

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4662273号
(P4662273)

(45) 発行日 平成23年3月30日(2011.3.30)

(24) 登録日 平成23年1月14日(2011.1.14)

(51) Int.Cl. F I
G06F 13/00 (2006.01) G06F 13/00 353A

請求項の数 15 (全 66 頁)

(21) 出願番号	特願2006-83185 (P2006-83185)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成18年3月24日 (2006. 3. 24)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2007-257479 (P2007-257479A)	(74) 代理人	100079359 弁理士 竹内 進
(43) 公開日	平成19年10月4日 (2007. 10. 4)	(72) 発明者	中島 耕太 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
審査請求日	平成20年8月6日 (2008.8.6)	(72) 発明者	久門 耕一 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
		審査官	須藤 竜也

最終頁に続く

(54) 【発明の名称】 通信装置、方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

パケットを作成して送信する送信部と、パケットを受信する受信部とを備えた通信をサポートする通信装置に於いて、

前記送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して送信するパケット送信部と、

転送先から再送要求を受信した際に、要求された転送データからパケットを作成して送信するパケット再送部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、

を備え、

前記受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域

管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、

前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。

【請求項 2】

パケットを作成して送信する送信部と、パケットを受信する受信部とを備えた通信をサポートする通信装置に於いて、

前記送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して送信するパケット送信部と、

転送先から受信不許可通知を受信した際に、前記パケット送信部によるパケット転送を中断する転送中断部と、

転送先から再送要求を受信した際に、要求された転送データからパケットを作成して送信するパケット再送部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、

を備え、

前記受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、

前記パケット破棄部で前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知部と、

前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。

【請求項 3】

パケットを作成して送信する送信部と、パケットを受信する受信部とを備えた通信をサポートする通信装置に於いて、

前記送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して繰返し送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、

を備え、

前記受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別し

10

20

30

40

50

た場合に、受信パケットを破棄するパケット破棄部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。

【請求項 4】

パケットを作成して送信する送信部と、パケットを受信する受信部とを備えた通信をサポートする通信装置に於いて、

前記送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管理情報に記録するバッファ転送部と、

前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動させるデータ移動部と、

前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。

【請求項 5】

パケットを作成して送信する送信部と、パケットを受信する受信部とを備えた通信をサポートする通信装置に於いて、

転送先に許可を問合せることなく転送し、不許可の場合には再送要求を受けて再送する第 1 転送モード処理部と、

転送先に問合せることなく転送し、不許可の場合は不許可通知を受けて転送を中断し、再送要求を受けて再送する第 2 転送モード処理部と、

転送完了通知を受けるまで転送先に許可を問合せることなく繰返し転送する第 3 転送モード処理部と、

転送先に許可を問合せることなく転送し、不許可の場合は一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する第 4 転送モード処理部と、

転送先に許可を問合せ、許可通知を受けて転送する第 5 転送モード処理部と、

前記第 1 乃至第 5 転送モード処理部のいずれか 1 つを、ネットワーク負荷と受信側のメモリ使用量の少なくともいずれか一方に基づいて選択してデータ転送を実行させる転送モード選択制御部と、

を備えたことを特徴とする通信装置。

【請求項 6】

請求項 5 記載の通信装置に於いて、

前記第 1 転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して送信するパケット送信部と、

転送先から再送要求を受信した際に、要求された転送データからパケットを作成して送信するパケット再送部と、

10

20

30

40

50

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記第1転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合
に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不
許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別し
た場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記
録するパケット破棄部と、

前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した
場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信
する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知
部と、

を備えたことを特徴とする通信装置。

【請求項7】

請求項5記載の通信装置に於いて、

前記第2転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送デー
タからパケットを作成して送信するパケット送信部と、

転送先から受信不許可通知を受信した際に、パケット送信部によるパケット転送を中断
する転送中断部と、

転送先から再送要求を受信した際に、要求された転送データからパケットを作成して送
信するパケット再送部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記第2転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合
に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不
許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別し
た場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記
録するパケット破棄部と、

前記パケット破棄部で前記受信領域の転送不許可を判別した場合に、転送元に受信不許
可通知を送信する不許可通知部と、

前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した
場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信
する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知
部と、

を備えたことを特徴とする通信装置。

【請求項8】

請求項5記載の通信装置に於いて、

前記第3転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送デー
タからパケットを作成して繰返し送信するパケット送信部と、

10

20

30

40

50

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記第3転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合
に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不
許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別し
た場合に、受信パケットを破棄するパケット破棄部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知
部と、

を備えたことを特徴とする通信装置。

【請求項9】

請求項5記載の通信装置に於いて、

前記第4転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送デー
タからパケットを作成して送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記第4転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合
に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不
許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別し
た場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域
管理情報に記録するバッファ転送部と、

前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情
報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動させ
るデータ移動部と、

前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了
通知部と、

を備えたことを特徴とする通信装置。

【請求項10】

請求項5記載の通信装置に於いて、

前記第5転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可確認を送信して受信許可通知を受信し
た場合に、転送データからパケットを作成して送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

前記第5転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理部と、

転送元から受信確認通知を受信した時に前記転送領域管理情報を参照して受信許可通知
又は受信不許可通知を送信する確認応答部と、

前記受信領域に受信パケットを転送するパケット受信部と、

前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する転送
完了通知部と、

10

20

30

40

50

を備えたことを特徴とする通信装置。

【請求項 1 1】

請求項 5 記載の通信装置に於いて、前記転送モード設定部は、
前記受信側のメモリ使用率が少ない場合は前記第 4 転送モード処理部を選択し、
前記受信側のメモリ使用率が多い場合は、ネットワーク負荷の低い順に、前記第 3 転送
モード処理部、第 1 転送モード処理部及び第 5 転送モード処理部を順次選択することを特
徴とする通信装置。

【請求項 1 2】

パケットを作成して送信する送信ステップと、パケットを受信する受信ステップとを備
えた通信をサポートする通信方法に於いて、

コンピュータに、

パケット送信部がデータ転送要求を受けた際に、転送先に受信許可の有無を問合せること
なく、転送データからパケットを作成して送信するパケット送信ステップと、

パケット再送部が転送先から再送要求を受信した際に、要求された転送データからパケ
ットを作成して送信するパケット再送ステップと、

送信完了部が転送先から転送完了通知を受信してパケット送信を正常終了する送信完了
ステップと、

転送領域管理部が受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を
管理する転送領域管理ステップと、

パケット受信部がパケット受信時に前記転送領域管理情報を参照して受信領域の転送許
可を判別した場合に、前記受信領域に受信パケットを転送すると共に前記受信領域の転送
許可を転送不許可に変更するパケット受信ステップと、

パケット受信時にパケット破棄部が前記転送領域管理情報を参照して前記受信領域の転
送不許可を判別した場合に、受信パケットを破棄すると共に前記受信パケットの破棄を前
記転送領域管理情報に記録するパケット破棄ステップと、

前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判
別した場合に、再送要求部が前記転送領域管理情報のパケット破棄の記録に基づいて転送
元に再送要求を送信する再送要求ステップと、

完了通知部が受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信
する完了通知ステップと、

を実行させることを特徴とする通信方法。

【請求項 1 3】

パケットを作成して送信する送信ステップと、パケットを受信する受信ステップとを備
えた通信をサポートする通信方法に於いて、

コンピュータに、

パケット送信部がデータ転送要求を受けた際に、転送先に受信許可の有無を問合せること
なく、転送データからパケットを作成して送信するパケット送信ステップと、

転送中断部が転送先から受信不許可通知を受信した際に、前記パケット送信ステップに
よるパケット転送を中断する転送中断ステップと、

パケット再送部が転送先から再送要求を受信した際に、要求された転送データからパケ
ットを作成して送信するパケット再送ステップと、

送信完了部が転送先から転送完了通知を受信してパケット送信を正常終了する送信完了
ステップと、

転送領域管理部が受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を
管理する転送領域管理ステップと、

パケット受信部がパケット受信時に前記転送領域管理情報を参照して受信領域の転送許
可を判別した場合に、前記受信領域に受信パケットを転送すると共に前記受信領域の転送
許可を転送不許可に変更するパケット受信ステップと、

パケット破棄部がパケット受信時に前記転送領域管理情報を参照して前記受信領域の転
送不許可を判別した場合に、受信パケットを破棄すると共に前記パケットの破棄を前記転

10

20

30

40

50

送領域管理情報に記録するパケット破棄ステップと、

不許可通知部が前記パケット破棄ステップで前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知ステップと、

前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、再送要求部が前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、

完了通知部が受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を実行させることを特徴とする通信方法。

【請求項 14】

パケットを作成して送信する送信ステップと、パケットを受信する受信ステップとを備えた通信をサポートする通信方法に於いて、

コンピュータに、

パケット送信部がデータ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して繰返し送信するパケット送信ステップと、

送信完了部が転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、

転送領域管理部が受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信部がパケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

パケット破棄部がパケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄ステップと、

完了通知部により受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を実行することを特徴とする通信方法。

【請求項 15】

パケットを作成して送信する送信ステップと、パケットを受信する受信ステップとを備えた通信をサポートする通信方法に於いて、

コンピュータに、

パケット送信部がデータ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからパケットを作成して送信するパケット送信ステップと、

送信完了部が転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、

転送領域管理部が受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信部がパケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

バッファ転送部がパケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ領域への転送を前記転送領域管理情報に記録するバッファ転送ステップと、

データ移動部が前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動させるデータ移動ステップと、

完了通知部が前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を実行させることを特徴とする通信方法。

【発明の詳細な説明】

10

20

30

40

50

【技術分野】

【0001】

本発明は、RDMA (Remote Direct Memory Access) 方式によりデータを高速転送する通信装置、方法及びプログラムに関し、特に受信先の許可を問合せることなく投機的にRDMAパケットを転送する通信装置、方法及びプログラムに関する。

【背景技術】

【0002】

従来、大規模なPCクラスタや並列分散サーバシステムでは、ノード間やストレージ・ホスト間の接続を高性能化するため、高い通信性能を求められる部分においては、ソケット通信に代表されるセンド/レシーブ (send/receive) 型通信の代わりに、リード/ライト (Read/Write) 型通信のRDMA方式が用いられている。

10

【0003】

センド/レシーブ型通信では、転送元ノードが転送元のメモリ領域を指定し、転送先ノードが転送先のメモリ領域を指定することにより通信を行うが、RDMA方式は、転送起動側が転送元及び転送先の双方のメモリ領域を指定するため、転送の際に必要な中間バッファを削減でき、中間バッファと転送元及び転送先のメモリ領域間で発生するコピー処理を削減できる。

【0004】

特に、インフィニバンド (InfiniBand) や10ギガビットイーサネット (R) (10GbEthernet) のように、1ギガバイト/秒 (1GB/s) クラスの高速なネットワークを用いる際には、コピー処理の削減により大幅に通信性能を改善することができる。

20

【0005】

RDMA通信では、転送起動側が転送元及び転送先の双方のメモリ領域を指定し、起動処理を行う。この際、通信対象となるノードには、転送の開始及び終了が通知されない。したがって、実際には、RDMA通信を起動する際には、通信対象となる転送先のノードに対して、転送先のメモリ領域への転送の開始許可を求めるための通信処理を事前に行う必要がある。

【0006】

従来のRDMA通信のデータ転送手順は次のようになる。

(1) 送信側

送信側は、RDMA起動可能になったら、受信側に受信許可の有無を問い合わせる。受信許可を受信してから、RDMAデータを転送する。最後に、ack (転送完了通知) を受信すると転送を完了する。

(2) 受信側

受信側は受信許可の有無の問い合わせに対し、受信可能になった時点で、受信許可を転送する。その後、RDMAデータを受信完了すると、ackを送信側へ返す。

30

【特許文献1】特開2004-192179号公報

【特許文献2】特開2005-044353号公報

40

【発明の開示】

【発明が解決しようとする課題】

【0007】

しかしながら、このような従来のRDMA通信にあっては、RDMA通信を開始する前に、転送先に対し転送許可を求めるための通信処理が必要となり、その分、データ転送に時間がかかるという問題がある。この転送先に対し転送許可を求めるための通信処理は、特に高スループット、高レイテンシであるネットワークにおいて大きな問題となる。

【0008】

例えば、近年の10ギガビットイーサネット (なおイーサネットは登録商標である。以下、省略する) で使用するNIC (ネットワークインタフェースカード) の実装では、ス

50

ループットは1ギガバイト/秒以上であるが、レイテンシは10～20 μ s程度と大きい。

【0009】

このような場合に、例えば1キロバイトのデータといった短いデータをRDMA転送する場合を考えると、データ転送にかかる時間は、RDMA転送自体は1 μ s以下となるが、事前に必要となる転送許可を求める通信処理で発生する時間は10～20 μ sであるため、レイテンシがボトルネックとなり、スループットを十分に発揮できないという問題がある。

【0010】

本発明は、転送先に対し転送許可を求めるための通信処理を不要にしてスループットを高めるようにしたRDMA通信をサポートする通信装置、方法及びプログラムを提供することを目的とする。

【課題を解決するための手段】

【0011】

本発明はRDMA通信をサポートする通信装置を提供する。

【0012】

(第1転送モード処理)

本発明は、RDMAパケットを作成して送信する送信部と、RDMAパケットを受信する受信部とを備えた第1転送モードのRDMA通信をサポートする通信装置を提供する。

【0013】

第1転送モード処理は、転送先に許可を問合せことなく投機的に転送し、不許可の場合には再送要求を受けて再送する処理である。

【0014】

ここで、送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せことなく、転送データからRDMAパケットを作成して投機的に送信するパケット送信部と、

転送先から再送要求を受信した際に、要求された転送データからRDMAパケットを作成して送信するパケット再送部と、

転送先から転送完了通知(ack)を受信してパケット送信を正常終了する送信完了部と、

を備え、

受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、受信領域に前記受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共にパケット破棄を転送領域管理情報に記録するパケット破棄部と、

パケット破棄部で受信パケットを破棄した後に受信領域の転送許可を判別した場合に、転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知(ack)を送信する完了通知部と、

を備えたことを特徴とする。

【0015】

(第2転送モード処理)

本発明の別の形態にあっては、RDMAパケットを作成して送信する送信部と、RDMAパケットを受信する受信部とを備えた第2転送モードのRDMA通信をサポートする通

10

20

30

40

50

信装置を提供する。

第2転送モードは、転送先に問合せることなく投機的に転送し、不許可の場合は不許可通知を受けて転送を中断し、再送要求を受けて再送する処理である。

【0016】

送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に送信するパケット送信部と、

転送先から受信不許可通知を受信した際に、パケット送信部によるパケット転送を中断する転送中断部と、

転送先から再送要求を受信した際に、要求された転送データからRDMAパケットを作成して送信するパケット再送部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、受信領域に前記受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共にパケット破棄を転送領域管理情報に記録するパケット破棄部と、

パケット破棄部で受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知部と、

パケット破棄部で受信パケットを破棄した後に受信領域の転送許可を判別した場合に、転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする。

【0017】

(第3転送モード処理)

本発明の別の形態にあつては、RDMAパケットを作成して送信する送信部と、RDMAパケットを受信する受信部とを備えた第3転送モードのRDMA通信をサポートする通信装置を提供する。第3転送モードは、転送完了通知を受けるまで転送先に許可を問合せることなく繰返し転送する処理である。

【0018】

送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に繰返し送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、受信領域に前記受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する転送完了通知部と、
を備えたことを特徴とする。

【 0 0 1 9 】

(第 4 転送モード処理)

本発明の別の形態にあっては、R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた第 4 転送モードの R D M A 通信をサポートする通信装置を提供する。

第 4 転送モードは、転送先に許可を問合せることなく投機的に転送し、不許可の場合は一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する処理である。

10

【 0 0 2 0 】

送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、

受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

20

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、受信領域に受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共にバッファ転送を転送領域管理情報に記録するバッファ転送部と、

バッファ転送後に受信領域の転送許可を判別した場合に、転送領域管理情報のバッファ転送の記録に基づいてバッファ領域のデータを受信領域に移動するデータ移動部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

30

を備えたことを特徴とする。

【 0 0 2 1 】

(転送モード選択)

本発明の別の形態にあっては、R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた R D M A 通信をサポートする通信装置に於いて、

転送先に許可を問合せることなく投機的に転送し、不許可の場合には再送要求を受けて再送する第 1 転送モード処理部と、

転送先に問合せることなく投機的に転送し、不許可の場合は不許可通知を受けて転送を中断し、再送要求を受けて再送する第 2 転送モード処理部と、

転送完了通知を受けるまで転送先に許可を問合せることなく投機的に繰返し転送する第 3 転送モード処理部と、

40

転送先に許可を問合せることなく投機的に転送し、不許可の場合は一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する第 4 転送モード処理部と、

転送先に許可を問合せ、許可通知を受けて転送する第 5 転送モード処理部と、

第 1 乃至第 5 転送モード処理部のいずれか 1 つを、ネットワーク負荷と受信側のメモリ使用率の少なくともいずれか一方に基づいて選択して R D M A データ転送を実行させる転送モード設定部と、

を備えたことを特徴とする。

【 0 0 2 2 】

ここで第 1 乃至第 4 転送モード処理部は、前述した第 1 転送モード乃至第 4 転送モード

50

の受信部と送信部と同じものである。

【0023】

また第5転送モード処理部は送信部と受信部を備え、
第5転送モード処理部の送信部は、
データ転送要求を受けた際に、転送先に受信許可確認を送信して受信許可通知を受信した場合に、転送データからRDMAパケットを作成して送信するパケット送信部と、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、
第5転送モード処理部の受信部は、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
転送元から受信確認通知を受信した時に転送領域管理情報を参照して受信許可通知又は受信不許可通知を送信する確認応答部と、
受信領域に受信パケットを転送するパケット受信部と、
受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
を備えたことを特徴とする。

10

【0024】

転送モード設定部は転送先の送信部に設けられ、選択した転送モードを送信部及び受信部に通知して第1乃至第5転送モード処理部のいずれかによるRDMAデータ転送を実行させる。

20

【0025】

また転送モード設定部は受信部に設けられ、選択した転送モードを送信部及び受信部に通知して第1乃至第5転送モード処理部のいずれかによるRDMAデータ転送を実行させるようにしても良い。

【0026】

転送モード設定部は、
受信側のメモリ使用率が少ない場合は第4転送モード処理部を選択し、
受信側のメモリ使用率が多い場合は、ネットワーク負荷の低い順に、第3転送モード処理部、第1転送モード処理部及び第5転送モード処理部を順次選択する。

30

【0027】

本発明の通信装置は、更に、論理アドレスと物理アドレスのアドレス変換情報を持ち、アドレス変換情報内に転送領域管理情報を含ませて一体化する。

【0028】

本発明はRDMA通信をサポートする通信方法を提供する。

【0029】

(第1転送モード処理方法)

RDMAパケットを作成して送信する送信ステップと、RDMAパケットを受信する受信ステップとを備えた第1転送モードのRDMA通信をサポートする通信方法に於いて、送信ステップは、

40

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に送信するパケット送信ステップと、

転送先から再送要求を受信した際に、要求された転送データからRDMAパケットを作成して送信するパケット再送ステップと、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を備え、

受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、

50

受信領域に前記受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共にパケット破棄を転送領域管理情報に記録するパケット破棄ステップと、

パケット破棄ステップで受信パケットを破棄した後に受信領域の転送許可を判別した場合に、転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を備えたことを特徴とする。

【 0 0 3 0 】

(第 2 転送モード処理方法)

本発明の別の形態にあつては、R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた第 2 転送モードの R D M A 通信をサポートする通信方法に於いて、

送信ステップは、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、

転送先から受信不許可通知を受信した際に、パケット送信ステップによるパケット転送を中断する転送中断ステップと、

転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送ステップと、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、

を備え、

受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、受信領域に前記受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共にパケット破棄を転送領域管理情報に記録するパケット破棄ステップと、

パケット破棄ステップで受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知ステップと、

パケット破棄ステップで受信パケットを破棄した後に受信領域の転送許可を判別した場合に、転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を備えたことを特徴とする。

【 0 0 3 1 】

(第 3 転送モード処理方法)

本発明の別の形態にあつては、R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた第 3 転送モードの R D M A 通信をサポートする通信方法に於いて、

送信ステップは、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に繰返し送信するパケット送信ステップと、

10

20

30

40

50

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を備え、

受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、
受信領域に受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパ
ケット受信ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に
、受信パケットを破棄するパケット破棄ステップと、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する転送完了
通知ステップと、

を備えたことを特徴とする。

【 0 0 3 2 】

(第 4 転送モード処理方法)

本発明の別の形態にあつては、R D M A パケットを作成して送信する送信ステップと、
R D M A パケットを受信する受信ステップとを備えた第 4 転送モードの R D M A 通信をサ
ポートする通信方法に於いて、

送信ステップは、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送デー
タから R D M A パケットを作成して投機的に送信するパケット送信ステップと、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を備え、

受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域
管理ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送許可を判別した場合に、
受信領域に受信パケットを転送すると共に受信領域の転送許可を転送不許可に変更するパ
ケット受信ステップと、

パケット受信時に転送領域管理情報を参照して受信領域の転送不許可を判別した場合に
受信パケットをバッファ領域に転送すると共にバッファ転送を転送領域管理情報に記録す
るバッファ転送ステップと、

バッファ転送後に受信領域の転送許可を判別した場合に、転送領域管理情報のバッファ
転送の記録に基づいてバッファ領域のデータを受信領域に移動するデータ移動ステップと
、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知
ステップと、

を備えたことを特徴とする。

【 0 0 3 3 】

(転送モード選択方法)

本発明の別の形態にあつては、R D M A パケットを作成して送信する送信ステップと、
R D M A パケットを受信する受信ステップとを備えた R D M A 通信をサポートする通信方
法に於いて、

転送先に許可を問合せることなく投機的に転送し、不許可の場合には再送要求を受けて
再送する第 1 転送モード処理ステップと、

転送先に問合せることなく投機的に転送し、不許可の場合は不許可通知を受けて転送を
中断し、再送要求を受けて再送する第 2 転送モード処理ステップと、

転送完了通知を受けるまで転送先に許可を問合せることなく投機的に繰返し転送する第
3 転送モード処理ステップと、

転送先に許可を問合せることなく投機的に転送し、不許可の場合には一時バッファに保存

10

20

30

40

50

し、許可が得られたら一時バッファから受信領域に転送する第4転送モード処理ステップと、

転送先に許可を問合せ、許可通知を受けて転送する第5転送モード処理ステップと、

第1乃至第5転送モード処理ステップのいずれか1つを、ネットワーク負荷と受信側のメモリ使用量の少なくともいずれか一方に基づいて選択してRDM Aデータ転送を実行させる転送モード設定ステップと、

を備えたことを特徴とする。

【0034】

更に、本発明は、RDM A通信をサポートする通信装置のコンピュータに、前述した通信方法の各ステップを実行させる通信プログラムを提供する。

【発明の効果】

【0035】

本発明によれば、転送先に受信許可の有無を問合せることなく投機的にRDM Aパケットを転送するため、もし受信領域が受信可能であれば、受信したRDM Aパケットのデータ本体を受信領域に格納してデータ転送を完了でき、事前に転送許可を求める通信処理を必要としないため、転送許可を求める通信処理に必要としていた10~20μs程度のレイテンシを不要にすることができる。このため転送許可通信処理のレイテンシより短い転送時間で済むデータサイズのRDM Aパケット転送のスループットを高めることができ、通信処理全体としてみた場合の通信性能を向上できる。

【0036】

また投機的なRDM Aパケットの転送に対し転送先で受信不可能であった場合についても、受信可能となった際の再送要求により再送させることで(第1及び第2転送モード)、受信許可の有無を問合せる従来方式よりは短い転送時間とすることができる。

【0037】

またネットワークの帯域に余剰があることを条件に、転送完了通知(ack)が得られるまで転送データ分のRDM Aパケットを繰返し転送することで(第3転送モード)、受信可能になれば、その時点から有効に受信され、従来方式よりは短い転送時間とすることができる。

【0038】

更に、転送先のメモリ領域に余剰が場合には、受信不可能な状態で受信したRDM Aパケットのデータ本体を一時バッファに保存し、受信許可となったらRDM Aパケットを有効に受信し且つ一時バッファのデータを受信領域に移動させることで(第4転送モード)、従来方式よりは短い転送時間とすることができる。

【0039】

更に本発明にあっては、ネットワーク負荷や受信側のメモリ使用率に基づいて最適な転送モードとなるように、投機的転送を行う第1乃至第4転送モード、従来方式である第5転送モードのいずれか1つを動的に選択してRDM Aパケット転送を行うことで、各転送モードの利点を最大限に生かした効率的なデータ転送を行い、全体的に見た通信性能を向上させることができる。

【発明を実施するための最良の形態】

【0040】

(通信装置の構成)

図1はホスト間通信に適用された本発明の一実施形態の説明図である。図1において、ホスト10-1, 10-2に対しては、それぞれ10Gbイーサネットのネットワークインタフェースカード(以下「NIC」という)12-1, 12-2が設けられ、ネットワーク18を介して接続している。

【0041】

NIC12-1には送信部14-1と受信部16-1が設けられ、同様にNIC12-2にも送信部14-2と受信部16-2が設けられる。NIC12-1, 12-2はネットワーク18を介してホスト10-1, 10-2間でのRDM Aパケット転送をサポート

10

20

30

40

50

する。

【 0 0 4 2 】

図 2 は本実施形態における N I C の送信部と受信部の機能構成のブロック図である。図 2 においては、ホスト 1 0 - 1 側を送信ノード（転送元）としており、ホスト 1 0 - 1 に設けた N I C 1 2 - 1 には送信部 1 4 - 1 としての機能が設けられている。一方、ホスト 1 0 - 2 は受信側となる受信ノード（転送先）であり、ホスト 1 0 - 2 に設けた N I C 1 2 - 2 は受信部 1 6 - 2 としての機能を示している。

【 0 0 4 3 】

送信側の N I C 1 2 - 1 には、送信要求受付部 2 6、転送領域管理部 2 8、送信制御部 3 0、受信制御部 3 2、モード選択制御部 3 4、ネットワーク負荷情報収集部 3 6、1 0 G b イーサネット M A C（「1 0 G b E M A C」と表記）3 8 が設けられる。

10

【 0 0 4 4 】

送信要求受付部 2 6 は送信要求を受け付ける送信要求受付レジスタを備える。この送信要求受付レジスタは P C I を経由してホスト 1 0 - 1 のアドレス空間にマッピングできる。このためホスト 1 0 - 1 の送信要求部 2 0 は、アドレス空間にマッピングされた送信要求受付レジスタにパラメータを書き込むことにより送信要求を発行する。また送信要求部 2 0 は、ホスト 1 0 - 1 のメモリ空間である送信領域 2 2 に配置している送信データ 2 4 を指定して送信要求を行う。

【 0 0 4 5 】

転送領域管理部 2 8 は、論理アドレスと物理アドレスの変換を行うアドレス変換表と、転送領域の受信許可の有無を管理する転送領域管理表を備える。本実施形態にあつては、アドレス変換表と転送領域管理表は一体化されている。即ち、転送領域管理表の情報をアドレス変換表に含ませることにより転送領域管理表の情報参照をアドレスの参照と同時に実行可能とし、これにより情報参照の性能を向上することができる。

20

【 0 0 4 6 】

送信制御部 3 0 は、送信要求受付部 2 6 や受信制御部 3 2 からの送信要求に基づき、ホスト 1 0 - 1 のメモリ上の送信領域 2 2 にある送信データ 2 4 を D M A 転送により取得する。この際に転送領域管理部 2 8 から送信データ 2 4 の論理アドレスから変換した物理アドレスを取得する。

【 0 0 4 7 】

このようにしてホスト 1 0 - 1 側から取得した送信データにヘッダを付加した R D M A データパケットを作成し、1 0 G b イーサネット M A C 3 8 に転送する。この場合に、モード選択制御部 3 4 からの指示に従い、選択された転送モードの転送手順に従って送信処理を実行する。

30

【 0 0 4 8 】

受信制御部 3 2 は受信側の N I C 1 2 - 2 より

- (1) a c k パケット
- (2) 再送制御パケット
- (3) 受信許可通知パケット
- (4) 受信不許可通知パケット
- (5) 転送モード通知パケット

40

を受信する。

【 0 0 4 9 】

a c k パケットを受信したときには転送処理を正常終了する。再送制御パケットを受信したときには再送を送信制御部 3 0 に通知する。受信許可通知パケットを受信したときには受信許可を送信制御部 3 0 に通知する。受信不許可通知パケットを受信したときには不許可となっている受信領域に対する R D M A 送信処理を停止する。転送モード通知パケットを受信したときにはモード選択制御部 3 4 に通知する。

【 0 0 5 0 】

モード選択制御部 3 4 は、本実施形態で準備されている第 1 乃至第 5 転送モードのうち

50

の最適な転送モードを選択し、それに基づいた転送手順を行うように送信制御部 30 に通知する。転送モードの選択を行う際には、ネットワーク負荷情報収集部 36 からのネットワーク負荷情報と、受信制御部 32 で受信された受信側の NIC 12 - 2 によるモード選択指示情報を用いる。

【 0 0 5 1 】

ここで本実施形態における第 1 乃至第 5 転送モードは次の内容を持つ。

(1) 第 1 転送モードは、転送先に許可を問い合わせることなく投機的に転送し、不許可の場合には再送要求を受けて再送する。

(2) 第 2 転送モードは、転送先に許可を問い合わせることなく投機的に転送し、不許可の場合には不許可通知を受けて転送を中断し、その後、再送要求を受けて再送する。

(3) 第 3 転送モードは、ack パケットを受け取るまで投機的に転送先に許可を問い合わせることなく繰り返し転送する。

(4) 第 4 転送モードは、転送先に許可を問い合わせることなく投機的に転送し、不許可の場合には転送先で一時的にバッファに保存し、許可が得られたら一時バッファから受信領域に転送する。

(5) 第 5 転送モードは、従来方式の転送モードであり、転送先に許可を問い合わせ、許可通知を受けて転送する。

【 0 0 5 2 】

これら第 1 乃至第 5 転送モードの詳細は後の説明で明らかにする。

【 0 0 5 3 】

次に受信側の NIC 12 - 2 を説明する。受信側の NIC 12 - 2 には、受信部 16 - 2 として機能する構成として、受信制御部 40、転送領域管理部 42、送信制御部 44、モード選択制御部 46、ネットワーク負荷情報収集部 48、メモリ使用率情報収集部 49 及び 10Gbイーサネット MAC 50 が設けられている。

【 0 0 5 4 】

受信制御部 40 は、RDMA データパケットや従来方式である第 5 転送モードにおける受信許可確認通知の制御パケットを受信する。RDMA データパケットを受信した場合、ヘッダを解析し、受信領域が受信可能であるか否かを転送領域管理部 42 に問い合わせる。この問合せに対し転送領域管理部 42 は、受信領域の受信許可の有無を通知し、同時に論理アドレスから変化した物理アドレスも通知する。

【 0 0 5 5 】

受信可能な場合は、DMA 転送によりホスト 10 - 2 の受信領域 54 に DMA データパケットから得られたデータ本体を転送する。受信不許可の場合には、モード選択制御部 46 の現在選択状態にある転送モードの指示に従い、パケットを破棄するか (第 1 乃至第 3 転送モード)、ホスト 12 - 2 の一時バッファ 56 の所定の位置に転送して保存する (第 4 転送モード)。

【 0 0 5 6 】

またモード選択制御部 46 より第 2 転送モードが指示されていた場合には、転送元に受信不許可通知を制御パケットにより転送するように送信制御部 44 に要求する。

【 0 0 5 7 】

転送領域管理部 42 は、送信側の NIC 12 - 1 の場合と同様、論理物理アドレス変換表と転送領域管理表を持ち、両者は一体化されている。受信制御部 40 からの受信領域に対する受信許可の有無の問合せに対し、転送領域管理表を参照して応答する。この際、受信許可の場合はアドレス変換表から論理アドレスに対応した物理アドレスを取得して返す。

【 0 0 5 8 】

このとき同時に受信許可を通知した受信領域の状態を受信不許可に変更して二重受信を防止する。

【 0 0 5 9 】

一方、受信制御部 40 からの問合せに対し受信領域が受信不許可の場合は、受信不許可

10

20

30

40

50

を応答すると同時に、受信不許可となっている受信領域に受信されたRDMAPケットデータが破棄されたことを記録する。更に第4転送モードが指示されている場合には、転送データ58が保存された一時バッファ56のアドレスを転送領域管理部42に通知して転送領域管理表に記録する。

【0060】

受信制御部40からの問合せに対し、一度、受信領域の不許可を応答した後に、ホスト10-2の領域許可部52の通知を受けて受信許可になった場合には、転送領域管理表を参照してデータが破棄されている転送要求があったかどうかを確認する。データ破棄となった転送要求の記録があった場合には、一時バッファ56への転送の有無を確認する。転送している場合には、これは第4転送モードの場合であることから、一時バッファ56の転送先のアドレスをホスト10-2に通知し、一時バッファ56から受信領域54へのデータ転送処理を転送制御部60に行わせる。

10

【0061】

一時バッファ56に保持していない場合には再送要求を送信制御部44に通知する。更に転送領域管理部42は、従来方式である第5転送モードの状態を受信許可確認通知を受信制御部40より受け取った場合は、対応する受信領域が受信許可となった時点で受信許可通知の送信を送信制御部44に依頼する。

【0062】

送信制御部44は、受信制御部40からのack通知受信や不許可通知の要求に従って、これらの制御パケットを生成し送信する。また転送領域管理部42からの再送要求や受信許可通知に従って、これらの制御パケットを生成して送信する。

20

【0063】

モード選択制御部46は、本実施形態でサポートしている第1乃至第5転送モードのうちの最適な転送モードを選択し、それに基づいた転送手順を行うように受信制御部40に通知する。また現在選択している転送モードを、送信制御部44を介して送信側のNIC12-1に通知する。

【0064】

モード選択制御部46で転送モードを選択する際には、メモリ使用率情報収集部49で収集している受信側のメモリ使用率と、ネットワーク負荷情報収集部48で収集しているネットワーク負荷情報を使用する。

30

【0065】

ネットワーク負荷情報収集部48は、10GbイーサネットMAC50からネットワーク負荷情報を取得し、収集する。またRDMAP転送開始から完了までの転送完了時間を計測する。メモリ使用率情報収集部49は、ホスト10-2上に確保可能な一時バッファ56の使用率を収集する。

【0066】

10GbイーサネットMAC50は、送信制御部44から渡されたパケットをネットワークへ送出し、またネットワークより受信したパケットを受信制御部40に渡す。またネットワークの負荷情報を取得し、ネットワーク負荷情報収集部48に受け渡す。

【0067】

次に、送信側のホスト10-1の機能を説明する。ホスト10-1には送信要求部20が設けられる。送信要求部20は、通信プロセスからの通信要求を受け、送信側のNIC12-1にRDMAP送信を要求する。

40

【0068】

一方、受信側のホスト10-2には領域許可部52と転送制御部60の機能が設けられる。領域許可部52は受信領域54を使用している通信プロセスからの領域許可を受け、受信側NIC12-2に許可を受けた領域についての受信許可を通知する。転送制御部60は受信側のNIC12-2の転送領域管理部42から通知を受け、一時バッファ56に転送保持されたデータを、受信許可となった受信領域54で移動させる。

【0069】

50

一時バッファ56から受信領域54へのデータ移動については、受信領域54へのデータ転送とページ単位でのアドレスオフセットを併せて行う。即ち一時バッファ56の保存データ58を移動させる際には、先頭ページと末尾ページについてはコピー処理を行うが、それ以外のページについてはページの再マッピングによりコピー処理なしで移動させ、コピー処理のオーバーヘッドを削減する。このデータ移動については後の説明で明らかにする。

【0070】

このような送信側および受信側のNIC12-1, 12-2としての機能は、各々に設けられたCPUとメモリを備えた通信処理用のプロセッサを含むハードウェア環境によるプログラムの実行により実現される。

10

【0071】

図3は本実施形態の受信側で使用する転送領域管理表の説明図である。図3において、転送領域管理表62は、領域ID64、論理アドレス66、物理アドレス68、受信許可70、パケット破棄情報72、受信不許可通知情報74及び一時バッファ保存アドレス75で構成されている。

【0072】

このように前半の論理アドレスを物理アドレスに変換するアドレス変換表と、受信領域の許可の有無、更に転送モードに応じたパケット破棄状態、受信不許可通知情報、更には一時バッファアドレスなどを記録することで一体化し、受信制御部40からの問合せに対し、転送領域管理部42による転送領域管理表62の参照によるアドレス変換、及び受信領域に関する受信許可などの情報の取得を効率的にできるようにしている。

20

【0073】

図4は本実施形態における送信制御部の機能構成のブロック図であり、具体的には図2の送信側のNIC12-1に設けている送信制御部30の機能構成を示している。

【0074】

送信制御部30には、第1転送モード送信部76、第2転送モード送信部78、第3転送モード送信部80、第4転送モード送信部82、第5転送モード送信部84が設けられており、いずれかの転送モードに対応した送信部が図2のモード選択制御部34からの指示により選択されて送信制御を行う。

【0075】

第1転送モード送信部76には、パケット送信部86、パケット再送部88及び送信完了部90が設けられる。パケット送信部86はデータ転送要求を受けた際に、転送先に受信許可の有無を問い合わせることなく転送データからRDMAデータパケットを作成して投機的に送信させる。

30

【0076】

パケット再送部88は、転送先から再送要求を受信した際に、要求された転送データからRDMAデータパケットを作成して送信する。送信完了部90は、転送先から転送完了通知であるackを受信してデータ転送を正常終了する。

【0077】

第2転送モード送信部78は、パケット送信部92、転送中断部94、パケット再送部96及び送信完了部98を設けている。パケット送信部92は、データ転送要求を受けた際に、転送先に受信許可の有無を問い合わせることなく転送データからRDMAデータパケットを作成して投機的に送信する。

40

【0078】

転送中断部94は、転送先から受信不許可通知を受信した際に、パケット送信部92によるパケット転送を中断する。パケット再送部96は、転送先から再送要求を受信した際に、要求された転送データからRDMAパケットを作成して再送する。送信完了部98は、転送先から転送完了通知であるackを受信してデータ転送を正常終了する。

【0079】

第3転送モード送信部80は、パケット送信部100と送信完了部102を備える。パ

50

ケット送信部 100 は、データ転送要求を受けた際に、転送先に受信許可の有無を問い合わせることなく転送データから R D M A データパケットを作成し、投機的に繰り返し送信する。送信完了部 106 は、転送先から転送完了通知である a c k を受信してデータ転送を正常終了する。

【 0 0 8 0 】

このため第 4 転送モード送信部 82 にあっては、送信完了部 106 で転送先から a c k を受信するまで、転送対象となった送信データをパケット送信部 104 から繰り返し転送することになる。

【 0 0 8 1 】

第 5 転送モード送信部 84 は、受信許可確認部 108、パケット送信部 110 及び送信完了部 112 を備える。第 5 転送モード送信部 84 は従来方式であり、このため受信許可確認部 108 は、データ転送要求を受けた際に、転送先に受信許可確認を送信する。パケット送信部 110 は、転送先から受信許可通知を受信した際に、転送データから R D M A パケットを作成して送信する。送信完了部 112 は、転送先から転送完了通知である a c k を受信してデータ転送を正常終了する。

10

【 0 0 8 2 】

図 5 は本実施形態における受信制御部のブロック図であり、具体的には図 2 の受信側の N I C 12 - 2 に設けた受信制御部 40 の機能構成である。

【 0 0 8 3 】

図 5 において、受信制御部 40 は、第 1 転送モード受信部 114、第 2 転送モード受信部 116、第 3 転送モード受信部 118、第 4 転送モード受信部 120 及び第 5 転送モード受信部 124 が設けられる。この第 1 乃至第 5 転送モード受信部は、図 2 のモード選択制御部 46 からの指示により、いずれか 1 つの転送モードが選択されて受信制御を行うことになる。

20

【 0 0 8 4 】

第 1 転送モード受信部 114 には、パケット受信部 126、パケット破棄部 128、再送要求部 130 及び完了通知部 132 が設けられている。パケット受信部 126 は、パケット受信時に転送領域管理部 42 に問い合わせた受信領域の転送許可を判別した際に、その受信領域に受信した R D M A データパケットのデータ本体を転送すると共に、受信領域の転送許可を転送不許可に変更し、2 重受信を防止する。

30

【 0 0 8 5 】

パケット破棄部 128 は、パケット受信時に転送領域管理部 42 に問い合わせた転送不許可を判別した際に、受信パケットを破棄すると共に、パケット破棄を図 3 の転送領域管理表 62 のパケット破棄情報 72 に記録する。

【 0 0 8 6 】

再送要求部 130 は、パケット破棄部 128 でパケットを破棄した後に受信領域の転送許可を判別した際に、転送領域管理表に記録したパケット破棄情報に基づいて転送元に再送要求を送信する。完了通知部 132 は、受信領域に対する受信パケットによるデータ本体の転送完了を認識して、転送完了通知となる a c k パケットを送信する。

【 0 0 8 7 】

第 2 転送モード受信部 116 は、パケット受信部 134、パケット破棄部 136、不許可通知部 138、再送要求部 140 及び完了通知部 142 を設けている。パケット受信部 134 は、パケット受信時に問い合わせた転送領域管理表の参照で得られた受信領域の転送許可を判別した場合、その受信領域に受信パケットのデータ本体を転送すると共に、転送領域管理表における受信領域の転送許可を転送不許可に変更して 2 重受信を防止する。

40

【 0 0 8 8 】

パケット破棄部 136 は、パケット受信時に問い合わせた転送領域管理表に基づく受信領域の転送不許可を判別した際に、受信パケットを破棄すると共に、パケット破棄を転送領域管理表に記録する。

【 0 0 8 9 】

50

受信不許可通知部 1 3 8 は、パケット破棄部 1 3 6 で受信領域の転送不許可を判別した際に、転送元に受信不許可通知を送信する。再送要求部 1 4 0 は、パケット破棄部で受信パケットを破棄した後に受信領域の転送領域を判別した際に、転送領域管理表のパケット破棄の記録に基づいて転送元に再送要求を送信する。更に完了通知部 1 4 2 は、受信領域に対する受信パケットのデータ本体の転送完了を認識して、転送完了通知となる a c k パケットを送信する。

【 0 0 9 0 】

第 3 転送モード受信部 1 1 8 は、パケット受信部 1 4 4、パケット破棄部 1 4 6 及び完了通知部 1 4 8 を設けている。パケット受信部 1 4 4 は、パケット受信時に問い合わせた転送領域管理表の参照で得られた受信領域の転送許可を判別した際に、受信領域に受信パケットのデータ本体を転送すると共に、受信領域の転送許可を転送不許可に変更して 2 重受信を防止する。

10

【 0 0 9 1 】

パケット破棄部 1 4 6 は、パケット受信時に転送領域管理表の参照で通知された受信領域の転送不許可を判別した際に、受信パケットを破棄する。この場合には再送要求は行わないことから、受信パケットの破棄を転送領域管理表に記録することはない。完了通知部 1 4 8 は、受信領域に対する受信パケットのデータ本体の転送完了を認識して、転送完了通知である a c k パケットを送信する。

【 0 0 9 2 】

このような第 3 転送モード受信部 1 1 8 の機能は、図 4 の第 3 転送モード送信部 8 0 より受信許可の有無を問い合わせることなく投機的に繰り返し送られてくる転送データにつき、受信許可が得られてから完了するまでパケットを受信し、受信領域にデータ転送する単純な処理となる。

20

【 0 0 9 3 】

第 4 転送モード受信部 1 2 0 には、パケット受信部 1 5 0、バッファ転送部 1 5 2、データ移動部 1 5 4 及び完了通知部 1 5 6 が設けられる。

【 0 0 9 4 】

パケット受信部 1 5 0 は、パケット受信時に転送領域管理表の参照で得られた受信領域の転送許可を判別した場合、受信領域に受信パケットのデータ本体を転送すると共に、受信領域の転送許可を転送不許可に変更して 2 重受信を防止する。

30

【 0 0 9 5 】

バッファ転送部 1 5 2 は、パケット受信時に転送領域管理表を参照して受信領域の転送不許可を判別した際に、受信パケットから得られたデータ本体を一時バッファのバッファ領域に転送して保存すると共に、一時バッファへの保存を転送領域管理表 6 2 に記録する。

【 0 0 9 6 】

データ移動部 1 5 4 は、一時バッファに対するバッファ転送後に受信領域の転送許可を判別した際に、転送領域管理表のバッファ転送の記録即ち保存アドレスに基づき、一時バッファのバッファ領域の保存データを受信領域に移動する。完了通知部 1 5 6 は、受信領域に対する受信パケットのデータ本体の転送完了を認識して、転送完了通知である a c k

40

【 0 0 9 7 】

第 5 転送モード受信部 1 2 4 には、確認応答部 1 5 8、パケット受信部 1 6 0、完了通知部 1 6 2 が設けられる。第 5 転送モード受信部は従来方式の受信機能を持つ。確認応答部 1 5 8 は、転送元からの確認通知を受信したときに、転送領域管理表を参照して得られた受信許可または受信不許可の通知を送信する。

【 0 0 9 8 】

パケット受信部 1 6 0 は、転送元に対する受信許可の通知に対応して受信された受信パケットのデータ本体を受信領域に転送する。完了通知部 1 6 2 は、受信領域に対する受信パケットのデータ本体の転送完了を認識して、転送完了通知である a c k パケットを送信

50

する。

【 0 0 9 9 】

図 6 は本実施形態における送信処理のフローチャートである。図 6 において、送信処理は、ステップ S 1 でホストからの送信要求の有無をチェックしており、送信要求があると、ステップ S 2 に進み、転送モード選択処理を実行し、ステップ S 3 で選択された転送モードを判定し、ステップ S 4 ~ S 8 に示す第 1 乃至第 5 のいずれかの転送モード送信処理を選択して送信処理を行う。このようなステップ S 1 ~ S 8 の処理を、ステップ S 9 で停止指示があるまで繰り返す。なおステップ S 4 ~ S 8 の第 1 乃至第 5 転送モード送信処理の詳細は後の説明で明らかにする。

【 0 1 0 0 】

図 7 は本実施形態における受信処理のフローチャートである。図 7 において、受信処理は、ステップ S 1 で転送モード選択処理を実行した後、ステップ S 2 で転送モードを判定し、ステップ S 3 ~ S 7 の第 1 乃至第 5 転送モード受信処理のいずれか 1 つを選択して受信処理を行う。このようなステップ S 1 ~ S 7 の処理を、ステップ S 8 で停止指示があるまで繰り返す。

【 0 1 0 1 】

図 7 におけるステップ S 1 の転送モード設定処理の詳細、及びステップ S 3 ~ S 7 の第 1 乃至第 5 転送モード受信処理の詳細は後の説明で明らかにされる。

【 0 1 0 2 】

なお、本実施形態による転送モード選択処理は、図 6 のステップ S 2 のように送信側で選択して受信側に通知する方法と、図 7 のステップ S 1 のように受信側で選択して送信側に通知する方法とがあり、いずれか一方の転送モード選択処理を実装するが、本実施形態は、送信側で選択して受信側に通知する場合を例にとっている。

【 0 1 0 3 】

(第 1 転送モード)

図 8 は本実施形態における第 1 転送モードによる転送処理の説明図である。図 8 (A) は転送開始前に転送先の受信領域が受信可能となった場合の転送処理であり、図 8 (B) は転送開始後に受信可能となった場合の転送処理である。

【 0 1 0 4 】

ここで第 1 転送モードによる転送処理は、転送先に許可を問い合わせることなく投機的に転送し、不許可の場合には再送要求を受けて転送する処理である。

【 0 1 0 5 】

図 8 (A) の転送開始前に受信許可となっている場合にあっては、受信側 N I C 1 2 - 2 で受信領域について受信許可 1 6 4 を受けた後に、その受信領域に対する送信要求 1 6 6 が送信側 N I C 1 2 - 1 で発生すると、受信側に受信許可の有無を問い合わせることなく転送データから R D M A データパケットを作成し、R D M A 通信 1 6 8 を行う。

【 0 1 0 6 】

受信側 N I C 1 2 - 2 にあっては、R D M A データパケットを受信した際に、転送領域管理表から受信許可を判別し、許可された受信領域に受信パケットのデータ本体を格納する領域格納 1 7 0 を行い、格納終了で a c k パケット 1 7 2 を送信してデータ転送を終了させる。また、受信側 N I C 1 2 - 2 で転送領域管理表を参照して受信許可を認識した際に、2 重受信を防止するため、その受信領域の転送許可を不許可に変更している。

【 0 1 0 7 】

図 8 (B) は転送開始後に受信可能となった場合であり、送信側 N I C 1 2 - 1 で送信要求 1 7 4 が発生すると、受信側に問い合わせることなく投機的に R D M A 通信 1 7 6 を行ってデータパケットを転送する。

【 0 1 0 8 】

受信側 N I C 1 2 - 2 にあっては、R D M A データパケットのヘッダから転送領域管理表の受信領域を参照すると、このとき不許可にあることから、受信パケットを破棄 1 7 8 とする。その後、受信側 N I C 1 2 - 2 で対応する受信領域につき受信許可 1 8 0 が出さ

10

20

30

40

50

れたことが判別されると、転送領域管理表に記録された破棄パケットに関する情報を基に、送信側NIC12-1に対し再送要求制御パケット182を送信する。

【0109】

これを受けて送信側NIC12-1は、転送データからRDMAデータパケットを作成し、RDMA通信183により再送を行い、このとき受信側NIC12-2は受信許可となっていることから、受信してRDMAパケットのデータ本体を受信領域に転送する領域格納184を行い、格納終了でackパケット186を送信してデータ転送を完了する。

【0110】

図9は転送開始前に受信可能となっている場合の第1転送モードによる転送処理のタイムチャートであり、図2を参照して説明すると次のようになる。

10

【0111】

まず送信側処理にあつては、ステップS1でHOST10-1の送信要求部20が送信要求をNIC12-1に発行する。これを受けてステップS2でNIC12-1の送信要求受付部26が送信制御部30に送信処理要求を通知する。続いてステップS3でモード選択制御部34が第1転送モードによる転送を起動するようにモード選択が行われる。

【0112】

続いてステップS4で、送信制御部30はHOST10-1の送信領域22の転送データ24からRDMAデータパケットを生成し、10GbイーサネットMAC38に出力する。このとき送信元の物理アドレスを転送領域管理部28に問い合わせ取得する。続いてステップS5で10GbイーサネットMAC38が送信制御部30から送出されたRDMAデータパケットをネットワークに送信する。

20

【0113】

一方、受信側処理にあつては、ステップS101で10GbイーサネットMAC50がネットワークからRDMAデータパケットを受信し、受信制御部40に引き渡す。受信制御部40はステップS102でRDMAデータパケットのヘッダを解析し、ステップS103で転送領域管理部42に問い合わせ受信領域54の受信許可の有無と物理アドレスを取得する。

【0114】

ステップS104で受信許可を判別すると、ステップS105で受信領域54に受信パケットのデータ本体を転送し、転送完了後にステップS106でackパケットを生成し、ステップS107で10GbイーサネットMAC50からack制御パケットをネットワークに送信する。

30

【0115】

送信側処理にあつては、ステップS6でack制御パケットを10GbイーサネットMAC38で受信し、ステップS7で受信制御部32がパケット解析によりackパケットであることを認識して送信制御部30に通知し、一連のデータ転送処理を正常終了とする。

【0116】

図10は転送開始後に受信可能となった場合の第1転送モードによる転送処理のタイムチャートである。図10において、送信側処理のステップS1～S5及びこれに続く受信側処理のステップS101～S103は、図9と同じである。

40

【0117】

しかしながら、図10の場合には、ステップS103で転送領域管理表を参照した際に受信領域54が不許可であることから、これをステップS104で判別すると、ステップS105で受信パケットを破棄して転送領域管理表に記録する。続いてステップS106で受信領域54の受信許可の有無を監視しており、受信許可が得られるとステップS107に進み、転送領域管理表のパケット破棄の記録に基づき送信制御部44が送制御パケットを生成し、ステップS108で10GbイーサネットMAC50からネットワークに再送通知制御パケットを送信する。

【0118】

50

送信側処理にあつては、ステップS 6で10GbイーサネットMAC 38が制御パケットを受信し、ステップS 7で受信制御部32がパケット解析により再送要求の制御パケットであることを認識して送信制御部30に通知する。送信制御部30は、ステップS 8で再送要求に対応した転送データからRDMAデータパケットを生成し、ステップS 9で10GbイーサネットMAC 38に出力してネットワークにパケットを送信させる。

【0119】

受信側処理にあつては、ステップS 109で10GbイーサネットMAC 50がRDMAデータパケットを受信して受信制御部40に出力し、ステップS 110で受信制御部40がRDMAデータパケットのヘッダを解析し、転送領域管理部42に問い合わせして受信領域54の許可の有無及び物理アドレスを取得し、この場合には受信許可であることから、ステップS 111でホスト10-2の受信領域54に受信パケットから得られたデータ本体を転送する。

10

【0120】

続いてステップS 112で、データ転送が完了したら送信制御部44はackパケットを生成し、10GbイーサネットMAC 50に出力し、ステップS 113でネットワークにack制御パケットを送信する。

【0121】

送信側処理にあつては、ステップS 10で10GbイーサネットMAC 38が制御パケットを受信し、ステップS 11で受信制御部32がパケットを解析し、ackパケットの受信を確認して転送処理を正常終了する。

20

【0122】

図11は第1転送モード送信処理のフローチャートである。第1転送モード送信処理にあつては、ステップS 1でRDMAデータパケットを生成し、ステップS 2でネットワークにパケットを送信する。続いてステップS 3で再送要求制御パケットの受信の有無をチェックしており、受信側で受信領域が許可であれば、ステップS 4でackパケットが受信され、これにより一連の転送処理を終了する。

【0123】

一方、受信側で受信領域が不許可であった場合には、受信領域が許可に変わった際に再送要求制御パケットが送られてくることから、ステップS 3で再送要求制御パケットを受信すると、ステップS 5に進み、再送要求制御パケットの解析で得られた再送要求に対応して、再送用のRDMAデータパケットを作成し、ステップS 6でネットワークにパケットを送信し、ステップS 4でackパケットの受信を待って、一連のデータ転送処理を終了する。

30

【0124】

図12は第1転送モード受信処理のフローチャートである。図12において、第1転送モード受信処理は、ステップS 1でRDMAデータパケットの受信の有無をチェックしており、パケットを受信すると、ステップS 2でRDMAデータパケットを解析し、ステップS 3で転送領域管理表を参照し、物理アドレスと受信領域の許可の有無を取得する。

【0125】

続いてステップS 4で受信領域が受信許可であれば、ステップS 5に進み、受信パケットから得られたデータ本体を許可された受信領域に転送し、ステップS 6で転送完了を示すackパケットを生成して送信し、一連の受信処理を終了する。

40

【0126】

一方、ステップS 4で受信領域が受信不許可であった場合には、ステップS 7で受信パケットを破棄し、ステップS 8で受信パケットを破棄したことを転送領域管理表に記録する。続いてステップS 9で転送領域管理表に対する受信領域の受信許可の有無をチェックしており、受信不許可から受信許可に変わると、ステップS 10に進み、転送領域管理表の受信パケット破棄の記録を基に送信元に対し再送要求制御パケットを作成して送信する。

【0127】

50

この再送要求制御パケットの送信に対し、送信元からはR D M Aデータパケットが再送されてくることから、再送されたR D M AデータパケットをステップS 1で受信すると、ステップS 2 ~ S 6の処理を行って受信処理を完了する。

【 0 1 2 8 】

(第2転送モード)

図13は本実施形態における第2転送モードにおける第2転送モード処理の説明図である。第2転送モードは転送先に問い合わせることなく投機的に転送し、不許可の場合には許可通知を受けて転送を中断し、再送要求を受けて再送する処理である。

【 0 1 2 9 】

図13(A)は転送開始前に受信可能となっている場合の第2転送モードの転送処理である。この場合には受信側N I C 1 2 - 2において、受信許可190となった後にその受信領域に対し受信側N I C 1 2 - 1で送信要求192が発生すると、送信データからR D M Aデータパケットを作成し、R D M A通信194を行う。

【 0 1 3 0 】

受信側N I C 1 2 - 2にあつてはR D M Aデータパケットを受信してヘッダを解析し、受信領域管理表の参照から受信領域の受信許可を判別し、受信データのパケットのデータ本体を許可状態にある受信領域54に転送する領域格納196を実行し、格納終了でa c k制御パケット198を送信しデータ転送を終了する。

【 0 1 3 1 】

図13(B)は転送開始後に受信可能となった場合の第2転送モードによる転送処理である。この場合、送信側N I C 1 2 - 1で送信要求200が発生すると、受信側に問い合わせることなく投機的にR D M A転送206を行う。

【 0 1 3 2 】

受信側N I C 1 2 - 2にあつてはR D M Aデータパケットの受信でヘッダを解析し、転送領域管理表の参照で受信領域の不許可を認識し、送信側N I C 1 2 - 2に対し不許可通知制御パケット202を送信する。この不許可通知制御パケット202を受けて送信側N I C 1 2 - 1はそれまで行っていたR D M A通信206による転送中断204を行う。このとき受信されたデータ本体は破棄208となり、転送領域管理表にデータ破棄が記録される。

【 0 1 3 3 】

受信側N I C 1 2 - 2において受信パケットを破棄208にした後、受信許可210がその受信領域に対し行われると、再送要求制御パケット212を送信する。この再送要求制御パケット212を送信側N I C 1 2 - 1で解析し、再送要求の対象となった転送データについてR D M Aデータパケットを作成し、R D M A再送通信214を行う。再送分の転送が済むと残りの転送データについてR D M A通信216を行う。

【 0 1 3 4 】

このとき受信側N I C 1 2 - 2にあつては受信領域が受信許可210となっていることから再送された受信パケット及びこれに続く受信パケットからデータ本体を取得して受信領域にデータ転送して格納する領域格納218を行い、格納終了でa c kパケット220を送信して一連の処理を終了する。

【 0 1 3 5 】

この第2転送モードにあつては図13(B)のように受信側に問い合わせることなく投機的にR D M A通信を行った場合、受信領域が不許可の場合には不許可通知を返すことで転送を中断させ、不許可通知により破棄される無駄なパケットの喪失を止めることで、第1転送モードの場合に比べネットワーク帯域消費を削減することができる。

【 0 1 3 6 】

図14は転送開始前に受信可能となっている場合の第2転送モードによる転送処理のタイムチャートであり、この場合の転送処理は図9に示した第1転送モードの転送処理の場合と同じである。

【 0 1 3 7 】

10

20

30

40

50

図15は転送開始後に可能となった場合の第2転送モードによる転送処理のタイムチャートであり、図2を参照して説明すると次のようになる。

【0138】

図15において、ステップS1～S5及びこれに続く受信側のステップS101～S105の処理は、図10に示した第1転送モードの場合と同様であるが、第2転送モードにあってはステップS105で受信領域の不許可を認識してパケットを破棄した後、ステップS106で受信制御部40が不許可通知パケットを作成し、ステップS107で10GbイーサネットMAC50がネットワークに制御パケットを送信する。

【0139】

この制御パケットを送信処理側の10GbイーサネットMAC38がステップS6で受信し、ステップS7で送信制御部30がパケットを解析し、不許可通知パケットであることを認識し、ステップS8で転送データからのパケット生成を停止してRDMA通信を中断する。

【0140】

送信側でパケット転送を中断した後については、受信側においてステップS108で受信領域の許可が判別されるとステップS109で送信制御部44が再送制御パケットを生成し、ステップS110で制御パケットをネットワークに送信する。

【0141】

これを受けて送信側処理にあってはステップS9で10GbイーサネットMAC38がネットワークから制御パケットを受信し、ステップS10で受信制御部32がパケットを解析して再送要求を送信制御部30に通知し、通信制御部30はステップS11でRDMAデータパケットを生成し、ステップS12で10GbイーサネットMAC38がネットワークにRDMAデータパケットを再送する。

【0142】

受信側処理にあってはステップS111で10GbイーサネットMAC50が再送されたRDMAデータパケットを受信し、ステップS112で受信制御部40がパケットのヘッダを解析して転送領域管理表から受信領域の受信許可と物理アドレスを取得し、ステップS113でデータ本体をホスト10-2の受信領域54にDMAデータ転送する。

【0143】

続いて図16のステップS114でデータ転送完了に伴い送信制御部44がackパケットを作成し、ステップS115で10GbイーサネットMAC50がackパケットをネットワークに送信する。

【0144】

パケット送信を受けて受信側処理にあってはステップS13で10GbイーサネットMAC38がネットワークからパケットを受信し、ステップS14で受信制御部32が受信パケットを解析してackパケットであることを送信制御部30に通知し、データ転送完了を認識して正常終了とする。

【0145】

図17は第2転送モードの送信処理のフローチャートである。図17において、第2転送モード送信処理はステップS1でRDMAデータパケットを送信要求に基づいて生成した後、ステップS2でネットワークにパケットを送信する。

【0146】

続いてステップS3で受信領域不許可通知の有無をチェックしている。受信側の受信領域が受信許可であれば受信領域不許可通知はないことからステップS4に進み、ackパケットの受信を待って一連の送信処理を終了する。

【0147】

一方、受信側で受信領域が不許可であった場合には、受信領域の不許可通知を送ってこることから、これをステップS3で判別するとステップS5でパケット生成を停止し、RDMA通信を中断する。

10

20

30

40

50

【 0 1 4 8 】

その後、ステップ S 6 で転送先から再送制御パケットが受信されるとステップ S 7 で再送用の R D M A データパケットを生成し、ステップ S 8 でネットワークに送信する。ステップ S 4 で a c k パケットを受信すると転送中断後の再送によるデータ転送が終了したものと一連の処理を終了する。

【 0 1 4 9 】

図 1 8 は第 2 転送モード受信処理のフローチャートである。図 1 8 において、ステップ S 1 で R D M A データパケットの受信の有無をチェックしており、パケットを受信するとステップ S 2 でパケットを解析し、ステップ S 3 で転送領域管理表を参照して物理アドレスと受信領域の受信許可または不許可を取得する。

10

【 0 1 5 0 】

続いてステップ S 4 で受信領域が許可であった場合にはステップ S 5 に進み、受信パケットから得られたデータ本体を許可した受信領域に転送し、転送終了でステップ S 6 に進み、a c k パケットを生成して送信し一連の処理を終了する。なお、受信許可を認識した際には、転送領域管理表の許可を不許可に変更して 2 重受信を防止する。

【 0 1 5 1 】

一方、ステップ S 4 で受信不許可であった場合には、ステップ S 7 で受信パケットを破棄した後、ステップ S 8 で受信パケットの破棄を転送領域管理表に記録し、ステップ S 9 で受信領域の不許可通知を作成して転送元にパケット送信する。

【 0 1 5 2 】

その後、ステップ S 1 0 で受信領域の受信許可が判別されると、ステップ S 1 1 で転送領域管理表のパケット破棄の記録をもとに送信元に再送要求制御パケットを送信する。これに伴い送信元からは R D M A データパケットが再送されてくることから、これをステップ S 1 で受信し、ステップ S 2 ~ S 6 の処理を経て再送された R D M A データパケットについて受信処理を行うことになる。

20

【 0 1 5 3 】

(第 3 転送モード)

図 1 9 は本実施形態における第 3 転送モードによる転送処理の説明図である。第 3 転送モードは転送完了通知である a c k を受け取るまで転送先の許可を問い合わせることなく、投機的に転送データ分の R D M A データパケットを繰り返し転送する処理である。この第 3 転送モードによる処理はネットワークの帯域が十分に余っていることを想定しており、その結果、転送時間短縮のため a c k が受信されるまで何度も R D M A 転送を繰り返す。

30

【 0 1 5 4 】

図 1 9 (A) は転送開始前に受信可能となっている場合の第 3 転送モードの転送処理である。図 1 9 (A) において、送信側 N I C 1 2 - 1 で送信要求 2 3 2 が発生した際に、受信側 N I C 1 2 - 2 には、送信要求 2 3 2 の対象となった受信領域につきそれ以前に受信許可 2 3 0 が出されている。

【 0 1 5 5 】

この状態で転送データから R D M A パケットを生成して、R D M A 一回目通信 2 3 4 を行うと、受信側には受信領域が受信許可となっているため受信パケットの本体データを受信領域に転送して格納完了し、転送完了で a c k 制御パケット 2 3 6 を応答して一連の処理を終了する。尚、送信側は、a c k 制御パケット 2 3 6 を受け取るまで、R D M A 2 回目通信 2 3 5 を続行する。受信側は、R D M A 1 回目通信 2 3 4 を受信領域に転送した時点で、転送領域管理表を不許可とするため、R D M A 2 回目通信 2 3 5 は全て破棄 2 3 7 される。

40

【 0 1 5 6 】

図 1 9 (B) は転送開始後に受信可能となった場合の第 3 転送モードの転送処理である。図 1 9 (B) において、送信側 N I C 1 2 - 1 側で送信要求 2 3 8 が発生すると a c k パケットが受信されるまで転送データ分の R D M A データパケットを作成し、R D M A 1

50

回目通信 240、RDMA 2 回目通信 244、RDMA 3 回目通信 252 というように繰り返しデータ転送を行う。

【0157】

一方、受信側 NIC 12-2 においては RDMA 1 回目通信 240 の際には受信領域が不許可であるため受信パケットを破棄 242 とする。この例では RDMA 2 回目通信 244 のデータ転送を開始した直後に受信許可 248 となり、したがって受信許可 248 となった以降に受信された RDMA 2 回目通信 244 の受信パケットのデータ本体は受信許可となった受信領域に転送されて領域格納 250 となる。しかし、RDMA 2 回目通信 244 の先頭部分の転送データについては破棄 246 となっている。

【0158】

RDMA 2 回目通信 244 が終了すると、続いて RDMA 3 回目通信 252 が行われる。この RDMA 3 回目通信 252 につき通信許可 248 の直前で破棄した RDMA 2 回目通信 244 の先頭部分のデータが受信されるとデータ転送が完了し、この時点で ack 制御パケット 254 を応答し、その結果、送信側 NIC 12-1 においては ack 制御パケット 254 を受信した RDMA 3 回目通信 252 の途中で一連のデータ転送を完了することになる。尚、送信側は、ack 制御パケット 254 を受信するまで、RDMA 3 回目通信 252 を続行する。受信側は、ack 制御パケット 254 送信後の RDMA 3 回目通信 252 を破棄 256 とする。

【0159】

図 20 は第 3 転送モードによる転送処理のタイムチャートであり、図 2 を参照して説明すると、次のようになる。送信側処理においては、ステップ S1 でホスト 10-1 が送信要求を送信側の NIC 12-1 に発行する。これを受けて送信要求受付部 26 は送信制御部 30 にステップ S2 で送信処理要求を通知する。

【0160】

続いてステップ S3 でモード選択制御部 34 が第 3 転送モードによる転送を起動するように指示するモード選択を行っており、これを受けて送信制御部 30 はステップ S4 で RDMA データパケットを生成し、10Gbイーサネット MAC 38 に出力して、ステップ S5 でパケット送信をネットワークに対して行う。またステップ S4 で RDMA データパケットを生成する際には転送領域管理表を参照して転送元の物理アドレスを取得している。

【0161】

受信側処理においては、ステップ S101 で 10Gbイーサネット MAC 50 がネットワークから RDMA データパケットを受信し、ステップ S102 で受信制御部 40 がパケットを解析し、ステップ S103 で転送領域管理表による受信許可及び物理アドレスの取得を転送領域管理部 42 に問い合わせる。

【0162】

この問い合わせに対しステップ S104 で受信許可が判別されると、ステップ S105 で受信パケットのデータ本体を受信許可となっているホスト 10-2 の受信領域 54 に DMA データ転送し、ステップ S106 でデータ転送終了で送信制御部 44 が ack パケットを生成し、ステップ S107 で 10Gbイーサネット MAC 50 によりネットワークに ack 制御パケットを送信する。なお、ステップ S104 で受信許可が判別されると、転送領域管理表の受信許可を受信不許可に変更して 2 重受信を防止する。

【0163】

一方、ステップ S104 で転送領域管理表の参照で得られた受信領域の状態が不許可であった場合には、許可となるまでステップ S104 の待ち処理となり、その間に受信されたパケットについては破棄する。受信領域が許可となればその時点から受信パケットを有効に受信して受信領域にデータ転送する。

【0164】

送信側の処理においては、ステップ S6 で 10Gbイーサネット MAC 38 が ack パケットをネットワークから受信し、受信制御部 32 に受け渡す。受信制御部 32 はステッ

10

20

30

40

50

プ S 7 でパケットを解析し、a c k 受信を確認して送信制御部 3 0 に通知し、転送処理を完了する。

【 0 1 6 5 】

図 2 1 は第 3 転送モードの送信処理のフローチャートである。図 2 1 において、第 3 転送モードの送信処理は、ステップ S 1 で送信要求に伴い R D M A データパケットを作成し、ステップ S 2 でネットワークにパケットを送信する。続いてステップ S 3 で a c k パケットの受信の有無をチェックしており、a c k パケットを受信するまでステップ S 1 ~ S 2 の R D M A データパケットの生成とパケット送信を繰り返す。

【 0 1 6 6 】

図 2 2 は第 3 転送モードの受信処理のフローチャートである。図 2 2 において、第 3 転送モード受信処理は、ステップ S 1 で R D M A データパケットの受信を判別するとステップ S 2 でパケットを解析し、ステップ S 3 で転送領域管理表を参照して物理アドレスと受信領域の許可の有無を取得する。

10

【 0 1 6 7 】

続いてステップ S 4 で受信領域が受信許可であればステップ S 6 で受信パケットのデータ本体を許可された受信領域に転送し、データ転送が終了するとステップ S 7 で a c k パケットを作成しネットワークに送信する。

【 0 1 6 8 】

一方、ステップ S 4 で受信領域が不許可であった場合にはステップ S 5 で受信パケットを破棄し、ステップ S 4 で受信許可となるまでステップ S 1 ~ S 5 の処理を繰り返す。受信パケットを破棄している途中に受信許可となれば、ステップ S 6 で受信パケットのデータ本体を受信領域に転送し、転送完了でステップ S 7 において a c k パケットを生成し、ネットワークに送信し、一連の処理を終了する。

20

【 0 1 6 9 】

(第 4 転送モード)

図 2 3 は本実施形態における第 4 転送モードによる転送処理のタイムチャートである。本実施形態における第 4 転送モードは、転送先に許可を問い合わせることなく投機的に転送し、不許可の場合には一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する処理である。

【 0 1 7 0 】

30

この第 4 転送モードにおける転送処理は、受信側のホストにおいて十分にメモリがあまっている場合、受信領域が不許可であったときに受信パケットから得られたデータ本体を破棄せず一時バッファに保存し、受信領域が受信許可となった後に一時バッファの保存データを受信領域に転送し、これによって第 1 転送モードや第 2 転送モードにおける受信領域が不許可の場合の再送を不要とし、また第 3 転送モードのような受信完了まで繰り返しデータ転送を行うことによるネットワークの無駄な使用を回避できる。

【 0 1 7 1 】

第 4 転送モードにおける送信側の処理は、ホスト側からの送信要求により R D M A 起動可能となったら受信側に受信領域の受信許可の有無を問い合わせることなく投機的に R D M A 通信を開始し、a c k パケットが受信されれば送信を完了する。

40

【 0 1 7 2 】

一方、第 4 転送モードにおける受信側の処理は R D M A パケットを受信すると転送領域管理表を参照して受信領域が受信許可である場合には、受信パケットのデータ本体を受信領域に転送し格納する。このとき 2 重受信を防止するため受信領域の転送許可を不許可に変更する。受信領域の転送処理が完了すると a c k パケットを返送する。

【 0 1 7 3 】

一方、R D M A データパケット受信した際に受信領域が不許可であった場合には、予め登録されている一時バッファとして機能するバッファ領域に受信パケットのデータ本体を転送して保存する。このとき転送領域管理表にデータを保存した一時バッファのバッファ領域のアドレスを記憶する。

50

【 0 1 7 4 】

その後、受信領域に受信許可が得られると受信領域管理表に記憶したバッファ保存に対する情報に基づきバッファ領域の保存データをホスト側の処理として受信領域に移動し、移動完了で a c k パケットを返信する。

【 0 1 7 5 】

このような第 4 転送モードの処理を図 2 3 のタイムチャートにつき、図 2 を参照して説明すると次のようになる。

【 0 1 7 6 】

図 2 3 において、送信側処理にあつてはステップ S 1 でホスト 1 0 - 1 が送信要求を送信側 N I C 1 2 - 1 に発行すると、ステップ S 2 で送信要求受付部 2 6 が送信制御部 3 0 に対し送信処理要求を通知する。続いてステップ S 3 でモード選択制御部 3 4 が、この場合には第 4 転送モードを選択し、第 4 転送モードで転送を起動するように指示している。

【 0 1 7 7 】

送信制御部 3 0 はステップ S 4 で転送データから R D M A データパケットを生成し、ステップ S 5 で 1 0 G b イーサネット M A C 3 8 によりネットワークにパケットを送信する。このとき送信制御部 3 0 は転送領域管理表を参照して物理アドレスを取得し、物理アドレスから転送データを読み出してパケット転送している。

【 0 1 7 8 】

受信側処理にあつては、ステップ S 1 0 1 で 1 0 G b イーサネット M A C 5 0 が R D M A パケットをネットワークから受信して受信制御部 4 0 に受け渡す。受信制御部 4 0 はステップ S 1 0 2 で R D M A パケットを解析し、受信領域にデータ本体を転送する。ステップ S 1 0 3 で転送領域管理部 4 2 に問い合わせた転送領域管理表から受信領域の許可の有無と物理アドレスを取得する。

【 0 1 7 9 】

続いてステップ S 1 0 4 で受信領域が不許可であればステップ S 1 0 5 に進み、一時バッファ 5 6 に受信パケットのデータ本体を転送して保存する。このときデータ領域管理表に受信パケットのデータ本体を一時バッファ 5 6 に転送データ 5 8 として保存したことを記憶する。

【 0 1 8 0 】

続いてステップ S 1 0 6 で受信領域につき受信許可が得られると、ステップ S 1 0 7 に進み、転送領域管理表の一時バッファ 5 6 の記憶に基づきホスト 1 0 - 2 側の転送制御部 6 0 が一時バッファ 5 6 から転送データ 5 8 を読み出し、受信領域 5 4 にデータを移動させる。

【 0 1 8 1 】

データ移動が済むとステップ S 1 0 8 で送信制御部 4 4 は a c k パケットを生成し、ステップ S 1 0 9 で 1 0 G b イーサネット M A C 5 0 がネットワークにパケットを送信する。この a c k パケットは受信側処理におけるステップ S 6 で 1 0 G b イーサネット M A C 3 8 が受信して受信制御部 3 2 に受け渡し、受信制御部 3 2 はステップ S 7 で a c k パケットであることを認識して送信制御部 3 0 に通知し、データ転送処理を完了する。

【 0 1 8 2 】

図 2 4 は第 4 転送モードの送信処理のフローチャートである。図 2 4 において、第 4 転送モードの送信処理にあつては、ステップ S 1 で R D M A データパケットを送信要求に基づき生成し、ステップ S 2 でパケットをネットワークに送信する。続いてステップ S 3 で a c k パケットの送信を判別しており、a c k パケットの受信を判別すると一連の送信処理を完了する。

【 0 1 8 3 】

図 2 5 は第 4 転送モード受信処理のフローチャートである。図 2 5 において、第 4 転送モード受信処理は、ステップ S 1 で R D M A データパケットの受信の有無をチェックしており、パケットを受信するとステップ S 2 でパケットを解析し、ステップ S 3 で転送領域管理表を参照して物理アドレスと受信領域の受信許可の有無を取得する。

10

20

30

40

50

【 0 1 8 4 】

続いてステップ S 4 で受信領域が受信許可であればステップ S 5 で受信パケットのデータ本体を許可された受信領域に転送する。この場合、2重受信を防止するため転送領域管理表の受信許可を受信不許可に変更する。受信領域に対するデータ転送を終了するとステップ S 6 で a c k パケットを生成して、転送元に送信する。

【 0 1 8 5 】

一方、ステップ S 4 で受信領域が受信不許可であった場合にはステップ S 7 に進み、受信領域管理表に一時バッファへの転送を記憶して受信パケットのデータ本体を一時バッファに転送保存する。続いてステップ S 8 で受信領域の受信許可の有無をチェックしており、受信許可が判別されると、ステップ S 9 でそのとき受信された受信パケットのデータ本体を受信許可となった受信領域に転送して格納する。

10

【 0 1 8 6 】

またステップ S 1 0 で転送領域管理表の記録に基づき一時バッファから保存データを読み出し、受信許可となった受信領域にデータを移動する。受信領域に対する受信パケットのデータ本体の転送及びまたは一時バッファからのデータ移動によりデータ転送が完了すると、ステップ S 6 で a c k パケットを生成して転送元に送信することになる。

【 0 1 8 7 】

図 2 6 は第 4 転送モードの受信側における一時バッファから受信領域へのデータ移動の説明図である。受信側 N I C 1 2 - 2 において、第 4 転送モードによる受信処理で受信領域が受信不許可になった場合には、一時バッファ 5 6 に受信パケットのデータ本体を転送して保存する。

20

【 0 1 8 8 】

ここで一時バッファ 5 6 には例えば 5 ページの領域に分けて保存され、1 ページ目 2 6 0 - 1 の途中から 5 ページ目 2 6 0 - 5 の途中まで斜線部で示すように転送データ 5 8 が保存されている。このような一時バッファ 5 6 に対する転送データ 5 8 の保存状態において、受信側 N I C 1 2 - 2 で受信領域の受信許可が認識されると、一時バッファ 5 6 の転送データ 5 8 を受信領域 5 4 に移動する処理を行う。

【 0 1 8 9 】

この移動処理は、先頭となる第 1 ページ 2 6 0 - 1 と末尾となる第 5 ページ 2 6 0 - 5 についてはコピー処理を行うが、その間の 2 ページ目 2 6 0 - 2 ~ 4 ページ目 2 6 0 - 4 についてはページ全領域に転送データが格納されていることから、ページ単位でアドレス範囲をマッピングする処理を行い、このページ単位のマッピング処理によりコピー処理を不要とし、コピー処理のオーバーヘッドを削減する。

30

【 0 1 9 0 】

図 2 7 はアドレスマップ変換によるデータ移動の説明図である。図 2 7 において、受信領域 5 4 及び一時バッファ 5 6 は仮想アドレス（論理アドレス）で管理されており、例えば一時バッファ 5 6 の仮想アドレス $0 \times 8 0 0 0$ に転送データ 2 7 0 を保存しており、これはアドレス変換表 2 6 2 により物理メモリ 2 6 4 のアドレス $0 \times 6 0 0 0$ に転送データ 2 8 0 を格納していることを示している。

【 0 1 9 1 】

ここで、一時バッファ 5 6 の転送データ 2 7 0 を受信領域 5 4 の仮想アドレス $0 \times 1 0 0 0$ に移動する場合を考える。移動前において仮想アドレス $0 \times 1 0 0 0$ は物理メモリ 2 6 4 のアドレス $0 \times 3 0 0 0$ に対応している。このデータ移動前のアドレス変換表 2 6 2 は図 2 8 (A) の内容を持っている。

40

【 0 1 9 2 】

そこで仮想アドレス $0 \times 8 0 0 0$ の転送データ 2 7 0 を受信領域 5 4 の仮想アドレス $0 \times 1 0 0 0$ に移動するためにはアドレス変換表 2 6 2 における移動先となる仮想アドレス $0 \times 1 0 0 0$ の物理アドレス $0 \times 3 0 0 0$ を一時バッファ 5 6 の転送データ 2 7 0 に対応した物理アドレス $0 \times 6 0 0 0$ に書きかえる再マップを行う。即ち図 2 8 (A) のアドレス変換表 2 6 2 における仮想アドレス $0 \times 1 0 0 0$ に対応した物理アドレス $0 \times 3 0 0 0$

50

を図 28 (B) のように 0×6000 に書き替える。

【 0 1 9 3 】

その結果、図 27 の一時バッファ 56 の仮想アドレス 0×8000 の転送データ 270 に対応した物理メモリ 264 の転送データ 268 は、アドレス変換表 262 によるマッピングにより受信領域 54 の仮想アドレス 0×1000 にマッピングされ、これにより一時バッファ 56 から転送データ 270 を受信プロセスの仮想アドレス 0×1000 に移動して転送データ 272 とすることができる。

【 0 1 9 4 】

尚、一時バッファ 56 の仮想アドレス 0×8000 については図 28 (B) のように物理アドレス 0×3000 に当てており、この場合、物理アドレス 0×3000 は空き領域となる。

10

【 0 1 9 5 】

(第 5 転送モード)

図 29 は本実施形態における第 5 転送モードによる転送処理の説明図である。ここで第 5 転送モードは、転送先に許可を問い合わせ、許可通知を受けて転送する処理であり、これは従来の RDMA 転送そのものであるが、本実施形態にあつては最適転送モードの選択の中に前述した本実施形態に固有な第 1 乃至第 4 転送モードに加え、従来転送モードを第 5 転送モードとして加えている。

【 0 1 9 6 】

図 29 (A) は転送開始前に受信可能となった場合の第 5 転送モードの転送処理である。図 29 (A) において、送信側 NIC 12 - 1 で送信要求 282 が発生すると受信許可確認パケット 284 を受信側に送信する。このときそれ以前に受信許可 280 が得られるため、受信側 NIC 12 - 2 は許可パケット 286 を返信する。

20

【 0 1 9 7 】

これを受けて送信側 NIC 12 - 1 は転送データから RDMA データパケットを生成し、RDMA 通信 288 を行う。受信側 NIC 12 - 2 は受信領域が受信許可であることから、受信した RDMA データパケットから得られたデータ本体を受信領域に転送して格納する領域格納 290 を行い、格納終了で ack パケット 292 を送信してデータ転送を完了する。

【 0 1 9 8 】

30

図 29 (B) は転送開始後に受信許可となった場合である。送信側 NIC 12 - 1 は送信要求 294 が発生すると、受信側に対し、受信許可確認パケット 296 を送信する。このとき受信側 NIC 12 - 2 には受信領域の受信許可が得られていないことから、受信許可を待つ。その後、受信許可 298 が得られると許可パケット 300 を送信する。

【 0 1 9 9 】

この許可パケット 300 を受けて送信側 NIC 12 - 1 は RDMA 通信 302 を行い、受信側で受信パケットから得られたデータ本体を受信領域に転送して格納して領域格納 304 を終了すると ack パケット 306 を送信してデータ転送を完了する。

【 0 2 0 0 】

このように第 5 転送モードにあつては、送信要求が発生した際に受信側に受信領域の受信許可の有無を問い合わせ、受信許可の通知が得られたときに RDMA 転送を行っており、この受信領域の許可の有無を問い合わせる通信を必要とする分だけ第 1 乃至第 4 転送モードに比べるとレイテンシが大きくなる。

40

【 0 2 0 1 】

図 30 は第 5 転送モードによる転送処理のタイムチャートであり、図 2 を参照して説明すると次のようになる。送信側処理にあつては、ステップ S1 でホスト 10 - 1 が受信要求を NIC 12 - 1 に対し発行する。これを受けてステップ S2 で送信要求受付部 26 が送信制御部 30 に送信処理要求を通知する。

【 0 2 0 2 】

続いてステップ S3 でモード選択制御部 34 が、この場合には第 5 転送モードを選択し

50

て転送処理を起動するように指示している。続いてステップS4で送信制御部30は受信許可確認パケットを作成し、ステップS5で10GbイーサネットMAC38からネットワークにパケットを送信する。

【0203】

受信側処理にあつては、ステップS101でネットワークから10GbイーサネットMAC50がパケットを受信して受信制御部40に受け渡し、受信制御部40はステップS102でパケットを解析し、受信許可確認パケットであることからステップS103で転送領域管理部40に問い合わせ、転送領域管理表から受信領域の受信許可の有無を取得する。

【0204】

ステップS104で受信許可が得られると、ステップS105で受信制御部40は受信許可パケットを生成し、ステップS106で10GbイーサネットMAC50がネットワークにパケットを送信する。

【0205】

送信側の処理にあつては、ステップS6で10GbイーサネットMAC38がパケットを受信し、ステップS7で受信制御部32でパケットを解析して受信許可を送信制御部30に通知する。これを受けて送信制御部30はステップS8でRDMAデータパケットを生成し、ステップS9で10GbイーサネットMAC38がネットワークにパケットを送信する。

【0206】

受信側処理にあつては、ステップS107で10GbイーサネットMAC50がパケットを受信して受信制御部40に受け渡し、ステップS108で受信制御部40がパケットを解析し、転送領域管理部42に問い合わせして受信領域の受信許可の有無と物理アドレスを取得する。

【0207】

このとき受信領域は受信許可になっていることから、ステップS109で受信したRDMAデータパケットから得られたデータ本体をホスト10-2の受信領域54に転送する。このデータ転送が終了するとステップS110で送信制御部44がackパケットを生成し、10GbイーサネットMAC50からステップS111でネットワークにパケットを送信する。

【0208】

送信側処理にあつてはステップS10で10GbイーサネットMAC38がackパケットを受信して受信制御部32に引渡し、ステップS11で受信制御部32がパケットを解析してack受信を確認し、送信制御部30に通知して転送処理を完了する。

【0209】

図31は第5転送モード処理のフローチャートである。図31において、第5転送モード送信処理はステップS1で送信要求を受けると受信許可確認パケットを生成し、ステップS2でパケットをネットワークに送信する。

【0210】

続いてステップS3で転送先からの受信許可パケットの受信の有無をチェックしており、パケットを受信するとステップS4でRDMAデータパケット生成し、ステップS5でネットワークにパケットを送信する。そしてステップS6でackパケットの受信を確認すると一連の送信処理を終了する。

【0211】

図32は第5転送モードの受信処理のフローチャートである。第5転送モード受信処理は、ステップS1で受信許可確認パケットの受信の有無をチェックしており、これを受信するとステップS2で転送領域管理表を参照して受信領域の受信許可の有無を取得する。ステップS3で受信許可を判別するとステップS7で受信許可パケットを転送先に送信する。

【0212】

10

20

30

40

50

これに対し転送元からはR D M Aデータパケットが送られてくることからステップS 8でこの受信を判別し、ステップS 9でR D M Aデータパケットを解析し、転送領域管理表の参照で物理アドレスと受信領域の受信許可を取得し、ステップS 10でデータ本体を受信許可となっている転送領域に転送する。データ転送が終了するとステップS 11でackパケットを生成してネットワークに送信して一連の受信処理を終了する。

【0213】

(転送モードの選択)

本実施形態にあつては、R D M A転送を行う際に、ネットワーク負荷及び受信側のメモリ使用率に基づいて第1乃至第5転送モードの中から最適な転送モードを選択してR D M A転送を行っている。そこで、本実施形態における第1乃至第3転送モードを従来方式である第5転送モードと比較する。

10

【0214】

図33は転送開始前に受信可能となっている場合の第1乃至第3転送モードの転送完了時間を、従来方式である第5転送モードと対比して示した説明図である。図33(A)は従来方式の第5転送モードであり、図33(B)は第1転送モード、図33(C)は第2転送モード及び図33(D)は第3転送モードである。

【0215】

ここで第1乃至第3転送モード及び従来方式である第5転送モードの転送完了時間を T_1 、 T_2 、 T_3 、 T_5 とし、また転送完了時間を求めるためのパラメータ310として、片道レイテンシ、スループットの逆数 k 、データ長 L を設定する。なお、スループットの逆数 k にデータ長 L を掛け合わせた(kL)は、図33(A)~(D)における通常R D M A転送314、316、318、320の伝送時間となる。

20

【0216】

図33(A)~(D)においては、送信要求が発生する前に受信側で受信可能312となっており、このため第1乃至第3転送モードのいずれについても転送先に損失させることなく転送を開始して、通常R D M A転送316、318、320を行っている。

【0217】

一方、従来方式である図33(A)の第5転送モードにあつては、転送先に受信許可確認を行って許可通知を得られた後に通常R D M A314を行っている。なお図33(D)の第3転送モードにあつては、転送側からackが得られるまで繰り返しR D M A転送を行っているため、転送先からackを受け取るまでの転送が損失する転送322となっている。

30

【0218】

このように転送開始前に受信可能312となっている場合の第1乃至第3転送モード及び従来方式である第5転送モードの転送完了時間は次のようになる。

第1転送モード $T_1 = 2 + kL$

第2転送モード $T_2 = 2 + kL$

第3転送モード $T_3 = 2 + kL$

第5転送モード $T_5 = 4 + kL$

この場合には、本実施形態における第1乃至第3転送モードのいずれについても転送先に受信許可確認を行わずに投機的に転送しているため、受信許可確認のための通信時間(2)だけ第5転送モードに対し削減することができる。

40

【0219】

図34は図33における第5転送モード、第1乃至第3転送モードのそれぞれにおける転送完了時間 T_5 、 $T_1 \sim T_3$ を計算するための詳細を示している。

【0220】

次に転送先の受信許可が転送開始後に発行されている場合について、第1乃至第3転送モードの転送完了時間を従来方式である第5転送モードと対比して説明する。

【0221】

図35(A)~(D)は、パラメータ324に示すように、送信要求から受信許可が発

50

行されるまでの待ち時間 W が

$$W \quad k L$$

の場合である。この場合の第1乃至第3転送モードの転送完了時間 $T_1 \sim T_3$ は次のようになる。

$$\text{第1転送モード} \quad T_1 = W + 4 \quad + k L$$

$$\text{第2転送モード} \quad T_2 = W + 4 \quad + k L$$

$$\text{第3転送モード} \quad T_3 = W + 2 \quad + k L$$

また従来方式である第5転送モードの転送完了時間 T_5 は

$$\text{第5転送モード} \quad T_5 = W + 4 \quad + k L \text{ (固定)}$$

となる。

10

【0222】

図35(A)～(D)においては、RDMA通信328, 342, 344及び再送RDMA通信332, 336は、受信領域へ転送され、RDMA通信330, 334, 338, 340, 346は、転送先で受信パッケージが破棄される。

【0223】

図36(A)～(D)は、図35(A)～(D)に対応して転送完了時間 T_5 , $T_1 \sim T_3$ のそれぞれの計算内容の詳細を示している。

【0224】

図37は、パラメータ348に示すように、送信要求から受信許可が発行されるまでの待ち時間 W を

20

$$k L \quad W \quad k L - 2$$

とした場合である。この場合の図37(B)～(D)の第1乃至第3転送モードの転送完了時間は次のようになる。

$$\text{第1転送モード} \quad T_1 = 2W + 4$$

$$\text{第2転送モード} \quad T_2 = W + 4 \quad + k L$$

$$\text{第3転送モード} \quad T_3 = W + 2 \quad + k L$$

また図37(A)の従来方式である第5転送モードは、図35の場合と同様、

$$\text{第5転送モード} \quad T_5 = W + 4 \quad + k L \text{ (固定)}$$

であり、待ち時間 W が短くなっただけである。

【0225】

30

図37(A)～(D)においては、RDMA通信352, 355, 362, 366, 368及び再送RDMA通信356, 360は、受信領域へ転送され、RDMA通信354, 358, 364, 370, 372は、転送先で受信パッケージが破棄される。

図38(A)～(D)は図37(A)～(D)における転送完了時間 T_5 , $T_1 \sim T_3$ の計算内容を示している。

【0226】

ここで図38(B)の第1転送モードにあつては、受信可能350となるまでの待ち時間 W と受信可能350が発行された後の通常RDMA転送の転送時間 $k L$ は、転送先で破棄された損失分のデータを同一速度で転送するため、データ転送時間は待ち時間 W と等しくなり、したがって

40

$$T_1 = W + 4 \quad + k L = 2W + 4$$

となる。

【0227】

図39は、パラメータ374に示すように、送信要求から受信許可が発行されるまでの待ち時間 W が

$$k L - 2 \quad W \quad 2$$

の場合の第1乃至第3転送モードの転送完了時間を従来方式である第5転送モードと対比して示している。

【0228】

この場合の第1乃至第3転送モードの転送完了時間は次のようになる。

50

第1転送モード $T_1 = W + 2 + kL$

第2転送モード $T_2 = W + 4 + kL$

第3転送モード $T_3 = W + 2 + kL$

なお第5転送モードの転送完了時間は

第5転送モード $T_4 = W + 4 + kL$ (固定)

であり、待ち時間 W が短くなっただけである。

図39(A)～(D)においては、RDMA通信378, 382, 390, 394, 396及び再送RDMA通信384, 388は、受信領域へ転送され、RDMA通信380, 386, 392, 398は、転送先で受信パケットが破棄される。

【0229】

図40(A)～(D)は、図39(A)～(D)に対応したそれぞれの転送完了時間の計算内容を示している。

【0230】

図41は、パラメータ400に示すように、送信要求から受信許可が発行されるまでの待ち時間 W が

$2W$

の場合であり、この場合の第1乃至第3転送モードの転送完了時間は次のようになる。

第1転送モード $T_1 = W + 2 + kL$

第2転送モード $T_2 = 2W + 2 + kL$

第3転送モード $T_3 = W + 2 + kL$

なお第5転送モードの転送完了時間は

第5転送モード $T_4 = W + 4 + kL$ (固定)

であり、待ち時間 W が短くなっただけである。

【0231】

また図41(A)～(D)においては、RDMA転送404, 408, 414, 416, 422, 424及び再送RDMA転送410, 418は、受信領域へ転送され、RDMA通信406, 412, 420, 426は、転送先で受信パケットが破棄される。

【0232】

図42(A)～(D)は、図41(A)～(D)の転送完了時間の計算内容を示している。

【0233】

ここで図42(C)の第2転送モードの転送完了時間 T_2 は

$T_2 = 2W + 2 + kL$

であるが、これは図43に示すようにして算出されている。

【0234】

図43において、第2転送モードにおける転送完了時間 T_2 は

$T_2 = W + 4 + kL -$

である。ここで $2W$ は、図43より

$= 2W - W$

であり、これを代入すると

$T_2 = W + 4 + kL - (2W - W) = 2W + 2 + kL$

となる。

【0235】

以上の受信許可が転送開始後に発行された場合の待ち時間 W の相違に対する第1乃至第3転送モードの転送完了時間 $T_1 \sim T_3$ をまとめると次のようになる。

【0236】

(1) $W = kL$ の場合

第1転送モード $T_1 = W + 4 + kL$

第2転送モード $T_2 = W + 4 + kL$

第3転送モード $T_3 = W + 2 + kL$

10

20

30

40

50

【 0 2 3 7 】

(2) $k L \leq W \leq k L - 2$ の場合

第 1 転送モード $T 1 = 2 W + 4$
 第 2 転送モード $T 2 = W + 4 + k L$
 第 3 転送モード $T 3 = W + 2 + k L$

【 0 2 3 8 】

(3) $k L - 2 \leq W \leq 2$ の場合

第 1 転送モード $T 1 = W + 2 + k L$
 第 2 転送モード $T 2 = W + 4 + k L$
 第 3 転送モード $T 3 = W + 2 + k L$

10

【 0 2 3 9 】

(4) $2 \leq W$ の場合

第 1 転送モード $T 1 = W + 2 + k L$
 第 2 転送モード $T 2 = 2 W + 2 + k L$
 第 3 転送モード $T 3 = W + 2 + k L$

【 0 2 4 0 】

この第 1 乃至第 3 転送モードにおける転送完了時間 $T 1 \sim T 3$ の関係から、単純に転送時間を短縮するという観点から見ると第 3 転送モードが最も優れていることが分かる。しかしながら、第 3 転送モードにあっては、ack が得られるまで繰り返し R D M A 転送を行うことから、最もネットワーク帯域を消費するという問題がある。このため、ネットワーク負荷の状態に応じて最適な転送モードを選択する必要がある。

20

【 0 2 4 1 】

本実施形態にあっては、図 2 に示したように、ネットワーク情報は送信側のネットワーク負荷情報収集部 3 6 と受信側のネットワーク負荷情報収集部 4 8 の両方で収集可能である。

したがって

(1) 送信側でネットワーク負荷情報を収集する場合、ネットワーク収集情報に基づいて送信時に転送モードを選択する方法、

(2) 受信側で収集する場合は、収集情報に基づいて選択した転送モードを送信側に伝え、これに基づいて送信側で指定された伝送方式で送信処理を行う方法、

のいずれかをとることができる。本実施形態は (1) の送信側で転送モードを選択している。

30

【 0 2 4 2 】

更に本実施形態における第 1 乃至第 3 転送方式は、再送によるネットワーク帯域の増加があるものの、受信側のホスト上の一時バッファを必要としないのに対し、本実施形態の第 4 転送モードにあっては、再送によるネットワーク帯域の増加はないものの受信側のホスト上の一時バッファを必要とする。

【 0 2 4 3 】

そこで送信側で転送方式を決定する方法において、図 2 のように受信側の N I C 1 2 - 2 に設けているメモリ使用率情報収集部 4 9 でホスト 1 0 - 2 のメモリ使用率を収集し、これを送信側に通知し、メモリ使用率に基づいて転送モードを選択する。

40

【 0 2 4 4 】

即ち、メモリ使用率が少なく十分にメモリがある場合には一時バッファを使用する第 4 転送モードを選択し、メモリ使用率が大きく十分にメモリがない場合には第 1 乃至第 3 転送モードを選択し、これによって効率よく転送処理を行う。

【 0 2 4 5 】

更に、ネットワーク負荷が高く且つ受信側のホストのメモリ使用率が高い場合を考えると、第 1 乃至第 3 転送モードは無駄となるパケットの転送を伴い、ネットワーク負荷が高くなるため、好ましくはない。また一時バッファを使用する第 4 転送モードも、受信側のホストメモリ使用量が多いため好ましくない。このような場合には、従来方式である第 5

50

転送モードが最適な方式となる。

【0246】

そこで本実施形態にあつては、図2の受信側のNIC12-2のメモリ使用率情報収集部49がホスト10-2のメモリ使用率情報を収集してモード選択制御部46に通知し、モード選択制御部46においてメモリ使用率が少なければ第4転送モードを選択し、一方、メモリ使用率が多く第4転送モードを選択できない場合には、第4転送モードを選択できない点を送信側のNIC12-1のモード選択制御部34に通知する。

【0247】

これを受けて送信側のモード選択制御部34は、ネットワーク負荷情報収集部36から得られたネットワーク負荷情報に基づき、第1乃至第3転送モードあるいは従来方式である第5転送モードから最適な転送モードを選択する。なお、メモリ使用率を送信側のモード選択制御部34に通知し、モード選択制御部34で第1乃至第5転送モードのすべてについて選択制御しても良い。

10

【0248】

図44はネットワーク負荷情報と受信側のメモリ使用率に基づく転送モード選択のための転送モード管理表450の説明図である。図44において、転送モード管理表450は、ネットワーク負荷については、小さい順に3つの閾値 w_1 、 w_2 、 w_3 を設定し、ネットワーク負荷を4つの領域に分けている。一方、メモリ使用量については閾値 m を設定し、 m 未満と m 以上の2つの領域に分けている。

【0249】

この転送モード管理表450に基づき、次のようにして転送モードの選択を行う。

(1)メモリ使用率が閾値 m 未満の場合は、受信不許可の場合に一時バッファに格納する第4転送モードを選択する。

(2)メモリ使用率が閾値 m 以上の場合は、ネットワーク負荷に基づいて転送モードを選択する。

20

【0250】

またネットワーク負荷に基づく転送モードの選択は次のようになる。

(1)ネットワーク負荷が w_1 未満の場合は第3転送モードを選択する。

(2)ネットワーク負荷が w_1 以上 w_2 未満の場合は第1転送モードを選択する。

(3)ネットワーク負荷が w_2 以上 w_3 未満の場合は第2転送モードを選択する。

(4)ネットワーク負荷が w_3 以上の場合は従来方式である第5転送モードを選択する。

30

【0251】

ここでメモリ使用率については、受信側のホストにおける一時バッファとして使用可能なメモリ量とデータ転送量から閾値を求める。ネットワーク負荷に関しては、ネットワーク負荷情報を収集し、ネットワーク負荷実際の性能値を集計し、閾値の妥当性を評価することになる。

【0252】

図45は図44の転送モード管理表450に基づく図6のステップS2の送信側での転送モード選択処理の詳細を示したフローチャートであり、図2を参照して説明すると次のようになる。

40

【0253】

図45において、転送モード選択処理は、ステップS1で、モード選択制御部34はネットワーク負荷情報収集部36で収集されたネットワーク負荷情報と、このとき受信側のNIC12-2のメモリ使用率情報収集部49で収集されて転送されてきたメモリ使用率を読み込む。続いてステップS2でメモリ使用率は閾値 m 以上か否かチェックする。 m 未満であればステップS3に進み、受信側のホスト10-2の一時バッファ56を使用する第4転送モードを設定する。

【0254】

ステップS2でメモリ使用率が閾値 m 以上の場合には、ステップS4に進み、ネットワーク負荷が閾値 w_1 未満か否かチェックする。閾値 w_1 未満であればステップS5に進み

50

、転送データをackパケットが得られるまで繰り返し転送する第3転送モードを設定する。

【0255】

またステップS6でネットワーク負荷が閾値w1以上で 閾値w2未満の場合には、ステップS7で第1転送モードを設定する。またステップS8でネットワーク負荷が閾値w2以上で閾値w3未満の場合には、ステップS9で第2転送モードを設定する。更にステップS8でネットワーク負荷が閾値w3以上となった場合には、ステップS10に進み、従来方式である第5転送モードを設定する。

【0256】

なお図45の転送モード設定処理にあつては、送信側のNIC12-1のモード選択制御部34で転送モードを選択して選択指示を与えているが、第1転送モード、第2転送モード及び第4転送モードの選択については、受信側NIC12-2のモード選択制御部46においてモード選択を行う構成とすることもできる。この場合には、受信側の統計情報から得られたモード選択を送信側に通知する手順を省略することができる。

10

【0257】

なお本実施形態は、受信機側のメモリ使用率に基づき、第1乃至第4転送モード及び従来方式である第5転送モードの中から最適な転送モードを選択してRDMA転送を行う場合を例にとっているが、本実施形態における第1乃至第4転送モードのいずれか1つを固定的に選択してRDMA転送する構成とすることもできる。

【0258】

20

また本実施形態における第1乃至第4転送モードの少なくともいずれか1つと従来方式である第5転送モードにつき、ネットワーク負荷及び受信側のメモリ使用率から最適な転送モードを選択してRDMA転送する構成としてもよい。

【0259】

また本実施形態は、受信側及び送信側に設けたNICにおける送信部及び受信部として実行するプログラムを提供するものであり、このプログラムは図6，図7，図11，図12，図17，図18，図21，図22，図24，図25，図31，図32、更に図45のフローチャートに示した内容を備える。

【0260】

また本発明は、本実施形態のRDMA通信プログラムを記憶したコンピュータ読取可能な記録媒体を提供するものであり、この記録媒体とはCD-ROM、フロッピー(R)ディスク、DVDディスク、光磁気ディスク、ICカードなどの可搬型記憶媒体や、コンピュータシステムの内外に備えられたハードディスクドライブなどの記憶装置の他、回線を介してプログラムを保持するデータベース、あるいは他のコンピュータシステム並びにデータベースや、更に回線上の伝送媒体を含むものである。

30

【0261】

また上記の実施形態は10GbイーサネットをサポートするNICを例に取るものであったが、これに限定されず、同等のRDMA通信をサポートするインタフェースネットワークカードや通信機器であれば、そのまま本発明を適用することができる。

【0262】

40

また本発明は本実施形態に限定されず、その目的と利点を損なうことのない適宜の変形を含み、更に本実施形態に示した数値による限定は受けない。

【0263】

ここで本発明の特徴をまとめて列挙すると次の付記のようになる。

(付記)

(付記1)

RDMAパケットを作成して送信する送信部と、RDMAパケットを受信する受信部とを備えたRDMA通信をサポートする通信装置に於いて、前記送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データ

50

から R D M A パケットを作成して投機的に送信するパケット送信部と、
 転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送部と、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
 を備え、
 前記受信部は、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、
 前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
 を備えたことを特徴とする通信装置。(1)

10

【 0 2 6 4 】

(付記 2)

R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた R D M A 通信をサポートする通信装置に於いて、
 前記送信部は、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信部と、
 転送先から受信不許可通知を受信した際に、前記パケット送信部によるパケット転送を中断する転送中断部と、
 転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送部と、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
 を備え、
 前記受信部は、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、
 前記パケット破棄部で前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知部と、
 前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
 を備えたことを特徴とする通信装置。(2)

20

30

40

50

【 0 2 6 5 】

(付 記 3)

R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた R D M A 通信をサポートする通信装置に於いて、
前記送信部は、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に繰返し送信するパケット送信部と、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、
前記受信部は、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄部と、
受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
を備えたことを特徴とする通信装置。(3)

10

20

【 0 2 6 6 】

(付 記 4)

R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた R D M A 通信をサポートする通信装置に於いて、
前記送信部は、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信部と、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
を備え、
前記受信部は、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管理情報に記録するバッファ転送部と、
前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動させるデータ移動部と、
前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
を備えたことを特徴とする通信装置。(4)

30

40

【 0 2 6 7 】

(付 記 5)

R D M A パケットを作成して送信する送信部と、R D M A パケットを受信する受信部とを備えた R D M A 通信をサポートする通信装置に於いて、
転送先に許可を問合せることなく投機的に転送し、不許可の場合には再送要求を受けて再送する第 1 転送モード処理部と、

50

転送先に問合せることなく投機的に転送し、不許可の場合は不許可通知を受けて転送を中断し、再送要求を受けて再送する第2転送モード処理部と、
 転送完了通知を受けるまで転送先に許可を問合せることなく投機的に繰返し転送する第3転送モード処理部と、
 転送先に許可を問合せることなく投機的に転送し、不許可の場合は一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する第4転送モード処理部と、
 転送先に許可を問合せ、許可通知を受けて転送する第5転送モード処理部と、
 前記第1乃至第5転送モード処理部のいずれか1つを、ネットワーク負荷と受信側のメモリ使用量の少なくともいずれか一方に基づいて選択してR D M Aデータ転送を実行させる転送モード選択制御部と、
 を備えたことを特徴とする通信装置。(5)

【0268】

(付記6)

付記5記載の通信装置に於いて、
 前記第1転送モード処理部の送信部は、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからR D M Aパケットを作成して投機的に送信するパケット送信部と、
 転送先から再送要求を受信した際に、要求された転送データからR D M Aパケットを作成して送信するパケット再送部と、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
 を備え、

前記第1転送モード処理部の受信部は、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、
 前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
 を備えたことを特徴とする通信装置。(6)

【0269】

(付記7)

付記5記載の通信装置に於いて、
 前記第2転送モード処理部の送信部は、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからR D M Aパケットを作成して投機的に送信するパケット送信部と、
 転送先から受信不許可通知を受信した際に、パケット送信部によるパケット転送を中断する転送中断部と、
 転送先から再送要求を受信した際に、要求された転送データからR D M Aパケットを作成して送信するパケット再送部と、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
 を備え、
 前記第2転送モード処理部の受信部は、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管

10

20

30

40

50

理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄部と、

前記パケット破棄部で前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知部と、

前記パケット破棄部で受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。(7)

【0270】

(付記8)

付記5記載の通信装置に於いて、

前記第3転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に繰返し送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、

を備え、

前記第3転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄部と、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、

を備えたことを特徴とする通信装置。(8)

【0271】

(付記9)

付記5記載の通信装置に於いて、

前記第4転送モード処理部の送信部は、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に送信するパケット送信部と、

転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、

を備え、

前記第4転送モード処理部の受信部は、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信部と、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管

10

20

30

40

50

理情報に記録するバッファ転送部と、
 前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動させるデータ移動部と、
 前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知部と、
 を備えたことを特徴とする通信装置。(9)

【0272】

(付記10)

付記5記載の通信装置に於いて、
 前記第5転送モード処理部の送信部は、
 データ転送要求を受けた際に、転送先に受信許可確認を送信して受信許可通知を受信した場合に、転送データからRDMAパケットを作成して送信するパケット送信部と、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了部と、
 を備え、
 前記第5転送モード処理部の受信部は、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理部と、
 転送元から受信確認通知を受信した時に前記転送領域管理情報を参照して受信許可通知又は受信不許可通知を送信する確認応答部と、
 前記受信領域に受信パケットを転送するパケット受信部と、
 前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する転送完了通知部と、
 を備えたことを特徴とする通信装置。(10)

【0273】

(付記11)

付記5記載の通信装置に於いて、前記転送モード設定部は送信部に設けられ、選択した転送モードを転送元の送信部及び受信部に通知して前記第1乃至第5転送モード処理部のいずれかによるRDMAデータ転送を実行させることを特徴とする通信装置。

【0274】

(付記12)

付記5記載の通信装置に於いて、前記転送モード設定部は受信部に設けられ、選択した転送モードを送信部及び受信部に通知して前記第1乃至第5転送モード処理部のいずれかによるRDMAデータ転送を実行させることを特徴とする通信装置。

【0275】

(付記13)

付記5記載の通信装置に於いて、前記転送モード設定部は、
 前記受信側のメモリ使用率が少ない場合は前記第4転送モード処理部を選択し、
 前記受信側のメモリ使用率が多い場合は、ネットワーク負荷の低い順に、前記第3転送モード処理部、第1転送モード処理部及び第5転送モード処理部を順次選択することを特徴とする通信装置。(11)

【0276】

(付記14)

付記1乃至5記載のいずれかに記載の通信装置に於いて、更に、論理アドレスと物理アドレスのアドレス変換情報を持ち、前記アドレス変換情報内に前記転送領域管理情報を含ませて一体化したことを特徴とする通信装置。

【0277】

(付記15)

RDMAパケットを作成して送信する送信ステップと、RDMAパケットを受信する受信ステップとを備えたRDMA通信をサポートする通信方法に於いて、

10

20

30

40

50

前記送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
 転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、
 前記受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
 前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

【 0 2 7 8 】

(付 記 1 6)

R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた R D M A 通信をサポートする通信方法に於いて、
 前記送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
 転送先から受信不許可通知を受信した際に、前記パケット送信ステップによるパケット転送を中断する転送中断ステップと、
 転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、
 前記受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
 前記パケット破棄ステップで前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知ステップと、
 前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ス

10

20

30

40

50

テップと、
を備えたことを特徴とする通信方法。

【 0 2 7 9 】

(付 記 1 7)

R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた R D M A 通信をサポートする通信方法に於いて、

前記送信ステップは、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に繰返し送信するパケット送信ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を備え、

10

前記受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄ステップと、

受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

20

を備えたことを特徴とする通信方法。

【 0 2 8 0 】

(付 記 1 8)

R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた R D M A 通信をサポートする通信方法に於いて、

前記送信ステップは、

データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を備え、

30

前記受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、

パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、

パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管理情報に記録するバッファ転送ステップと、

40

前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動するデータ移動ステップと、

前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、

を備えたことを特徴とする通信方法。

【 0 2 8 1 】

(付 記 1 9)

R D M A パケットを作成して送信する送信ステップと、R D M A パケットを受信する受信ステップとを備えた R D M A 通信をサポートする通信方法に於いて、

50

転送先に許可を問合せることなく投機的に転送し、不許可の場合には再送要求を受けて再送する第1転送モード処理ステップと、
 転送先に問合せることなく投機的に転送し、不許可の場合は不許可通知を受けて転送を中断し、再送要求を受けて再送する第2転送モード処理ステップと、
 転送完了通知を受けるまで転送先に許可を問合せることなく投機的に繰返し転送する第3転送モード処理ステップと、
 転送先に許可を問合せることなく投機的に転送し、不許可の場合は一時バッファに保存し、許可が得られたら一時バッファから受信領域に転送する第4転送モード処理ステップと、
 転送先に許可を問合せ、許可通知を受けて転送する第5転送モード処理ステップと、
 前記第1乃至第5転送モード処理ステップのいずれか1つを、ネットワーク負荷と受信側のメモリ使用量の少なくともいずれか一方に基づいて選択してRDM Aデータ転送を実行させる転送モード設定ステップと、
 を備えたことを特徴とする通信方法。(12)

【0282】

(付記20)

付記19記載の通信方法に於いて、
 前記第1転送モード処理ステップの送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDM Aパケットを作成して投機的に送信するパケット送信ステップと、
 転送先から再送要求を受信した際に、要求された転送データからRDM Aパケットを作成して送信するパケット再送ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、
 前記第1転送モード処理ステップの受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
 前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

【0283】

(付記21)

付記19記載の通信方法に於いて、
 前記第2転送モード処理ステップの送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDM Aパケットを作成して投機的に送信するパケット送信ステップと、
 転送先から受信不許可通知を受信した際に、前記パケット送信ステップによるパケット転送を中断する転送中断ステップと、
 転送先から再送要求を受信した際に、要求された転送データからRDM Aパケットを作成して送信するパケット再送ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、

を備え、
 前記第 2 転送モード処理ステップの受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
 前記パケット破棄ステップで前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知ステップと、
 前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

【 0 2 8 4 】

(付記 22)

付記 19 記載の通信方法に於いて、
 前記第 3 転送モード処理ステップの送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に繰返し送信するパケット送信ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、
 前記第 3 転送モード処理ステップの受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄ステップと、
 受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

【 0 2 8 5 】

(付記 23)

付記 19 記載の通信方法に於いて、
 前記第 4 転送モード処理ステップの送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、
 前記第 4 転送モード処理ステップの受信ステップは、
 受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許

10

20

30

40

50

可に変更するパケット受信ステップと、
 パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管理情報に記録するバッファ転送ステップと、
 前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動するデータ移動ステップと、
 前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

10

【0286】

(付記24)

付記19記載の通信方法に於いて、
 前記第5転送モード処理ステップの第5送信ステップは、
 データ転送要求を受けた際に、転送先に受信許可確認を送信して受信許可通知を受信した場合に、転送データからRDMAパケットを作成して送信するパケット送信ステップと、
 転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
 を備え、

前記第5転送モード処理ステップの第5受信ステップは、

受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
 転送元から受信確認通知を受信した時に前記転送領域管理情報を参照して受信許可通知又は受信不許可通知を送信する確認応答ステップと、
 前記受信領域に受信パケットを転送するパケット受信ステップと、
 前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
 を備えたことを特徴とする通信方法。

20

【0287】

(付記25)

付記19記載の通信方法に於いて、前記転送モード選択ステップは転送元の送信ステップに設けられ、選択した転送モードを転送元及び転送先に通知して前記第1乃至第5転送モード処理ステップのいずれかによるRDMAデータ転送を実行させることを特徴とする通信方法。

30

【0288】

(付記26)

付記19記載の通信方法に於いて、前記転送モード設定ステップは転送先の受信ステップに設けられ、選択した転送モードを転送元及び転送先に通知して前記第1乃至第5転送モード処理ステップのいずれかによるRDMAデータ転送を実行させることを特徴とする通信方法。

【0289】

(付記27)

付記19記載の通信方法に於いて、前記転送モード設定ステップは、
 前記受信側のメモリ使用率が少ない場合は前記第4転送モード処理ステップを選択し、
 前記受信側のメモリ使用率が多い場合は、ネットワーク負荷の低い順に、前記第3転送モード処理ステップ、第1転送モード処理ステップ及び第5転送モード処理ステップを順次選択することを特徴とする通信方法。

40

【0290】

(付記28)

付記15乃至19記載のいずれかに記載の通信方法に於いて、更に、論理アドレスと物理アドレスのアドレス変換情報を持ち、前記アドレス変換情報内に前記転送領域管理情報

50

を含ませて一体化したことを特徴とする通信方法。

【 0 2 9 1 】

(付記 2 9)

R D M A 通信をサポートする通信装置のコンピュータに、
送信ステップとして、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を実行させ、
前記受信ステップとして、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
を実行させることを特徴とする通信プログラム。

10

20

【 0 2 9 2 】

(付記 3 0)

R D M A 通信をサポートする通信装置のコンピュータに、
送信ステップとして、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データから R D M A パケットを作成して投機的に送信するパケット送信ステップと、
転送先から受信不許可通知を受信した際に、前記パケット送信ステップによるパケット転送を中断する転送中断ステップと、
転送先から再送要求を受信した際に、要求された転送データから R D M A パケットを作成して送信するパケット再送ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を実行させ、
受信ステップとして、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄すると共に前記パケット破棄を前記転送領域管理情報に記録するパケット破棄ステップと、
前記パケット破棄ステップで前記受信領域の転送不許可を判別した場合に、転送元に受信不許可通知を送信する不許可通知ステップと、
前記パケット破棄ステップで受信パケットを破棄した後に前記受信領域の転送許可を判別

30

40

50

した場合に、前記転送領域管理情報のパケット破棄の記録に基づいて転送元に再送要求を送信する再送要求ステップと、
受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
を実行させることを特徴とする通信プログラム。

【0293】

(付記31)

RDMA通信をサポートする通信装置のコンピュータに、
送信ステップとして、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に繰返し送信するパケット送信ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を実行させ、
受信ステップとして、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に、受信パケットを破棄するパケット破棄ステップと、
受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
を実行させることを特徴とする通信プログラム。

【0294】

(付記32)

RDMA通信をサポートする通信装置のコンピュータに、
送信ステップとして、
データ転送要求を受けた際に、転送先に受信許可の有無を問合せることなく、転送データからRDMAパケットを作成して投機的に送信するパケット送信ステップと、
転送先から転送完了通知を受信してパケット送信を正常終了する送信完了ステップと、
を実行させ、
受信ステップとして、
受信領域に対する転送許可又は転送不許可を含む転送領域管理情報を管理する転送領域管理ステップと、
パケット受信時に前記転送領域管理情報を参照して受信領域の転送許可を判別した場合に、前記受信領域に前記受信パケットを転送すると共に前記受信領域の転送許可を転送不許可に変更するパケット受信ステップと、
パケット受信時に前記転送領域管理情報を参照して前記受信領域の転送不許可を判別した場合に受信パケットをバッファ領域に転送すると共に前記バッファ転送を前記転送領域管理情報に記録するバッファ転送ステップと、
前記バッファ転送後に前記受信領域の転送許可を判別した場合に、前記転送領域管理情報のバッファ転送の記録に基づいて前記バッファ領域のデータを前記受信領域に移動するデータ移動ステップと、
前記受信領域に対する受信パケットの転送完了を認識して転送完了通知を送信する完了通知ステップと、
を実行させることを特徴とする通信プログラム。

【図面の簡単な説明】

【0295】

【図1】ホスト間通信に適用された本発明の一実施形態の説明図

【図 2】本実施形態におけるNICの送信部と受信部の機能構成のブロック図	
【図 3】本実施形態で使用する転送領域管理表の説明図	
【図 4】本実施形態における送信制御部のブロック図	
【図 5】本実施形態における受信制御部のブロック図	
【図 6】本実施形態における送信処理のフローチャート	
【図 7】本実施形態における受信処理のフローチャート	
【図 8】本実施形態における第1転送モードによる転送処理の説明図	
【図 9】転送開始前に受信可能となっている場合の第1転送モードによる転送処理のタイムチャート	
【図 10】転送開始後に受信可能となった場合の第1転送モードによる転送処理のタイムチャート	10
【図 11】第1転送モード送信処理のフローチャート	
【図 12】第1転送モード受信処理のフローチャート	
【図 13】本実施形態における第2転送モードによる転送処理の説明図	
【図 14】転送開始前に受信可能となっている場合の第2転送モードによる転送処理のタイムチャート	
【図 15】転送開始後に受信可能となった場合の第2転送モードによる転送処理のタイムチャート	
【図 16】図 15 に続く第2転送モードによる転送処理のタイムチャート	
【図 17】第2転送モード送信処理のフローチャート	20
【図 18】第2転送モード受信処理のフローチャート	
【図 19】本実施形態における第3転送モードによる転送処理の説明図	
【図 20】第3転送モードによる転送処理のタイムチャート	
【図 21】第3転送モード送信処理のフローチャート	
【図 22】第3転送モード受信処理のフローチャート	
【図 23】本実施形態における第4転送モードによる転送処理のタイムチャート	
【図 24】第4転送モード送信処理のフローチャート	
【図 25】第4転送モード受信処理のフローチャート	
【図 26】第4転送モードの受信側における一時バッファから受信領域へのデータ移動の説明図	30
【図 27】図 26 のデータ移動におけるアドレスマッピングの説明図	
【図 28】図 27 のアドレスマッピング前と後のアドレス変換表の説明図	
【図 29】本実施形態における第5転送モードによる転送処理の説明図	
【図 30】第5転送モードによる転送処理のタイムチャート	
【図 31】第5転送モード送信処理のフローチャート	
【図 32】第5転送モード受信処理のフローチャート	
【図 33】転送開始前に受信可能となっている場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図	
【図 34】図 33 の転送完了時間の計算の説明図	
【図 35】転送開始の待ち時間 W が $W < kL$ の場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図	40
【図 36】図 35 の転送完了時間の計算の説明図	
【図 37】転送開始の待ち時間 W が $kL < W < kL - 2$ の場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図	
【図 38】図 37 の転送完了時間の計算の説明図	
【図 39】転送開始の待ち時間 W が $kL - 2 < W < 2$ の場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図	
【図 40】図 39 の転送完了時間の計算の説明図	
【図 41】転送開始の待ち時間 W が $2 < W$ の場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図	50

【図 4 2】図 4 1 の転送完了時間の計算の説明図

【図 4 3】図 4 2 (C) における転送完了時間の計算詳細を示した説明図

【図 4 4】本実施形態で使用する転送モード管理表の説明図

【図 4 5】本実施形態における転送モード選択処理のフローチャート

【符号の説明】

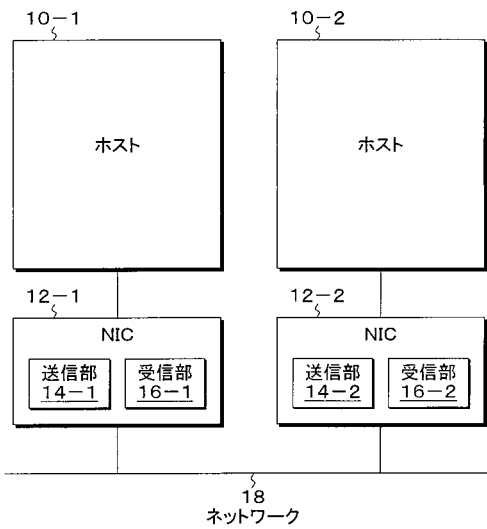
【 0 2 9 6 】

1 0 - 1 , 1 0 - 2 : ホスト	
1 2 - 1 , 1 2 - 2 : ネットワーク・インタフェース・カード (N I C)	
1 4 - 1 , 1 4 - 2 : 送信部	
1 6 - 1 , 1 6 - 2 : 受信部	10
1 8 : ネットワーク	
2 0 : 送信要求部	
2 2 : 送信領域	
2 4 : 送信データ	
2 6 : 送信要求受付部	
2 8 , 4 2 : 転送領域管理部	
3 0 , 4 4 : 送信制御部	
3 2 , 4 0 : 受信制御部	
3 4 , 4 6 : モード選択制御部	
3 6 , 4 8 : ネットワーク負荷情報収集部	20
3 8 , 5 0 : 1 0 G b イーサネット M A C	
4 9 : メモリ使用率情報収集部	
5 2 : 領域許可部	
5 4 : 受信領域	
5 5 : 転送データ	
5 6 : 一時バッファ	
5 8 : 保存データ	
6 0 : 転送制御部	
6 2 : 転送領域管理表	
6 4 : 領域 I D	30
6 6 : 論理アドレス	
6 8 : 物理アドレス	
7 0 : 受信許可	
7 2 : パケット破棄情報	
7 4 : 受信不許可通知情報	
7 5 : 一時バッファアドレス	
7 6 : 第 1 転送モード送信部	
7 8 : 第 2 転送モード送信部	
8 0 : 第 3 転送モード送信部	
8 2 : 第 4 転送モード送信部	40
8 4 : 第 5 転送モード送信部	
8 6 , 9 2 , 1 0 0 , 1 0 4 , 1 1 0 : パケット送信部	
8 8 , 9 6 : パケット再送部	
9 0 , 9 8 , 1 0 2 , 1 0 6 , 1 1 2 : 送信完了部	
9 4 : 転送中断部	
1 0 8 : 受信許可確認部	
1 1 4 : 第 1 転送モード受信部	
1 1 6 : 第 2 転送モード受信部	
1 1 8 : 第 3 転送モード受信部	
1 2 0 : 第 4 転送モード受信部	50

- 1 2 4 : 第 5 転送モード受信部
- 1 2 6 , 1 3 4 , 1 4 4 , 1 5 0 , 1 6 0 : パケット受信部
- 1 2 8 , 1 3 6 , 1 4 6 : パケット破棄部
- 1 3 0 , 1 4 0 : 再送要求部
- 1 3 2 , 1 4 2 , 1 4 8 , 1 5 6 , 1 6 2 : 完了通知部
- 1 3 8 : 不許可通知部
- 1 5 2 : バッファ転送部
- 1 5 4 : データ移動部
- 1 5 8 : 確認応答部

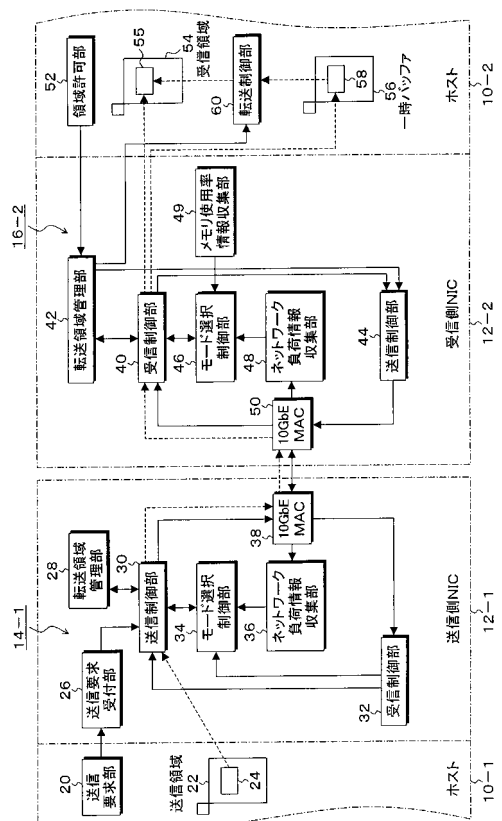
【 図 1 】

ホスト間通信に適用された本発明の一実施形態の説明図



【 図 2 】

本実施形態におけるNICの送信部と受信部の機能構成のブロック図



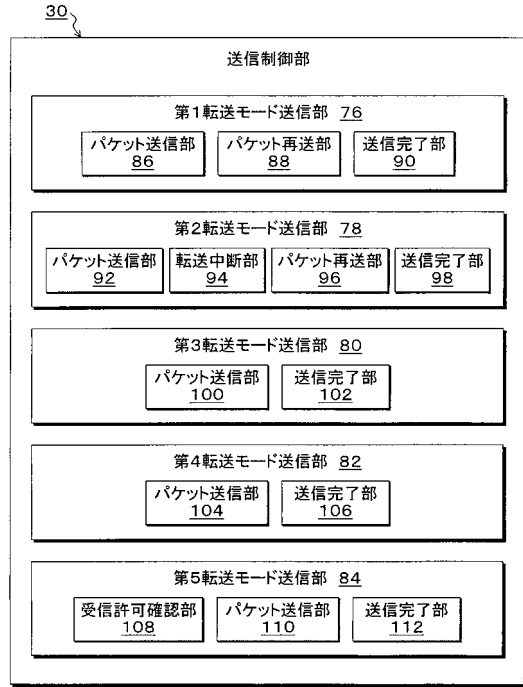
【図3】

本実施形態で使用する転送領域管理表の説明図

64		66		68		70		72		74		75	
領域ID	論理アドレス	物理アドレス	受信許可	パケット破棄情報	受信不許可通知情報	一時ハワアアドレス							
00	0x0000	0x3000	1										
01	0x1000	0x4000	1										
02	0x2000	0x5000	0										
03	0x3000	0x6000	1										
04	0x4000	0x7000	0										
05	0x5000	0x8000	0										
06	0x6000	0x8000	1										
07	0x7000	0x1000	1										
08	0x8000	0x2000	0										

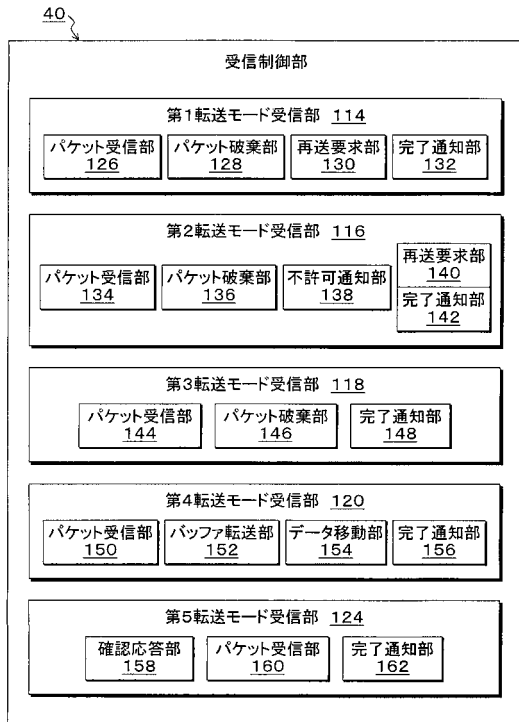
【図4】

本実施形態における送信制御部のブロック図



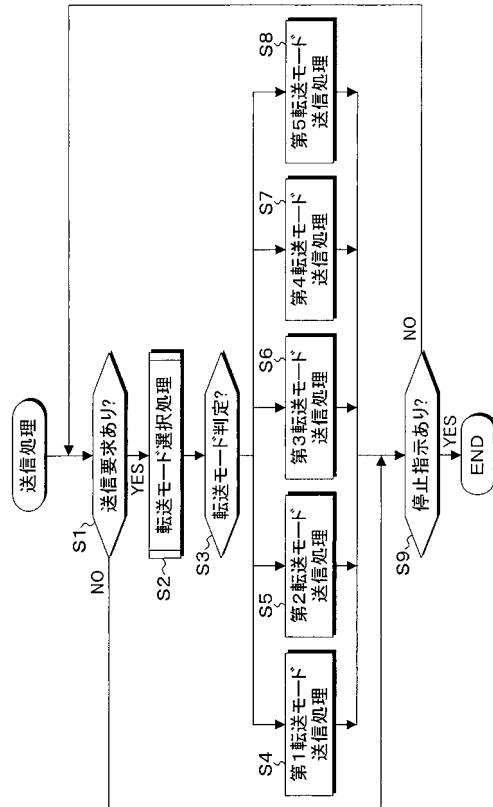
【図5】

本実施形態における受信制御部のブロック図



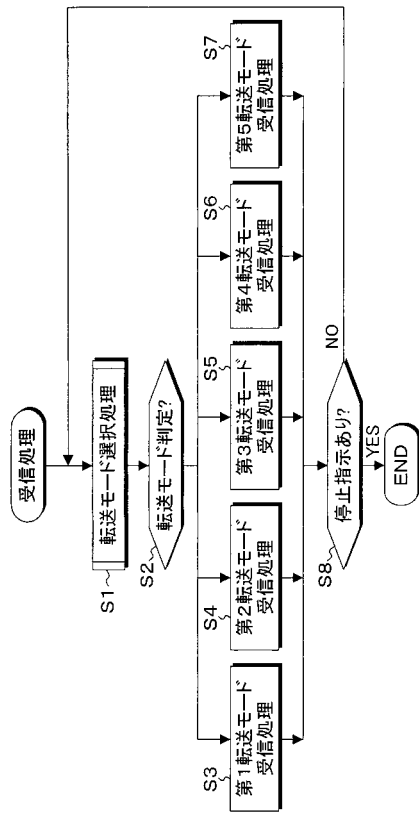
【図6】

本実施形態における送信処理のフローチャート



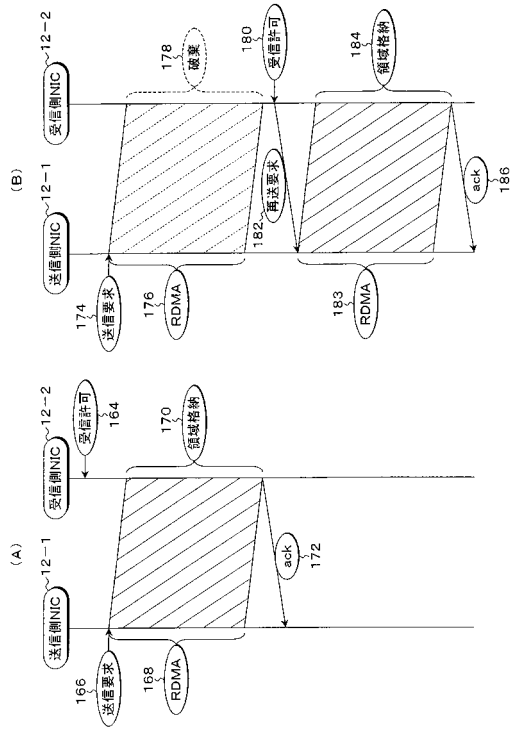
【 図 7 】

本実施形態における受信処理のフローチャート



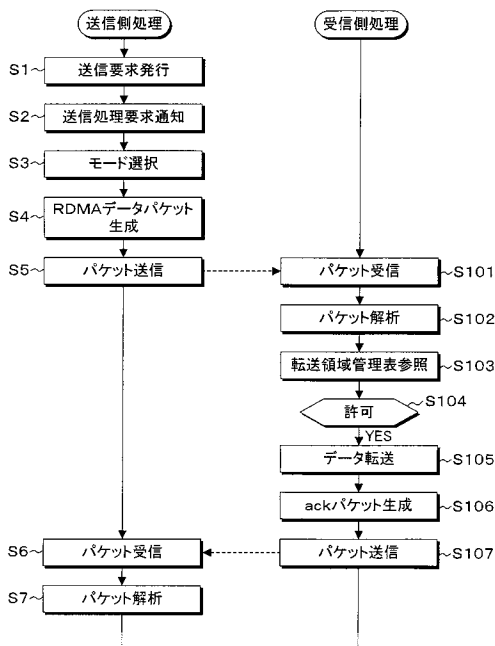
【 図 8 】

本実施形態における第1転送モードによる転送処理の説明図



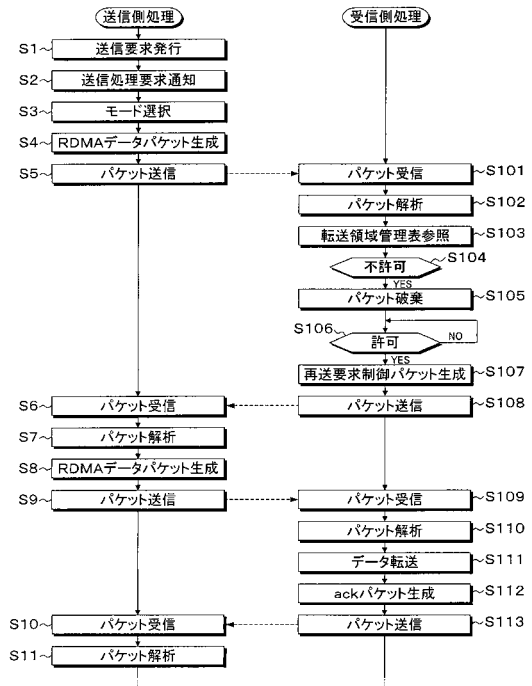
【 図 9 】

転送開始前に受信可能となっている場合の第1転送モードによる転送処理のタイムチャート

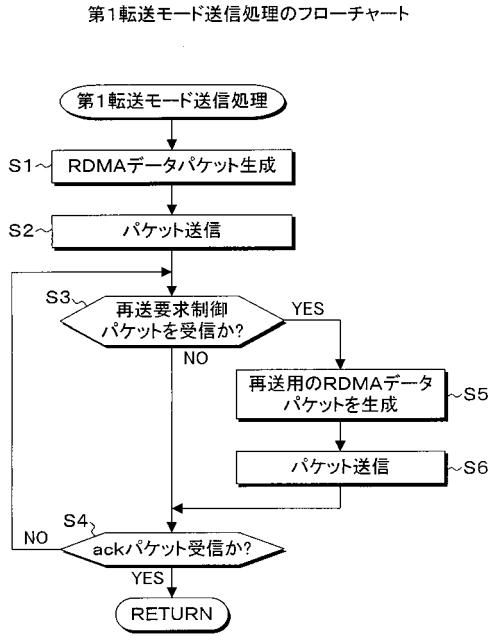


【 図 10 】

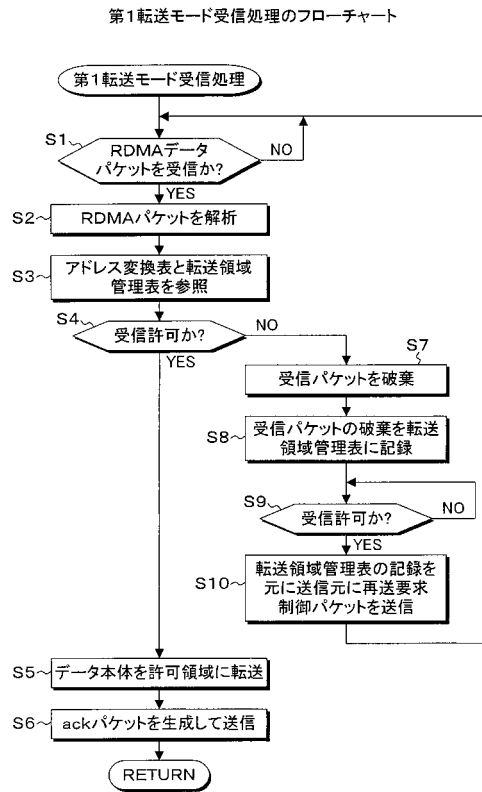
転送開始後に受信可能となった場合の第1転送モードによる転送処理のタイムチャート



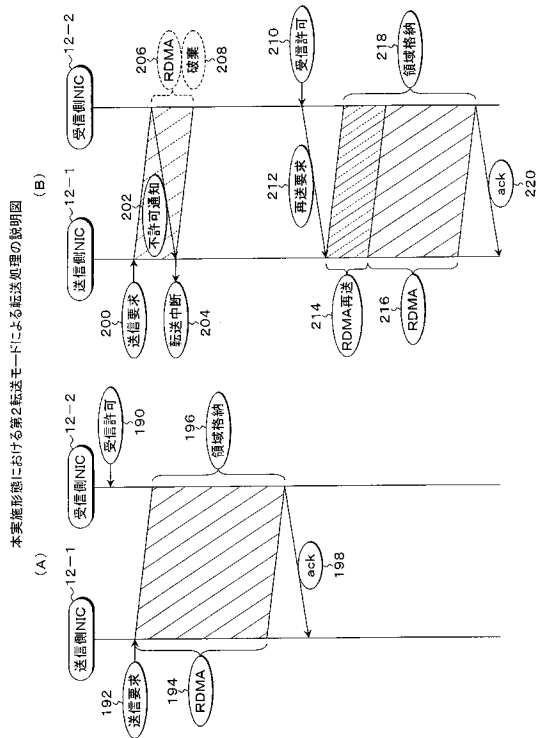
【図11】



【図12】

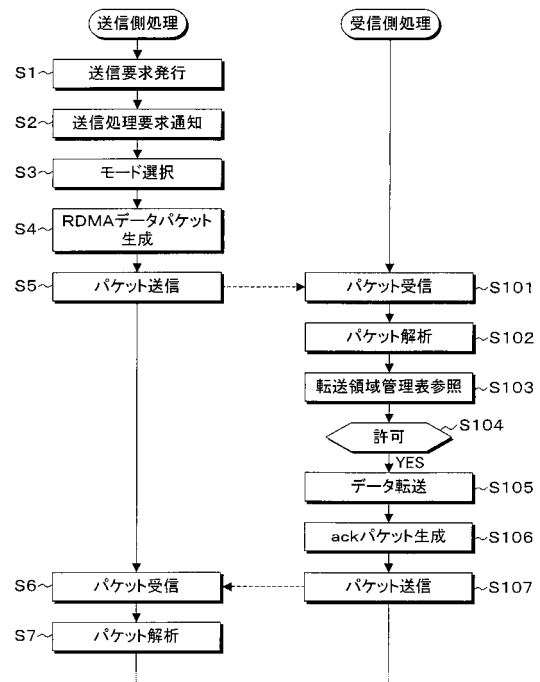


【図13】



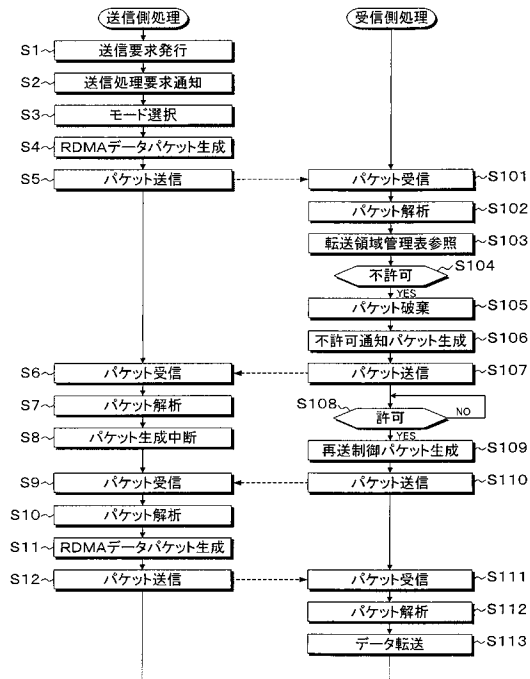
【図14】

転送開始前に受信可能となっている場合の第2転送モードによる転送処理のタイムチャート



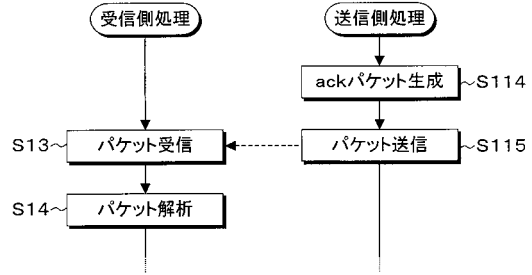
【図15】

転送開始後に受信可能となった場合の第2転送モードによる転送処理のタイムチャート



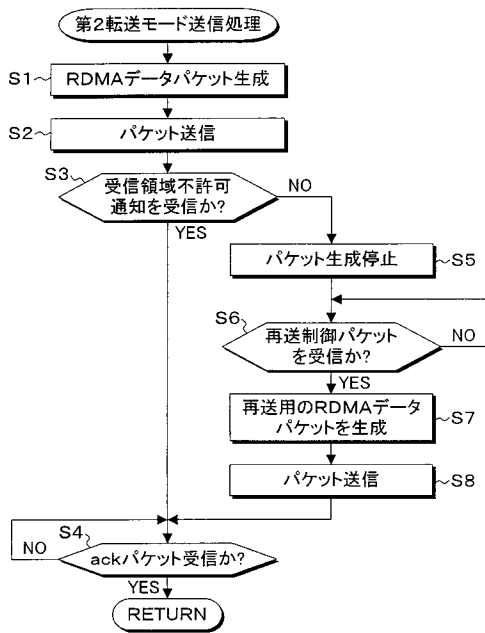
【図16】

図15に続く第2転送モードによる転送処理のタイムチャート



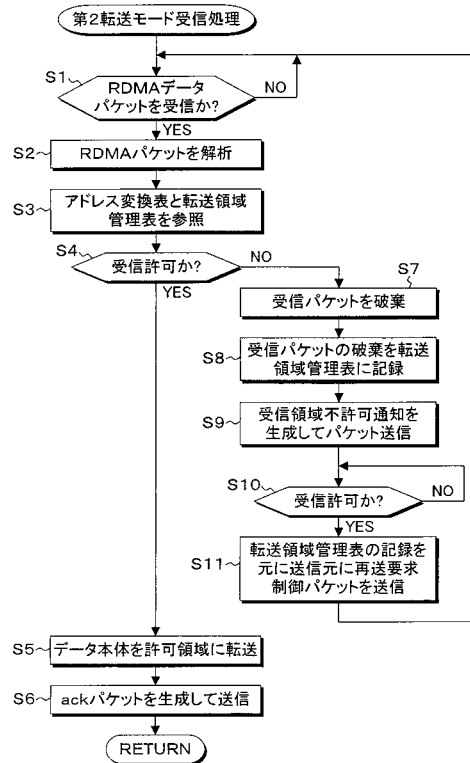
【図17】

第2転送モード送信処理のフローチャート

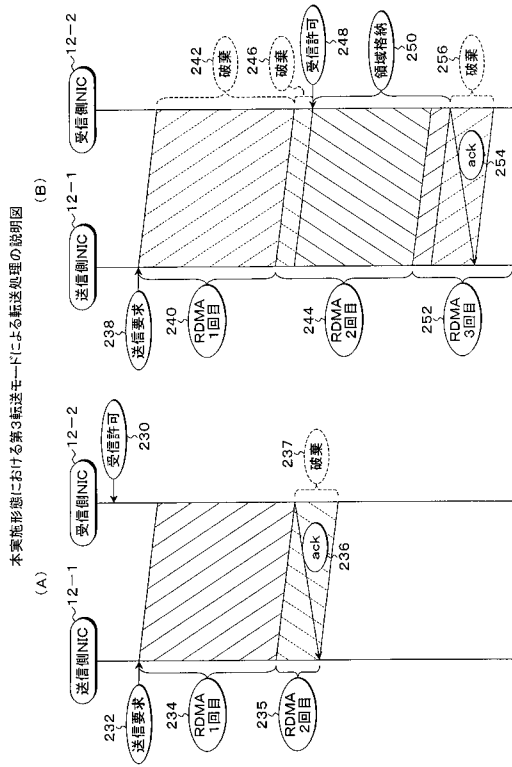


【図18】

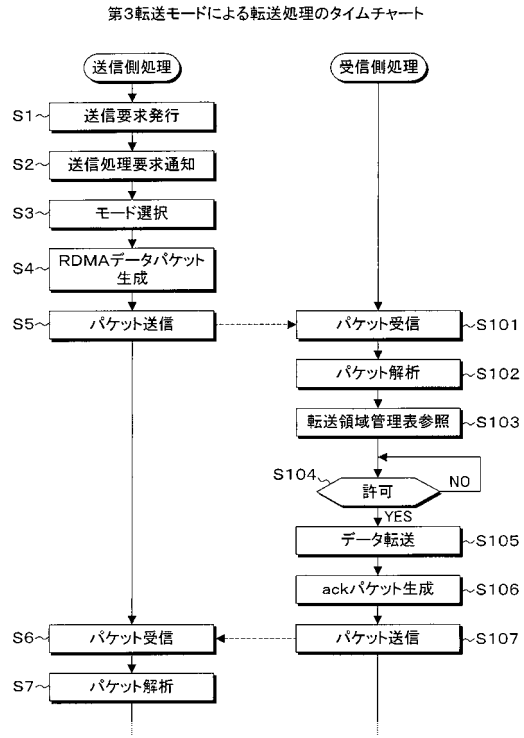
第2転送モード受信処理のフローチャート



【図19】

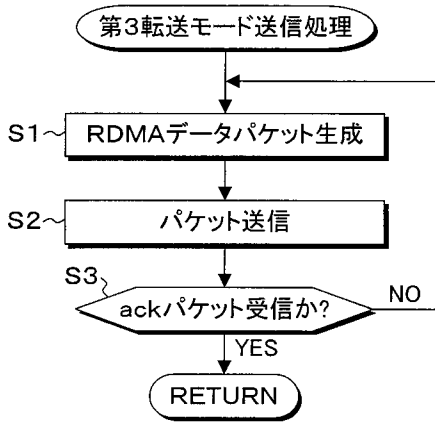


【図20】



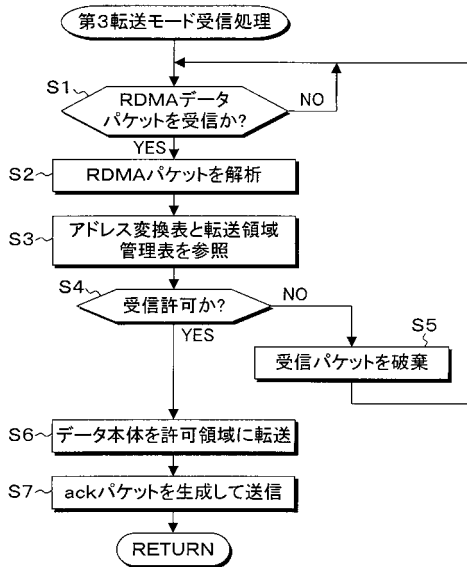
【図21】

第3転送モード送信処理のフローチャート

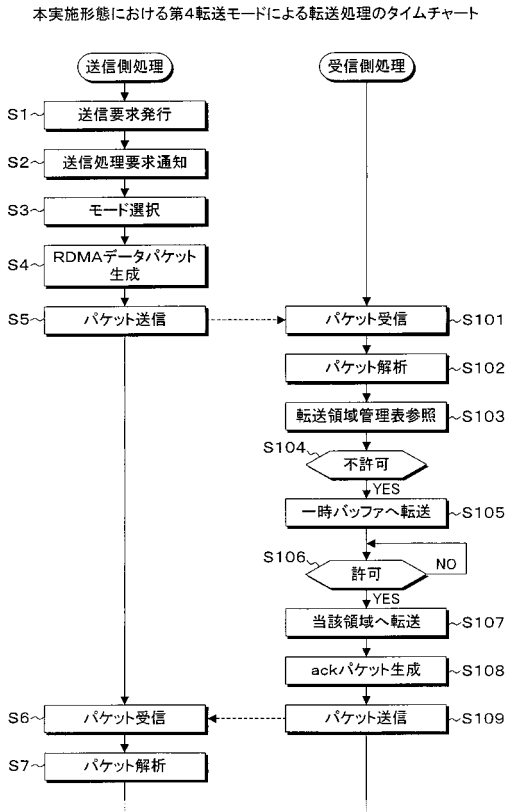


【図22】

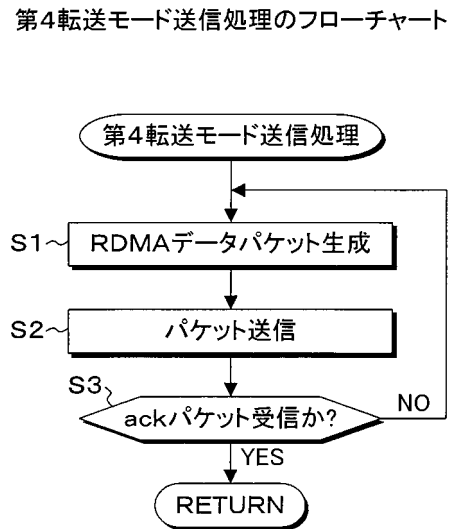
第3転送モード受信処理のフローチャート



【図 2 3】

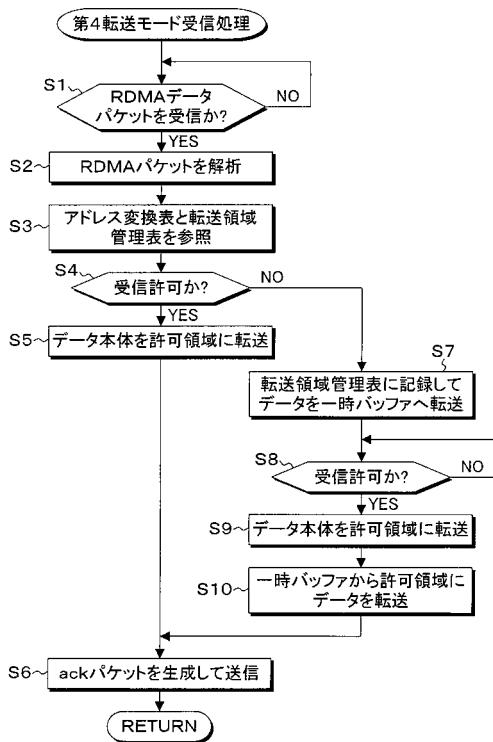


【図 2 4】



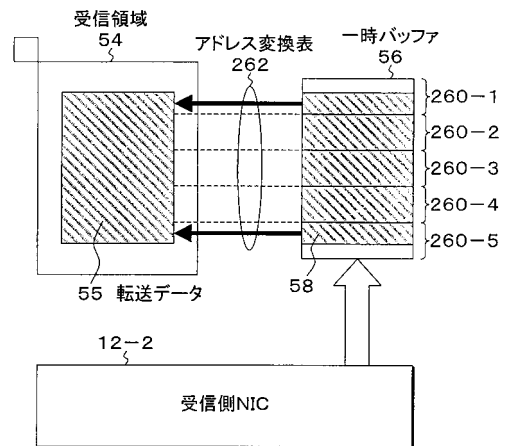
【図 2 5】

第4転送モード受信処理のフローチャート



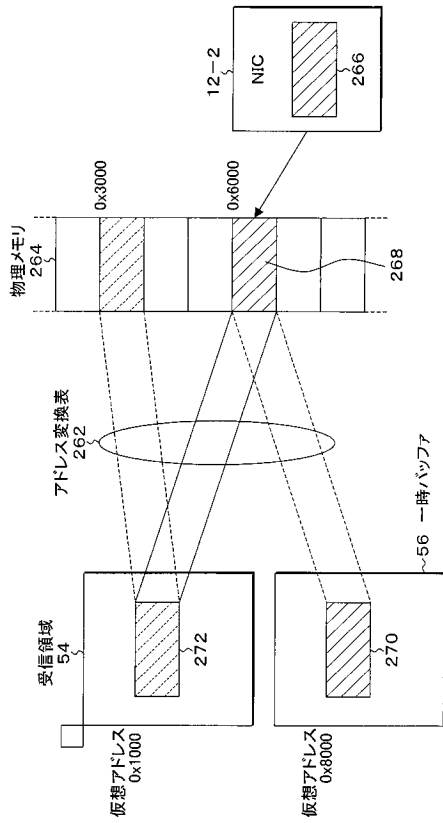
【図 2 6】

第4転送モードの受信側における一時バッファから受信領域へのデータ移動の説明図



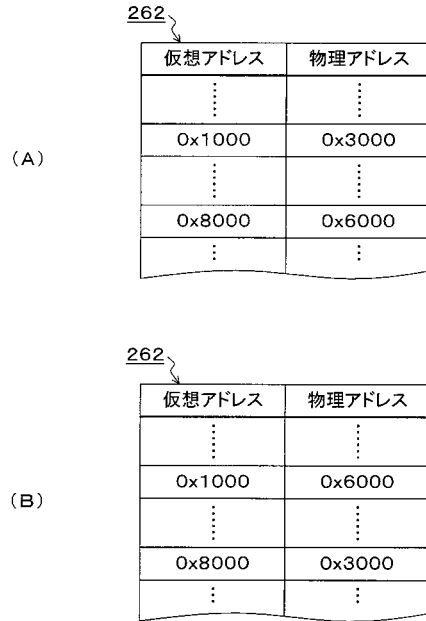
【図 27】

図26のデータ移動におけるアドレスマッピングの説明図



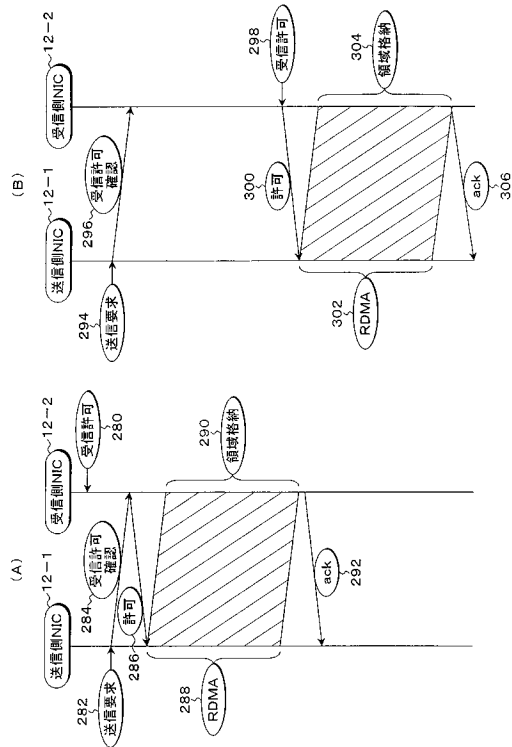
【図 28】

図27のアドレスマッピング前と後のアドレス変換表の説明図



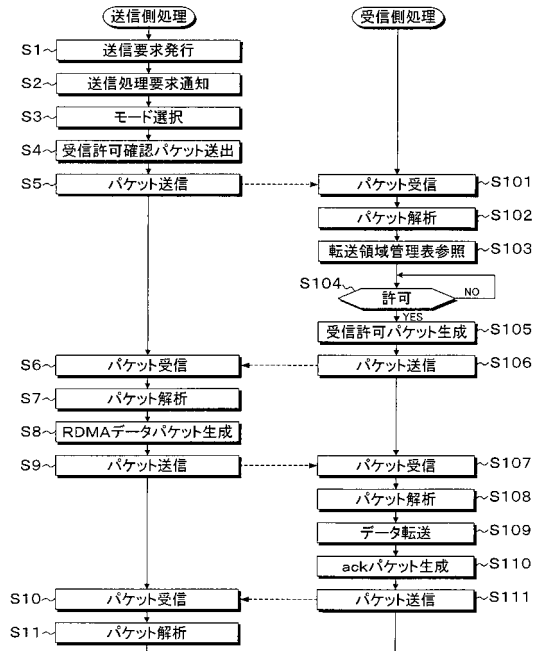
【図 29】

本実施形態における第5転送モードによる転送処理の説明図



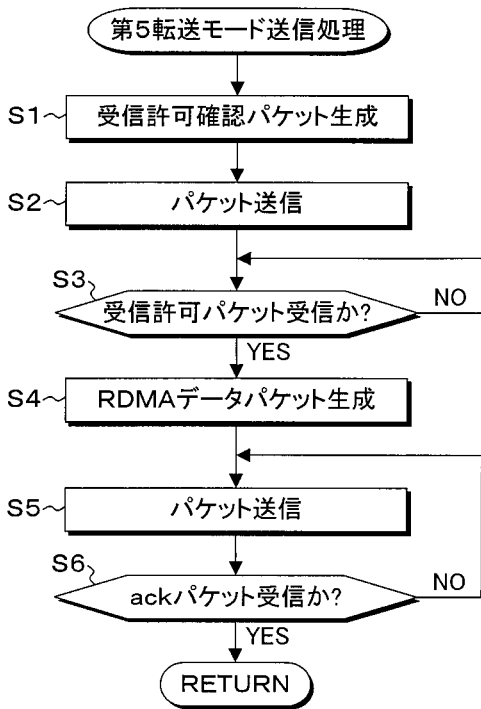
【図 30】

第5転送モードによる転送処理のタイムチャート



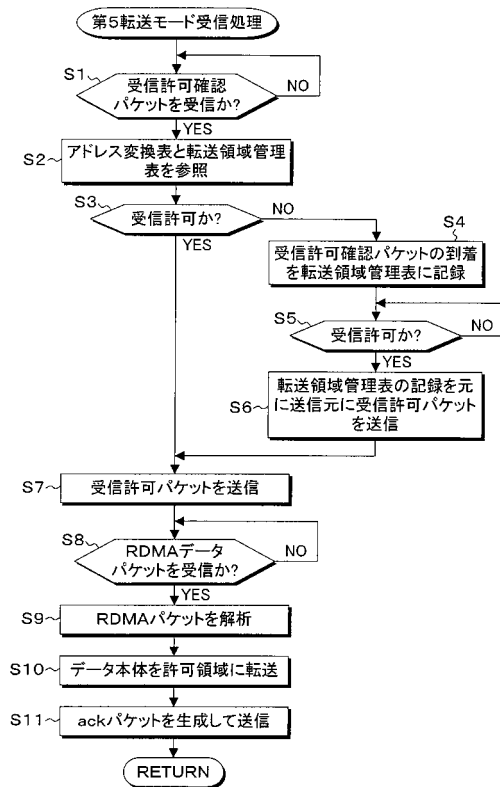
【図31】

第5転送モード送信処理のフローチャート



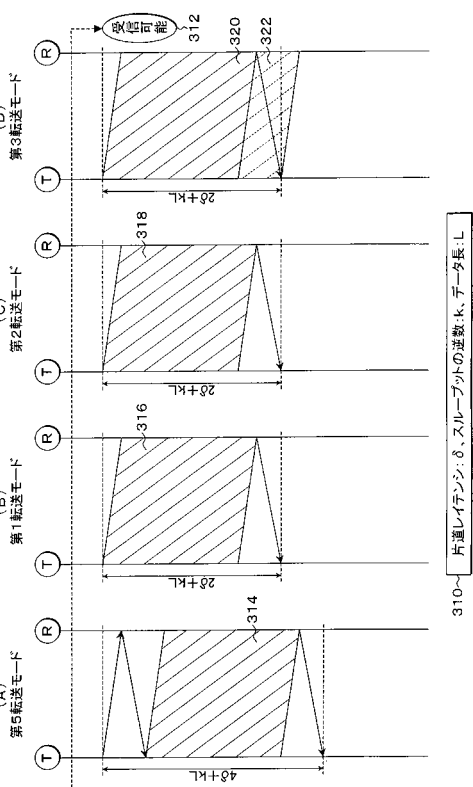
【図32】

第5転送モード受信処理のフローチャート



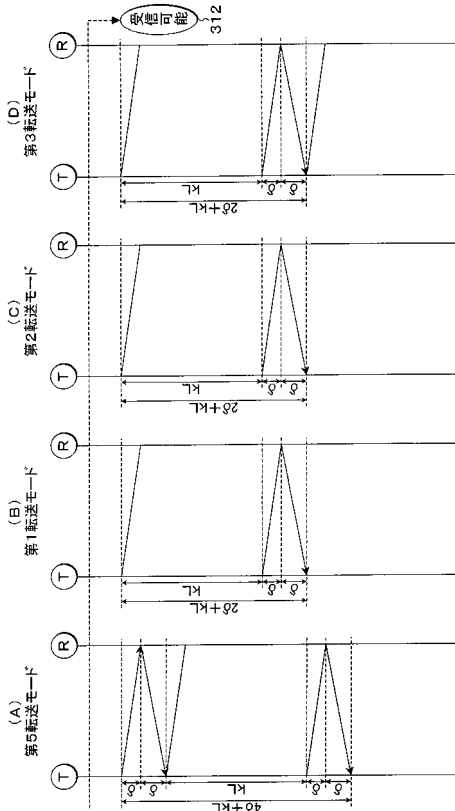
【図33】

転送開始前に受信可能となっている場合の第1乃至第3及び第5転送モードによる転送完了時間の説明図

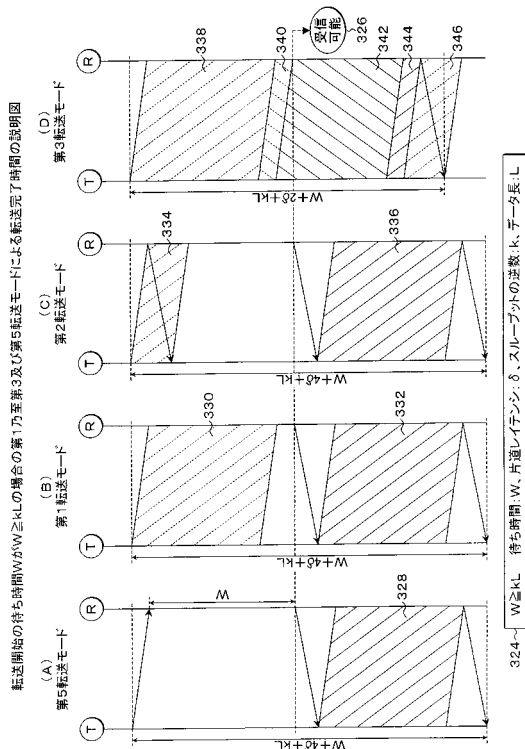


【図34】

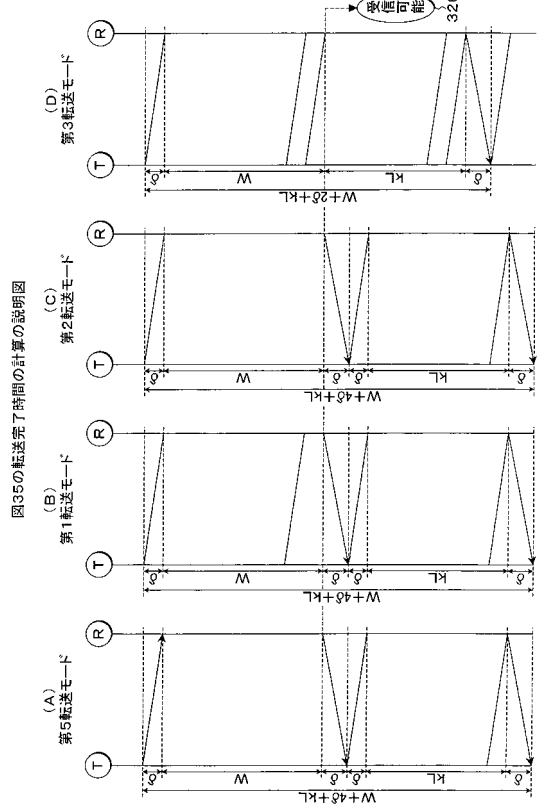
図33の転送完了時間の計算の説明図



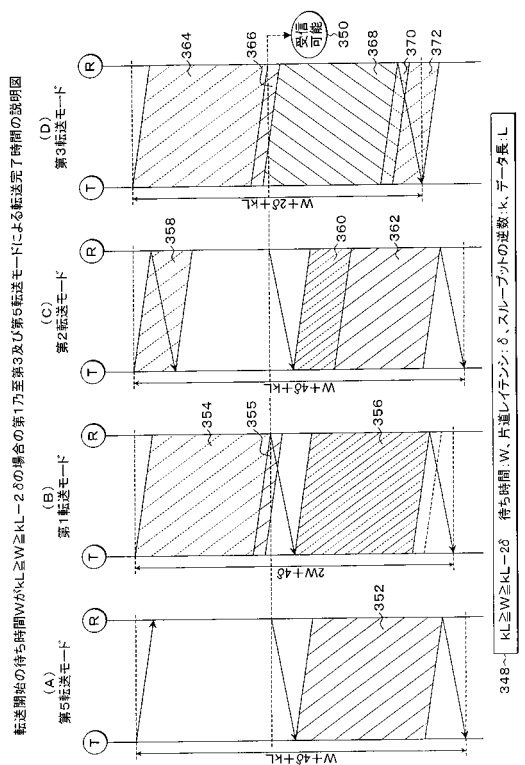
【 図 3 5 】



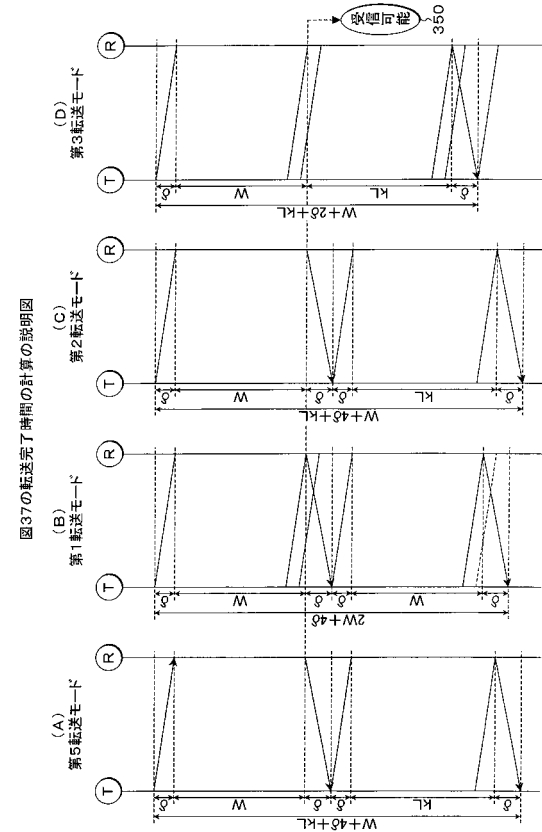
【 図 3 6 】



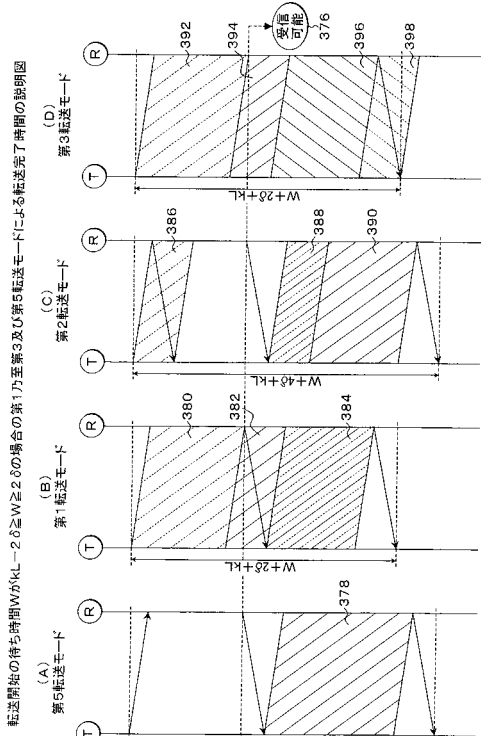
【 図 3 7 】



【 図 3 8 】



【 図 39 】



374 ~ $kL - 2\delta \geq W \geq 2\delta$ 待ち時間: W、片道レイテンシ: δ 、スループットの逆数: k、データ長: L

【 図 40 】

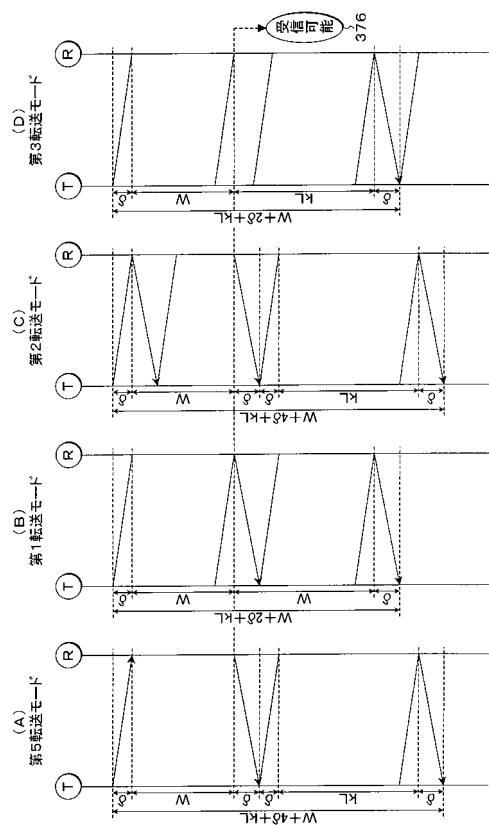
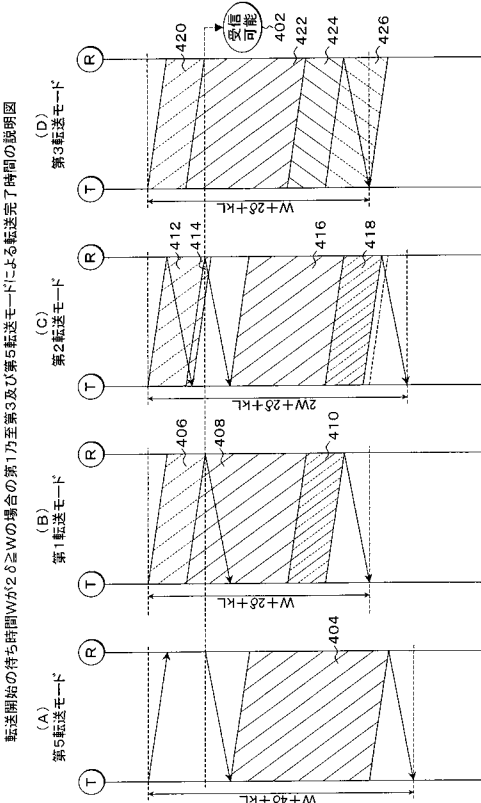


図39の転送完了時間の計算の説明図

【 図 41 】



400 ~ $2\delta \geq W$ 待ち時間: W、片道レイテンシ: δ 、スループットの逆数: k、データ長: L

【 図 42 】

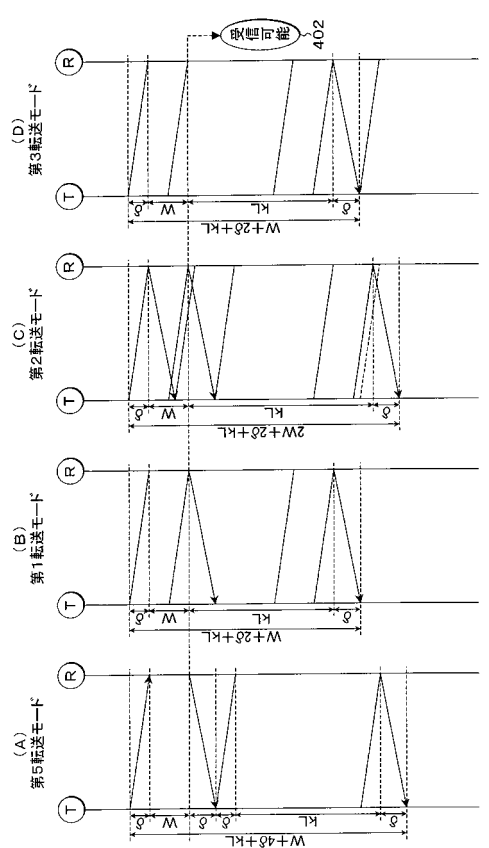
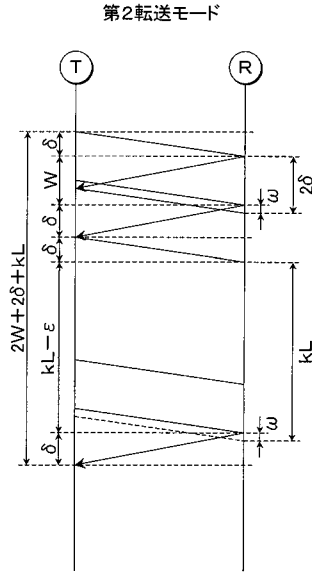


図41の転送完了時間の計算の説明図

【 図 4 3 】

図42(C)における転送完了時間の計算詳細を示した説明図



【 図 4 4 】

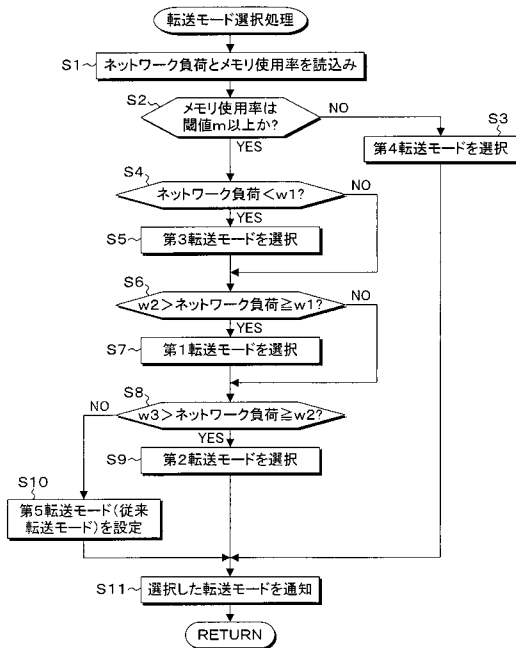
本実施形態で使用する転送モード管理表の説明図

転送モード管理表
450

ネットワーク負荷	メモリ使用率	
	m未満	m以上
w1未満	第4転送モード	第3転送モード
w1以上w2未満		第1転送モード
w2以上w3未満		第2転送モード
w3以上		第5転送モード

【 図 4 5 】

本実施形態における転送モード選択処理のフローチャート



フロントページの続き

(56)参考文献 特開平07-219916(JP,A)

特開平11-004256(JP,A)

特開平11-007434(JP,A)

住元真司 他, 10Gb Ethernetを用いた高性能通信機構の設計, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 2004年 8月 1日, Vol.2004 No.81, p.121~126

(58)調査した分野(Int.Cl., DB名)

G06F 13/00