

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4687925号
(P4687925)

(45) 発行日 平成23年5月25日 (2011.5.25)

(24) 登録日 平成23年2月25日 (2011.2.25)

(51) Int. Cl. F I
HO 4 L 12/56 (2006.01) HO 4 L 12/56 2 0 0 Z
GO 6 F 15/173 (2006.01) GO 6 F 15/173 6 5 0 S

請求項の数 10 (全 14 頁)

(21) 出願番号	特願2008-31545 (P2008-31545)	(73) 特許権者	000168285 エヌイーシーコンピュータテクノ株式会社 山梨県甲府市大津町1088-3
(22) 出願日	平成20年2月13日 (2008.2.13)	(74) 代理人	100102864 弁理士 工藤 実
(65) 公開番号	特開2009-194510 (P2009-194510A)	(72) 発明者	曾田 泰広 山梨県甲府市大津町1088-3 エヌイーシーコンピュータテクノ株式会社内
(43) 公開日	平成21年8月27日 (2009.8.27)		
審査請求日	平成21年3月12日 (2009.3.12)	審査官	安藤 一道

最終頁に続く

(54) 【発明の名称】 優先調停システム及び優先調停方法

(57) 【特許請求の範囲】

【請求項1】

複数のCPUと、
 前記複数のCPUからアクセスされる複数のハードウェア資源と、
 前記複数のCPUの各々に対応する記憶領域に設定され、自CPUと前記複数のハードウェア資源の各々とのレイテンシを示すレイテンシ情報を格納する複数のルーティングテーブルと、
 前記複数のCPUと前記複数のハードウェア資源との間で交換されるパケットのルーティングを行う複数のクロスバとを具備し、
 前記複数のCPUの各々は、送信パケットを前記複数のハードウェア資源のうちの送信先ハードウェア資源に送信するとき、前記複数のルーティングテーブルのうち自CPUに対応するルーティングテーブルにおいて前記送信先ハードウェア資源に対応する前記レイテンシ情報を前記送信パケットに付加し、
 前記複数のクロスバの各々は、複数のパケットを受信したとき、受信した複数のパケットの各々の前記レイテンシ情報を参照してレイテンシが大きいパケットを優先的に前記送信先ハードウェア資源の方に送信する優先調停処理を実行する
 優先調停システム。

10

【請求項2】

請求項1に記載された優先調停システムであって、
 前記レイテンシ情報は前記複数のCPUの各々から前記複数のハードウェア資源の各々

20

への片道レイテンシを示し、

前記複数のハードウェア資源の各々は、返信パケットを前記複数のCPUのうちの送信先CPUに返信するとき、前記複数のCPUの各々への片道レイテンシを示す復路レイテンシ情報を格納するハードウェア資源側ルーティングテーブルを参照して、前記送信先CPUに対応する前記復路レイテンシ情報を前記返信パケットに付加し、

前記複数のクロスバの各々は、前記複数のハードウェア資源の側から複数のパケットを受信したとき、受信した複数のパケットの各々の前記復路レイテンシ情報を参照して、レイテンシが大きいパケットを優先的に前記送信先CPUの方に送信する

優先調停システム。

【請求項3】

請求項1または2に記載された優先調停システムであって、

前記複数のルーティングテーブルの各々は、前記複数のCPUのうちの自ルーティングテーブルに対応するCPUと前記複数のハードウェア資源の各々との間でのパケット転送のルートを示すルーティング情報を格納する

優先調停システム。

【請求項4】

請求項3に記載された優先調停システムであって、

前記ルーティング情報は前記複数のCPUの各々から前記複数のハードウェア資源の各々へのルートを示し、

前記複数のハードウェア資源側ルーティングテーブルの各々は、前記複数のハードウェア資源のうちの自ハードウェア資源側ルーティングテーブルに対応するハードウェア資源から前記複数のCPU資源の各々へのパケット転送のルートを示す復路ルーティング情報を格納する

優先調停システム。

【請求項5】

請求項1から4のいずれかに記載された優先調停システムであって、

前記複数のクロスバの各々は、受信したパケットの前記レイテンシ情報を、前記優先調停処理による待ち時間を加算した値に更新する

優先調停システム。

【請求項6】

複数のCPUの各々に対応する記憶領域に、自CPUと複数のハードウェア資源の各々とのレイテンシを示すレイテンシ情報を格納する複数のルーティングテーブルを設定するステップと、

前記複数のCPUの各々が、送信パケットを前記複数のハードウェア資源のうちの送信先ハードウェア資源に送信するとき、前記複数のルーティングテーブルのうち自CPUに対応するルーティングテーブルにおいて前記送信先ハードウェア資源に対応する前記レイテンシ情報を前記送信パケットに付加するステップと、

前記複数のクロスバの各々が、複数のパケットを受信したとき、受信した複数のパケットの各々の前記レイテンシ情報を参照してレイテンシが大きいパケットを優先的に前記送信先ハードウェア資源の方に送信する優先調停処理を実行するステップ

とを具備する優先調停方法。

【請求項7】

請求項6に記載された優先調停方法であって、

前記レイテンシ情報は前記複数のCPUの各々から前記複数のハードウェア資源の各々への片道レイテンシを示し、

更に、前記複数のハードウェア資源の各々が、返信パケットを前記複数のCPUのうちの送信先CPUに返信するとき、前記複数のCPUの各々への片道レイテンシを示す復路レイテンシ情報を格納するハードウェア資源側ルーティングテーブルを参照して、前記送信先CPUに対応する前記復路レイテンシ情報を前記返信パケットに付加するステップと

、

10

20

30

40

50

前記複数のクロスバの各々が、前記複数のハードウェア資源の側から複数のパケットを受信したとき、受信した複数のパケットの各々の前記復路レイテンシ情報を参照して、レイテンシが大きいパケットを優先的に前記送信先CPUの方に送信するステップとを具備する優先調停方法。

【請求項 8】

請求項 6 または 7 に記載された優先調停方法であって、

前記複数のルーティングテーブルの各々は、前記複数のCPUのうちの自ルーティングテーブルに対応するCPUと前記複数のハードウェア資源の各々との間でのパケット転送のルートを示すルーティング情報を格納し、

更に、前記複数のCPUの各々が、前記複数のルーティングテーブルのうち自CPUに対応するルーティングテーブルにおいて前記送信先ハードウェア資源に対応する前記ルーティング情報を前記送信パケットに付加するステップと、

前記複数のクロスバの各々が、受信したパケットを前記ルーティング情報に設定されたルーティング先に転送するステップ

とを具備する優先調停方法。

【請求項 9】

請求項 8 に記載された優先調停方法であって、

前記ルーティング情報は前記複数のCPUの各々から前記複数のハードウェア資源の各々へのルートを示し、

前記複数のハードウェア資源側ルーティングテーブルの各々は、前記複数のハードウェア資源のうちの自ハードウェア資源側ルーティングテーブルに対応するハードウェア資源から前記複数のCPU資源の各々へのパケット転送のルートを示す復路ルーティング情報を格納し、

更に、前記複数のハードウェア資源の各々が、前記複数の復路ルーティングテーブルのうち自ハードウェア資源に対応する復路ルーティングテーブルにおいて送信先CPUに対応する前記復路ルーティング情報を前記返信パケットに付加するステップと、

前記複数のクロスバの各々が、受信した前記返信パケットを前記復路ルーティング情報に設定されたルーティング先に転送するステップ

とを具備する優先調停方法。

【請求項 10】

請求項 6 から 9 のいずれかに記載された優先調停方法であって、

更に、前記複数のクロスバの各々が、受信したパケットのレイテンシ情報を、前記優先調停処理による待ち時間を加算した値に更新する

優先調停方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、パケットの調停に関する。

【背景技術】

【0002】

複数のCPUからアクセス先の共有資源（メモリ、IO等）までの経路上に、調停回路を有する複数のクロスバが存在する大規模SMP（Symmetric Multi Processing）コンピュータシステムが知られている。

【0003】

以下に、出願人が知り得た先行技術文献を記載する。

【特許文献 1】特開 2005 - 340959 号公報

【特許文献 2】特開平 9 - 321795 号公報

【特許文献 3】特開平 10 - 135952 号公報

【特許文献 4】特開平 11 - 312093 号公報

【発明の開示】

【発明が解決しようとする課題】**【0004】**

物理的に経路の遠くに位置する共有資源へのアクセスは、複数のLSIを経由して行われる。そのため、レイテンシが長くなってしまふという欠点が存在する。さらに、経路途中の調停回路において競合が発生した場合、該リクエストが待たされることがある。この場合、該リクエストの完了までにかかるレイテンシは更に増大する。

【0005】

仮に、CPU内に、経路の遠くに位置する共有資源にアクセスを要求する先行リクエストと、該先行リクエストの完了を待って発行される後続リクエストが存在した場合、先行リクエストのレイテンシの増大は、後続リクエストの発行を遅らせる原因となる。その結果、コンピュータシステム全体の性能を低下させてしまふという問題が存在した。

10

【課題を解決するための手段】**【0006】**

本発明による優先調停システムは、複数のCPUと、複数のCPUからアクセスされる複数のハードウェア資源と、複数のCPUの各々に対応する記憶領域に設定され、自CPUと複数のハードウェア資源の各々とのレイテンシを示すレイテンシ情報を格納する複数のルーティングテーブルと、複数のCPUと複数のハードウェア資源との間で交換されるパケットのルーティングを行う複数のクロスバとを備える。複数のCPUの各々は、送信パケットを複数のハードウェア資源のうちの送信先ハードウェア資源に送信するとき、複数のルーティングテーブルのうち自CPUに対応するルーティングテーブルにおいて送信先ハードウェア資源に対応するレイテンシ情報を送信パケットに付加する。複数のクロスバの各々は、複数のパケットを受信したとき、受信した複数のパケットの各々のレイテンシ情報を参照してレイテンシが大きいパケットを優先的に送信先ハードウェア資源の方に送信する優先調停処理を実行する。

20

【0007】

本発明による優先調停方法は、複数のCPUの各々に対応する記憶領域に、自CPUと複数のハードウェア資源の各々とのレイテンシを示すレイテンシ情報を格納する複数のルーティングテーブルを設定するステップと、複数のCPUの各々が、送信パケットを複数のハードウェア資源のうちの送信先ハードウェア資源に送信するとき、複数のルーティングテーブルのうち自CPUに対応するルーティングテーブルにおいて送信先ハードウェア資源に対応するレイテンシ情報を送信パケットに付加するステップと、複数のクロスバの各々が、複数のパケットを受信したとき、受信した複数のパケットの各々のレイテンシ情報を参照してレイテンシが大きいパケットを優先的に送信先ハードウェア資源の方に送信する優先調停処理を実行するステップとを備える。

30

【発明の効果】**【0008】**

本発明により、長経路のパケットや経路上の競合により待たされたパケットによるレイテンシの低下を抑制する優先調停装置及び優先調停方法が実現される。

【発明を実施するための最良の形態】**【0009】**

以下、本発明を実施するための最良の形態について説明する。最初に、本実施の形態について概略的に説明する。

40

【0010】

本実施の形態の優先調停システムと優先調停方法は、リクエスタ(CPU)からアクセス先の共有資源(メモリ、IO等のハードウェア資源)までの経路上に、調停回路を有する複数のクロスバが存在する大規模SMPコンピュータシステムに適用される。

【0011】

リクエスタが発行するリクエストパケットのヘッダ部分に、レイテンシ値を記載するフィールド(レイテンシフィールド)が設けられる。アクセスリクエストの行き先情報から、行き先の共有資源までのルーティング情報と最短所要時間(最少レイテンシ)を求める

50

ことが出来るルーティングテーブルが設けられる。CPUは、このルーティングテーブルを参照して、パケットを生成する際にパケットヘッダのレイテンシフィールドに該レイテンシテーブルから求めた経路通過で所要されるレイテンシ値を設定する。

【0012】

パケットの転送経路上で複数のパケットが同一経路を同時に使用することにより競合が発生し、調停が必要となる場合がある。こうした場合、パケットヘッダのレイテンシ情報を比較して、パケットが所用するレイテンシ値が大きなパケットを優先的に通過させる機構が設けられる。この機構は、長経路のパケットのレイテンシを改善する。

【0013】

また、或るパケットが調停競合で一時的に待たされた際に、待たされた時間をパケットヘッダのレイテンシ値に加算する機構が設けられる。この機構は、該パケットを次の調停においては以前より高優先なパケットとすることを可能とする。その結果、該パケットの次調停競合における更なるレイテンシ悪化が抑制される。

【0014】

次に、こうした優先調停システムと優先調停方法を実現するための具体的な構成について説明する。図1に本実施の形態の構成を示す。図1を参照すると、本実施の形態における優先調停システム、優先調停方法が適用されるシステムは複数のLSIで構成される。このシステムは、リクエストパケットを送出する機能を有する4個のリクエストCPU1010、CPU1011、CPU1110、CPU1111と、リクエストのアクセス先となる4個の共有資源IOH1220、IOH1221、IOH1320、IOH1321と、リクエストCPU1010、CPU1011、CPU1110、CPU1111と共有資源IOH1220、IOH1221、IOH1320、IOH1321との間に位置し、パケットの通過する経路を形成するクロスバNC1000、NC1100、NC1200、NC1300とを備える。

【0015】

これらのLSIのうちの任意の2LSI間は、一方から他方にデータを転送する単方向の接続インタフェースと、他方から一方にデータを転送する単方向の接続インタフェースとの2本を束ねた双方向の接続インタフェースを介して接続される。

【0016】

リクエストCPU1010およびCPU1011は双方向の接続インタフェースを介してクロスバNC1000と接続し、共有資源IOH1220～IOH1321に対するアクセスリクエストをリクエストパケットとしてクロスバNC1000へ出力する機能および、共有資源から処理完了通知として送られるリプライパケットを受信する機能を有している。

【0017】

同様に、リクエストCPU1110およびCPU1111は双方向の接続インタフェースを介してクロスバNC1100と接続し、共有資源IOH1220～IOH1321へのアクセスリクエストパケットをクロスバNC1100へ出力する機能および、共有資源からのリプライパケットを受信する機能を有している。

【0018】

また、各リクエストCPU1010、CPU1011、CPU1110、CPU1111に対応する記憶領域には、自CPUからアクセスし得る共有資源毎に、パケット転送時にクロスバで使用されるルーティング情報と、自CPUから行き先共有資源までパケットを転送する際に最低必要な所要時間を表す最少レイテンシ情報(本実施の形態においては片道レイテンシを示す)とを保持するルーティングテーブルR1010、R1011、R1110、R1111が格納される。各リクエストCPU1010、CPU1011、CPU1110、CPU1111は、リクエストパケット出力時に、自CPUのルーティングテーブルを参照して、行き先の共有資源に対応するルーティング情報と最短レイテンシ値をパケットヘッダの行き先指定フィールドおよびレイテンシフィールドに設定する機能を有している。

10

20

30

40

50

【 0 0 1 9 】

図4はルーティングテーブルR1010、R1011、R1110、R1111のうちの任意の一つのルーティングテーブルRの構成および、パケットヘッダへの設定動作を示している。ルーティングテーブルR内には共有資源毎にルーティング情報およびレイテンシ情報（最短レイテンシ値）が保持されている。各リクエストCPU1010、CPU1011、CPU1110、CPU1111は、行き先となる共有資源を指定することで、対応するルーティング情報および最短レイテンシ値を取り出し、それらの値をパケットヘッダの行き先指定フィールドF1とレイテンシフィールドF2にそれぞれ設定できる。

【 0 0 2 0 】

共有資源IOH1220、IOH1221は、クロスバNC1200と双方向の接続インタフェースを介して接続し、全てのリクエストCPUからのリクエストパケットをクロスバNC1200経由で受信する機能を有している。同様に共有資源IOH1320、IOH1321も、クロスバNC1300と双方向の接続インタフェースを介して接続し、全てのリクエストCPUからのリクエストパケットをクロスバNC1300経由で受信する機能を有している。

10

【 0 0 2 1 】

リクエストパケットを受信した共有資源IOH1220～1321は、リクエストCPUからのリクエストパケットを解析し、適当な処理（共有資源に対するリードやライト等）を行う。その後、共有資源IOH1220～1321は、リクエストCPUを送信先CPUとして対するリプライパケットを生成し、接続するクロスバへ送出する。

20

【 0 0 2 2 】

共有資源IOH1220～1321は、リクエストCPUが有するルーティングテーブルと同等機能のハードウェア資源側ルーティングテーブルR1220、1221、1320、1321を内部に有し、リクエストCPU1010～1111によるリクエストパケット生成と同様の機能を有する。すなわち、共有資源IOH1220～1321は、リプライパケット出力時に、リクエスト発行元のリクエストCPUを送信先CPUとする。共有資源IOH1220～1321は、送信先CPUに対応するルーティング情報である復路ルーティング情報と最短レイテンシ値を示す復路レイテンシ情報を返信パケットのパケットヘッダの行き先指定フィールドおよびレイテンシフィールドに設定する。

30

【 0 0 2 3 】

クロスバNC1000は隣接するクロスバNC1100、NC1200および、リクエストCPU1010、CPU1011と双方向の接続インタフェースを介して接続し、LSI間でパケットの入出力を行う機能を有している。クロスバNC1000内の各接続インタフェースにはポート番号が割り当てられている。クロスバNC1100との接続インタフェースはポート2、クロスバNC1200との接続インタフェースはポート3、リクエストCPU1010との接続インタフェースはポート0、リクエストCPU1011との接続インタフェースはポート1と割りあてられる。

【 0 0 2 4 】

同様に、クロスバNC1100は隣接するクロスバNC1000、NC1300および、リクエストCPU1110、CPU1111と双方向の接続インタフェースを介して接続し、LSI間でのパケットの入出力を行う機能を有している。クロスバNC1100内の各接続インタフェースは、クロスバNC1000と同様にポート番号が割り当てられている。クロスバNC1000との接続インタフェースはポート3、クロスバNC1300との接続インタフェースはポート2、リクエストCPU1110との接続インタフェースはポート0、リクエストCPU1111との接続インタフェースはポート1と割りあてられる。

40

【 0 0 2 5 】

前述のクロスバと同様に、クロスバNC1200は隣接するクロスバNC1000、NC1300および、共有資源IOH1220、IOH1221と双方向の接続インタフェースを介して接続し、LSI間でのパケットの入出力を行う機能を有している。クロスバ

50

NC1200内の各接続インタフェースは、前述のクロスバと同様にポート番号が割り当てられている。クロスバNC1000との接続インタフェースはポート2、クロスバNC1300との接続インタフェースはポート3、共資源IOH1220との接続インタフェースはポート1、共有資源IOH1221との接続インタフェースはポート0と割り当てられる。

【0026】

前述のクロスバと同様に、クロスバNC1300は隣接するクロスバNC1100、NC1200および、共有資源IOH1320、IOH1321と双方向の接続インタフェースを介して接続し、LSI間でのパケットの入出力を行う機能を有している。クロスバNC1300内の各接続インタフェースは、前述のクロスバと同様にポート番号が割り当てられている。クロスバNC1200との接続インタフェースはポート2、クロスバNC1100との接続インタフェースはポート3、共資源IOH1320との接続インタフェースはポート1、共有資源IOH1321との接続インタフェースはポート0と割り当てられる。

10

【0027】

各クロスバNC1000～NC1300はリクエストパケットとリプライパケットの区別は行わず、パケットヘッダに存在するルーティング情報に従い、パケットを行き先IOH1220～1321またはCPU1010～CPU1111へ転送することを目的として動作する。各クロスバNC1000～NC1300は内部にパケット受信部および出力調停部を接続インタフェース毎に有する。各クロスバNC1000～NC1300が接続

20

【0028】

図2はクロスバ内部のパケット受信部および出力調停部の動作を示している。本実施の形態の構成である図1のクロスバ構成では、クロスバ1個あたりポートが4個存在する。そのため、実際は4個のパケット受信部と4個の出力調停部の図が適当ではあるが、図2ではこれを簡略化しており、2個のパケット受信部と1個の出力調停部について表している。

【0029】

各パケット受信部は、受信バッファと行き先解析回路C1を有している。受信バッファは、受信パケットを一旦保持する。行き先解析回路C1は、受信バッファに保持されたパケットのパケットヘッダの行き先指定フィールドを解析して、パケットの出力先ポートを決定し、該出力先ポートの出力調停回路C2に対して調停リクエストを出力する。また、各パケット受信部は、行き先解析回路C1で決定した行き先ポートの出力調停回路C2に対して、調停リクエストを出力すると同時に、パケットヘッダのレイテンシフィールドの値を出力する機能も有している。

30

【0030】

各パケット受信部は、出力調停部から使用許可信号(GNT信号)を受信したら、受信バッファからパケットを取出し、出力ポートへパケットを送出する機能を有している。

【0031】

各パケット受信部は、以下に説明するように、パケット毎に調停の結果に応じてレイテンシ値を書き換える機能を有する。各パケット受信部は、出力先ポートの出力調停回路C2に調停リクエストを出力しても、直ぐにGNT信号が出力されない場合は、受信バッファからパケットを取り出すことなく、GNT信号を受信するまで待ち状態(調停待ち状態)となる。各パケット受信部は、この間、パケットヘッダのレイテンシフィールド値に待ち時間を加算し、レイテンシフィールドを更新する。すなわち、パケット受信部が調停待ち状態となった場合、その間、ヘッダのレイテンシフィールドを加算し続け、該値をレイテンシフィールド値として出力調停部へ出力し続ける。

40

【0032】

クロスバ内の出力調停部は、調停回路C2を有している。調停回路C2は、クロスバ内

50

の各ポートの packets 受信部から調停リクエストを受信した際に、それら調停リクエストから使用許可ポートを一つだけに決定し、該ポートに対し出力ポートの使用許可通知として GNT 信号を出力する。調停回路 C2 は、各 packets 受信部から出力されている packets ヘッダ内のレイテンシフィールド値に基づいて出力ポートの使用を許可する優先調停機能を有している。この機能により、使用許可ポートを決定する際に、レイテンシ値が大きな調停リクエストに対して優先的に出力ポートの使用が許可される。

【0033】

また、クロスバ内の出力調停部は、調停回路 C2 が決定した使用許可ポートから packets を受信し、接続インタフェースに送出する機能を有する packets 送信部を有している。

【0034】

次に、以上の構成を備えたシステムの動作について説明する。図3は本実施の形態の動作例を示している。図3は図1の構成とほぼ同じ構成で、各要素の番号も千番台が三千番台になっているだけである。動作の説明に登場しない部分は削除した。以降では図3を参照し、長経路の packets 転送における優先調停の動作および、調停競合で遅れた packets に対する優先調停の動作の一例を説明する。

【0035】

[動作例の説明における前提]

本動作例では、リクエスタ CPU3010 から送信先のハードウェア資源である共有資源 IOH3320 までの経路は、クロスバ3000、クロスバ3200 およびクロスバ3300 経由とする(経路1)。

本動作例では、リクエスタ CPU3011 から共有資源 IOH3220 までの経路は、クロスバ3000 およびクロスバ3200 経由とする(経路2)。

本動作例では、リクエスタ CPU3110 から共有資源 IOH3220 までの経路は、クロスバ3100、クロスバ3300 およびクロスバ3200 経由とする(経路3)。

【0036】

経路の最短レイテンシ値は、本実施の形態の動作例においてはクロスバ1個当たりの通過レイテンシを10と設定される。即ち、経路1および経路3の最短レイテンシ値は30、経路2の最短レイテンシは20である。

また、本実施の形態の動作例において、調停待ち時間は一律15と設定される。

尚、packets ヘッダに設定するルーティング情報については、本実施の形態の動作説明では詳しく言及しないが、一般的なルーティング情報相当のものである。ルーティング情報は、クロスバ内で解析することにより、出力先のポートが判断できる情報を含んでいるものとして扱う。

【0037】

[長経路の packets 転送における優先調停]

リクエスタ CPU3010 は共有資源 IOH3220 へのリクエスト packets を生成し、クロスバ NC3000 へ送出する。この際、リクエスタ CPU3010 は、共有資源 IOH03220 までのルーティング情報および最短レイテンシ値(=30)を自 CPU のルーティングテーブルから求め、packets のヘッダの行き先指定フィールドと、レイテンシフィールドに設定する。リクエスタ CPU2010 からのリクエスト packets はクロスバ NC3000 のポート0の packets 受信部に受信される。クロスバ NC3000 のポート0の packets 受信部は、受信した packets の packets ヘッダ内の行き先指定フィールドを参照し、出力先はポート3と判断し、クロスバ NC3000 のポート3の出力調停部に対して調停リクエストを出力する。また、これと同時にリクエスタ CPU3010 からのリクエスト packets のレイテンシフィールド値(=30)もクロスバ NC3000 のポート3の出力調停部に対して出力する。

【0038】

クロスバ NC3000 のポート0の packets 受信部がリクエスタ CPU3010 からのリクエスト packets を受信すると同時に、クロスバ NC3000 のポート1の packets 受信部は、リクエスタ CPU3011 から共有資源 IOH3220 行きのリクエスト packets

10

20

30

40

50

トを受信する。クロスバNC3000のポート1の packets 受信部は、受信した packets の packets ヘッダ内の行き先指定フィールドを参照し、出力先はポート3と判断し、クロスバNC3000のポート3の出力調停部に対して調停リクエストを出力する。また、これと同時にリクエストCPU3011からのリクエスト packets のレイテンシフィールド値(=20)もクロスバNC3000のポート3の出力調停部に対して出力する。

【0039】

クロスバNC3000のポート3の出力調停部は、以下のように優先調停処理を実行する。出力調停部は、クロスバNC3000のポート0の packets 受信部とクロスバNC3000のポート1の packets 受信部から同時に調停リクエストを受信する。この時、クロスバNC3000のポート3の出力調停部はポート0とポート1から受信するレイテンシフィールド値(ポート0:ポート1=30:20)を比較し、より値の大きいポート0側のリクエストを優先的に通過させるリクエストであると判断して、クロスバNC3000のポート0に対してGNT信号を出力する。

10

【0040】

クロスバNC3000のポート3の出力調停部からGNT信号を受信したクロスバNC3000のポート0の packets 受信部は、該 packets をクロスバNC3000のポート3に送出する。この際、調停による待ち時間は無いため、 packets ヘッダのレイテンシフィールドの加算は行われない。

【0041】

クロスバNC3000のポート3から出力されたリクエストCPU3010からのリクエスト packets は、クロスバNC3200のポート2の packets 受信部に受信され、クロスバNC3200のポート3の出力調停部へ送出される。ここでは他の packets との競合は無いため、 packets は調停無しで通過する。調停で待たされることが無いため、 packets ヘッダのレイテンシフィールド値は加算されない。

20

【0042】

クロスバNC3200のポート3から送出されたリクエストCPU3010からのリクエスト packets は、クロスバNC3300のポート2の packets 受信部に受信され、クロスバNC3300のポート1の出力調停部へ出力される。ここでも他の packets との競合は無いため、 packets は調停無しで通過する。調停で待たされることが無いため、 packets のヘッダのレイテンシフィールド値は加算されない。

30

【0043】

クロスバNC3300のポート1から共有資源IOH3320に到達したリクエストCPU3010からのリクエスト packets の転送経路における総レイテンシは、該 packets ヘッダのレイテンシフィールドの値と一致する。そのため、本動作例の場合、該 packets は最短レイテンシ(レイテンシフィールドの値=30)で共有資源IOH3320への経路を通過する。

【0044】

本動作例では、リクエストCPU3010からのリクエスト packets の転送経路において、出力調停部での競合が1回しかなかった。しかし、この経路上の全てのクロスバ内で競合が起きた場合、計3回、調停競合が発生し、各調停部においてレイテンシ値の比較による優先調停が行われる。この場合、競合相手の packets が一度も調停待ちを行っていない packets であるならば、本構成例において最長経路であるクエストCPU3010からのリクエスト packets より優先な packets は存在しない。そのため、競合があったケースにおいても、該 packets は最短レイテンシで通過できる可能性が高い。

40

【0045】

[調停競合で遅れた packets の優先調停]

リクエストCPU3010は共有資源IOH3220へのリクエスト packets を生成し、クロスバNC3000へ送出する。この際、リクエストCPU3010は、共有資源IOH03220までのルーティング情報および最短レイテンシ値をルーティングテーブルRから求め、 packets のヘッダの行き先指定フィールドF1と、レイテンシフィールドF

50

2 に設定する。

【 0 0 4 6 】

リクエストCPU3011からのリクエストパケットはクロスバNC3000のポート1のパケット受信部に受信される。クロスバNC3000のポート1のパケット受信部は、受信したパケットのパケットヘッダ内の行き先指定フィールドF1を参照し、出力先はポート3と判断し、クロスバNC3000のポート3の出力調停部に対して調停リクエストを出力する。これと同時にリクエストCPU3011からのリクエストパケットのレイテンシフィールドF2の値 = 20もクロスバNC3000のポート3の出力調停部に対して出力する。

【 0 0 4 7 】

クロスバNC3000のポート1のパケット受信部がリクエストCPU3011からのリクエストパケットを受信すると同時に、クロスバNC3000のポート0のパケット受信部は、リクエストCPU3010から共有資源IOH3320行きのリクエストパケット(レイテンシフィールド値 = 30)を受信する。クロスバNC3000のポート0のパケット受信部は、ポート1のパケット受信部と同じく、クロスバNC3000のポート3の出力調停部に対して調停リクエストを出力する。

【 0 0 4 8 】

クロスバNC3000のポート3の出力調停部は、クロスバNC3000のポート0のパケット受信部とクロスバNC3000のポート1のパケット受信部から同時に調停リクエストを受信する。この時、クロスバNC3000のポート3の出力調停部はポート0とポート1から受信するレイテンシフィールド値(ポート0:ポート1 = 30:20)を比較し、値の大きいポート0側のリクエストを優先的に通過させるリクエストであると判断して、クロスバNC3000のポート0に対してGNT信号を出力する。クロスバNC3000のポート3の出力調停部は、ポート0のパケット受信部からポート3行きのパケット出力が完了するまで、ポート1に対するGNT信号の出力は抑止する。その間、クロスバNC3000のポート1のパケット受信部は調停待ち状態となる。

【 0 0 4 9 】

クロスバNC3000のポート1のパケット受信部は、調停待ち状態の期間中、CPU3011からのリクエストパケットのパケットヘッダのレイテンシフィールドに調停待ち時間を加算し続ける。また、該加算したレイテンシフィールド値を次回の調停のレイテンシ値として使用されるように、クロスバNC3000のポート3の出力調停部へ出力し続ける。

【 0 0 5 0 】

クロスバNC3000のポート3の出力調停部は、クロスバNC3000のポート0のパケット受信部がパケットの出力を完了すると、クロスバNC3000のポート1のパケット受信部に対してGNT信号を出力する。クロスバNC3000のポート3の出力調停部からGNT信号を受信したクロスバNC3000のポート1のパケット受信部は、調停による待ち時間(=15)をパケットヘッダのレイテンシフィールドに加算する。その結果、レイテンシ値 = 35の状態、該パケットがクロスバNC3000のポート3に出力される。

【 0 0 5 1 】

クロスバNC3000のポート3から出力されたリクエストCPU3011からのリクエストパケットは、クロスバNC3200のポート2のパケット受信部に受信される。クロスバNC3200のポート2のパケット受信部は、受信したリクエストCPU3011からのリクエストパケットのパケットヘッダを解析し、共有資源IOH3220と接続するクロスバNC3200のポート1の出力調停部へ調停リクエストおよびレイテンシ値 = 35を出力する。

【 0 0 5 2 】

クロスバNC3200のポート2のパケット受信部がリクエストCPU3011からのリクエストパケットを受信すると同時に、クロスバNC3200のポート3のパケット受

10

20

30

40

50

信部は、リクエスタCPU3110から共有資源IOH3220行きのリクエストパケット(レイテンシフィールド値=30)を受信する。クロスバNC3200のポート3のパケット受信部は、ポート2のパケット受信部と同じく、クロスバNC3200のポート1の出力調停部に対して調停リクエストを出力する。

【0053】

クロスバNC3200のポート1の出力調停部は、クロスバNC3200のポート2のパケット受信部とクロスバNC3200のポート1のパケット受信部から同時にリクエストを受信する。

【0054】

本実施の形態の構成上、クロスバNC3200のポート2のパケット受信部に存在するリクエスタCPU3011からのリクエストパケットから共有資源IOH3220までの距離よりも、クロスバNC3200のポート3のパケット受信部に存在するリクエスタCPU3110からのリクエストパケットから共有資源IOH3220までの距離の方が長距離である。そのため、最短レイテンシの値ではリクエスタCPU3110からのリクエストパケットの方が高優先のリクエストとなる。しかし、本実施の形態のケースでは、リクエスタCPU3011からのリクエストパケットのレイテンシ値の方が、クロスバNC3000で待たされた時間分が加算されている。従って、クロスバNC3200のポート1の出力調停部において、クロスバNC3200のポート2が高優先なポートと判断され、ポート2に対するGNT信号が先に出力される。

【0055】

クロスバNC3200のポート1の出力調停部からGNT信号を受信したクロスバNC3200のポート2のパケット受信部は、該パケットをポート1に送出する。この際、調停による待ち時間は無いため、パケットヘッダのレイテンシフィールドの値は加算されない。

【0056】

クロスバNC3200のポート1から送出され、共有資源IOH3220に到達したリクエスタCPU3011からのリクエストパケットの総レイテンシは、該パケットヘッダのレイテンシフィールドの値と一致する。そのため総レイテンシは、最短レイテンシにクロスバNC3000およびNC3200で待たされた時間を加えた値となる。本実施例の場合、調停待ち時間は、クロスバNC3000内での待ち時間の15だけであり、クロスバNC3200では優先的に通過したため待ち時間は無く、そのパケットは結果レイテンシ値=35で通過出来る。結果として、クロスバNC3000での調停待ち発生以降のレイテンシの低下は抑制されたことになる。

【0057】

以上で説明した機構により、転送経路の長いリクエストパケットは優先的にクロスバ内の調停部を通過できる。また、途中経路の競合で待たされたパケットも、その後の競合においては優先的に経路を通過できる可能性が高くなる。

【0058】

これらの効果から、ボトルネックとなりうる長経路のパケット、および、経路上の競合のために待たされたパケットのレイテンシの低下を抑制できる。その結果、大規模SMPコンピュータシステムの性能を改善することが可能となる。

【図面の簡単な説明】

【0059】

【図1】優先調停システム、優先調停方法が適用されるシステムの構成を示す。

【図2】クロスバの構成と動作を説明するための図である。

【図3】パケットの転送について説明するための図である。

【図4】ルーティングテーブルとパケットヘッダを示す。

【符号の説明】

【0060】

CPU1010、1011、1110、1111 リクエスタ

10

20

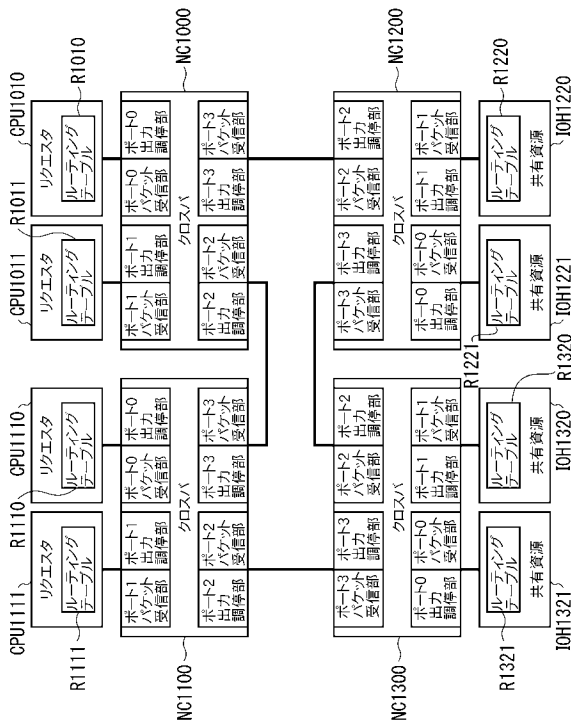
30

40

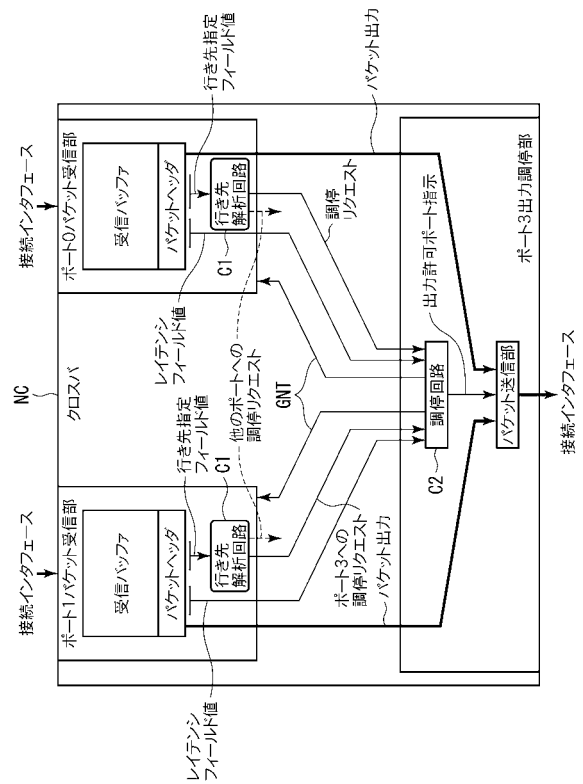
50

- R 1 0 1 0、1 0 1 1、1 1 1 0、1 1 1 1、1 2 2 0、1 2 2 1、1 3 2 0、1 3 2 1
- ルータイングテーブル
- NC 1 0 0 0、1 1 0 0、1 2 0 0、1 3 0 0 クロスバ
- IOH 1 2 2 0、1 2 2 1、1 3 2 0、1 3 2 1 共有資源
- C 1 行き先解析回路
- C 2 調停回路
- CPU 3 0 1 0、3 0 1 1、3 1 1 0 リクエスト
- R 3 0 1 0、3 0 1 1、3 1 1 0、3 2 2 0、3 3 2 0 ルータイングテーブル
- IOH 3 2 2 0、3 3 2 0 共有資源
- R ルータイング
- F 1 行き先指定フィールド
- F 2 レイテンシフィールド

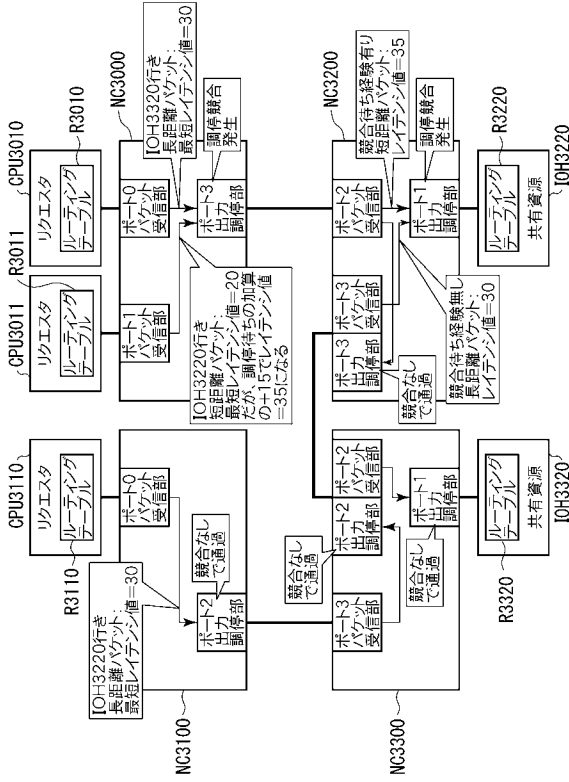
【 図 1 】



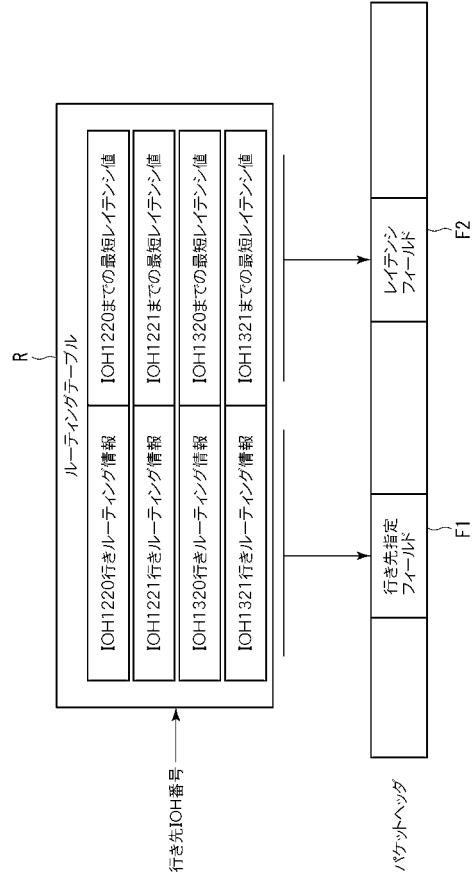
【 図 2 】



【 図 3 】



【 図 4 】



フロントページの続き

- (56)参考文献 特開2000-293495(JP,A)
特開2005-173859(JP,A)
特開平08-186577(JP,A)
特開2004-086304(JP,A)

(58)調査した分野(Int.Cl., DB名)

H04L 12/56
G06F 15/173