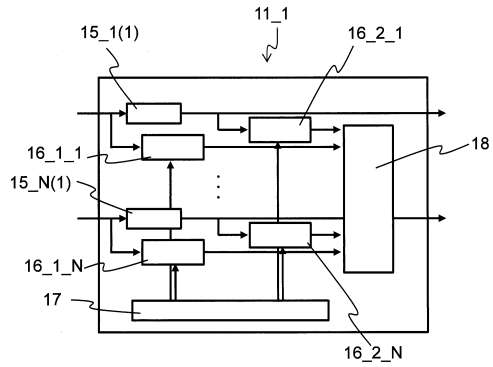
50

【図面】

【図 1】

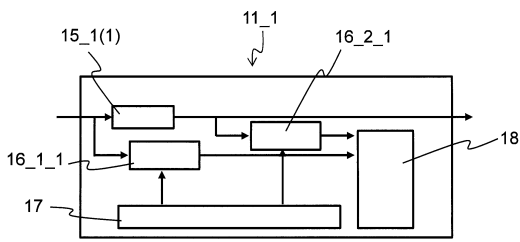


【図 2】

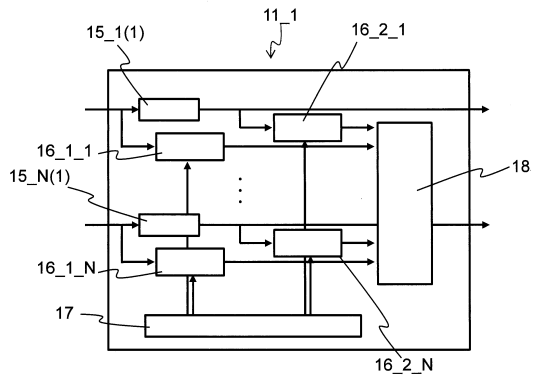


10

【図 3 A】

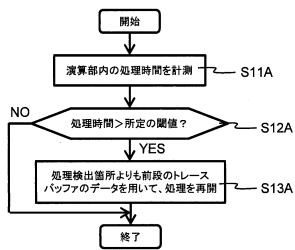


【図 3 B】

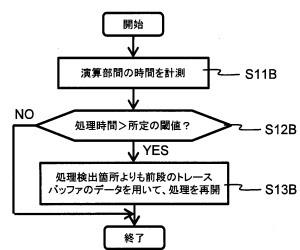


20

【図 4 A】



【図 4 B】

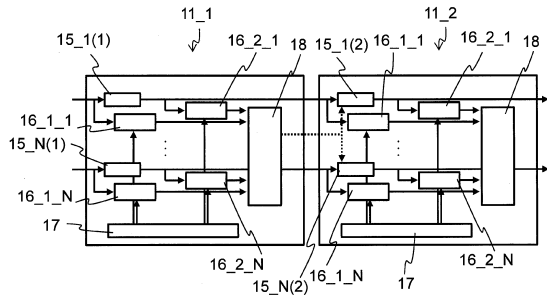


30

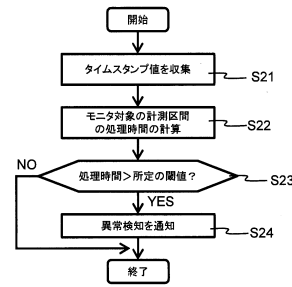
40

50

【図 4 C】

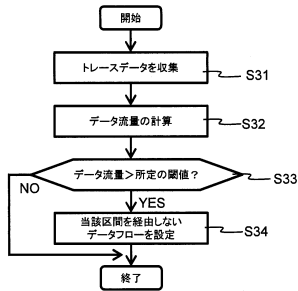


【図 5】



10

【図 6】



【図 7 A】

Timestamp:Dec,	Timestamp:0x,	Ins,	Evt,	Dec,	TID	EventData
6983,	1b47,	f0,	07,	HLR-	00000000,	00000000
6996,	1b54,	f1,	07,	HLR-	00000000,	00000000
401563,	6209d,	f0,	05,	H-R,	00000001,	00000000
401564,	6209c,	10,	05,	H-R,	00000001,	00000000
401656,	620f3,	11,	05,	H-R,	00000001,	00000000
402995,	62633,	f1,	05,	H-R,	00000000,	00000000
404019,	62a33,	f1,	06,	-LR-	00000000,	00000000
418315,	6620b,	f1,	05,	H-R,	00000000,	00000000
419339,	6660b,	f1,	06,	-LR-	00000000,	00000000
428625,	68a51,	f1,	05,	H-R,	00000000,	00000000
429649,	68e51,	f1,	06,	-LR-	00000000,	00000000
438909,	6b27d,	f1,	05,	H-R,	00000000,	00000000
439933,	6b67d,	f1,	06,	-LR-	00000000,	00000000
448499,	6d7f3,	f1,	05,	H-R,	00000000,	00000000
449523,	6dbf3,	f1,	06,	-LR-	00000000,	00000000
457367,	6fa97,	f1,	05,	H-R,	00000000,	00000000
458391,	6fe97,	f1,	06,	-LR-	00000000,	00000000
466369,	71dc1,	f1,	05,	H-R,	00000000,	00000000
467393,	721c1,	f1,	06,	-LR-	00000000,	00000000
475235,	74063,	f1,	05,	H-R,	00000000,	00000000
476259,	74463,	f1,	06,	-LR-	00000000,	00000000
484240,	76390,	f1,	05,	H-R,	00000000,	00000000
485264,	76790,	f1,	06,	-LR-	00000000,	00000000
493267,	786d3,	f1,	05,	H-R,	00000000,	00000000
494291,	78ad3,	f1,	06,	-LR-	00000000,	00000000
502432,	7aAa0,	f1,	05,	H-R,	00000000,	00000000
503456,	7aea0,	f1,	06,	-LR-	00000000,	00000000
511523,	7ce23,	f1,	05,	H-R,	00000000,	00000000
512547,	7d223,	f1,	06,	-LR-	00000000,	00000000
520549,	7f165,	f1,	05,	H-R,	00000000,	00000000
521573,	7f565,	f1,	06,	-LR-	00000000,	00000000
529688,	81518,	f1,	05,	H-R,	00000000,	00000000
530712,	81918,	f1,	06,	-LR-	00000000,	00000000
538726,	83866,	f1,	05,	H-R,	00000000,	00000000
539750,	83c66,	f1,	06,	-LR-	00000000,	00000000
547790,	85bce,	f0,	06,	-LR-	00000001,	00000000
547791,	85bcf,	10,	06,	-LR-	00000001,	00000000
547794,	85bd2,	11,	06,	-LR-	00000001,	00000000
547796,	85bd4,	f1,	05,	H-R,	00000000,	00000000
548820,	85fd4,	f1,	06,	-LR-	00000000,	00000000

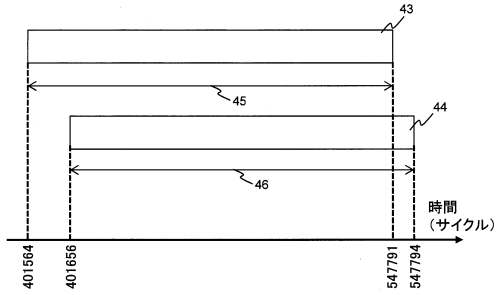
20

30

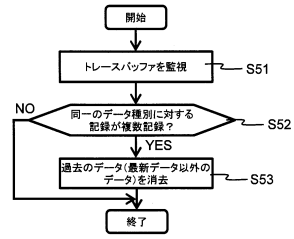
40

50

【 図 7 B 】

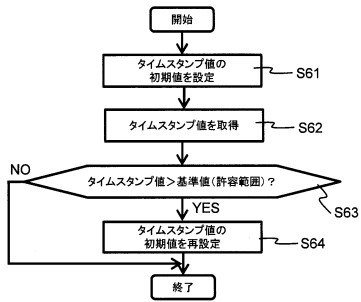


【 図 8 】



10

【 図 9 】



20

30

40

50

 フロントページの続き

- 東京都千代田区大手町一丁目 5 番 1 号 日本電信電話株式会社内
- (72)発明者 田仲 顕至
- 東京都千代田区大手町一丁目 5 番 1 号 日本電信電話株式会社内
- (72)発明者 伊藤 猛
- 東京都千代田区大手町一丁目 5 番 1 号 日本電信電話株式会社内
- (72)発明者 坂本 健
- 東京都千代田区大手町一丁目 5 番 1 号 日本電信電話株式会社内
- (72)発明者 村中 勇介
- 東京都千代田区大手町一丁目 5 番 1 号 日本電信電話株式会社内
- 審査官 真木 健彦
- (56)参考文献 国際公開第 2018/182782 (WO, A1)
 米国特許出願公開第 2012/0331354 (US, A1)
 特開 2012-252636 (JP, A)
 田仲 顕至 ほか, 分散深層学習を高速化させる FPGA Ring - Allreduce の
 検討, 第 82 回 (2020 年) 全国大会講演論文集 (1) コンピュータシステム ソフト
 ウェア科学・工学 データとウェブ, 日本, 一般社団法人情報処理学会, 2020 年 02 月 20 日
 , 7A-01, P. 1-31 ~ 1-32
- (58)調査した分野 (Int.Cl., DB 名)
- G 0 6 F 1 1 / 2 8 - 1 1 / 3 6
 G 0 6 F 1 1 / 0 7
 G 0 6 F 1 5 / 1 6 - 1 5 / 1 7 7
 G 0 6 F 1 5 / 8 2