

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6464911号
(P6464911)

(45) 発行日 平成31年2月6日(2019.2.6)

(24) 登録日 平成31年1月18日(2019.1.18)

(51) Int.Cl. F I
 HO 4 L 12/807 (2013.01) HO 4 L 12/807
 HO 4 L 29/08 (2006.01) HO 4 L 13/00 3 O 7 C

請求項の数 10 (全 32 頁)

(21) 出願番号 特願2015-93799 (P2015-93799)
 (22) 出願日 平成27年5月1日(2015.5.1)
 (65) 公開番号 特開2016-213581 (P2016-213581A)
 (43) 公開日 平成28年12月15日(2016.12.15)
 審査請求日 平成30年2月6日(2018.2.6)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100089118
 弁理士 酒井 宏明
 (72) 発明者 李 忠翰
 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 森田 充功

最終頁に続く

(54) 【発明の名称】 情報処理システム、情報処理システムの制御方法及び受信装置

(57) 【特許請求の範囲】

【請求項1】

パケットを送信する送信装置と、複数のスイッチ装置と、前記送信装置が送信したパケットを前記複数のスイッチ装置を介して受信する複数の受信装置と、前記送信装置と前記複数の受信装置を制御する制御装置とを有する情報処理システムにおいて、

各受信装置は、

前記送信装置が有する輻輳ウィンドウのサイズの増加を監視する監視部と、

前記送信装置に対する応答パケットを前記送信装置に送信してから、前記応答パケットに対する次のデータを前記送信装置から受信するまでの時間である往復遅延時間を測定する測定部と、

前記監視部が、前記輻輳ウィンドウのサイズの増加が定常状態になったことを監視した場合、パケットを受信するフローのうちボリュームが所定サイズ以上の対象フローの最小スループットを計算する計算部と、

過去に応答パケットに設定した受信ウィンドウのサイズである設定済ウィンドウサイズと、前記最小スループットから計算した受信ウィンドウのサイズである計算ウィンドウサイズとに基づいて、設定用の受信ウィンドウサイズである設定ウィンドウサイズを生成する生成部と、

前記設定ウィンドウサイズが設定された応答パケットを前記送信装置に送信する設定部とを有し、

前記制御装置は、

前記複数の受信装置から受信した平均往復遅延時間に基づいて、前記制御装置からの設定用の受信ウィンドウサイズである制御設定ウィンドウサイズを算出する算出部と、

前記制御設定ウィンドウサイズが設定された制御パケットを前記複数の受信装置の各々に送信する通信部と

を有することを特徴とする情報処理システム。

【請求項 2】

前記生成部は、対象フローのスループット及び前記設定済ウィンドウサイズから計算される値と前記設定済ウィンドウサイズとの比率に基づいて、前記設定済ウィンドウサイズを生成することを特徴とする請求項 1 記載の情報処理システム。

【請求項 3】

前記複数のスイッチ装置には、

前記送信装置に接続する第 1 のスイッチ装置と、

前記受信装置に接続する第 2 のスイッチ装置と、

前記第 1 のスイッチ装置と前記第 2 のスイッチ装置とに接続する第 3 のスイッチ装置と

、
前記第 1 のスイッチ装置と前記第 2 のスイッチ装置とに接続する第 4 のスイッチ装置とが含まれることを特徴とする請求項 1 又は 2 記載の情報処理システム。

【請求項 4】

前記算出部は、前記対象フローが通過するリンクの使用可能バンド幅のうちの最小値に基づいて前記制御設定ウィンドウサイズを算出することを特徴とする請求項 1、2 又は 3 記載の情報処理システム。

【請求項 5】

前記算出部は、前記使用可能バンド幅が前記最小値であるリンクにおいて前記対象フローにより使用されているバンド幅の総和を前記最小値に加えた値に前記平均往復遅延時間を乗じて前記対象フローの数で割った値を前記制御設定ウィンドウサイズとして算出することを特徴とする請求項 4 記載の情報処理システム。

【請求項 6】

送信装置及び受信装置が仮想スイッチを用いて通信する場合には、前記監視部、前記設定部、前記計算部、前記生成部及び前記設定部は、前記仮想スイッチに含まれることを特徴とする請求項 1 ~ 5 のいずれか 1 つ記載の情報処理システム。

【請求項 7】

前記生成部は、前記最小スループットとパケットのサイズのうち大きな方を前記計算ウィンドウサイズとすることを特徴とする請求項 1 ~ 6 のいずれか 1 つ記載の情報処理システム。

【請求項 8】

前記生成部は、前記対象フローが複数あり、前記比率が 0 である場合に、前記比率を強制的に 0 と 1 の間の数にすることを特徴とする請求項 2 ~ 7 のいずれか 1 つ記載の情報処理システム。

【請求項 9】

パケットを送信する送信装置と、複数のスイッチ装置と、前記送信装置が送信したパケットを前記複数のスイッチ装置を介して受信する複数の受信装置と、前記送信装置と前記複数の受信装置を制御する制御装置とを有する情報処理システムの制御方法において、

各受信装置が、

前記送信装置が有する輻輳ウィンドウのサイズの増加を監視し、

前記送信装置に対する応答パケットを前記送信装置に送信してから、前記応答パケットに対する次のデータを前記送信装置から受信するまでの時間である往復遅延時間を測定し、

前記輻輳ウィンドウのサイズの増加が定常状態になった場合、パケットを受信するフローのうちボリュームが所定サイズ以上の対象フローの最小スループットを計算し、

過去に応答パケットに設定した受信ウィンドウのサイズである設定済ウィンドウサイズ

10

20

30

40

50

と、前記最小スループットから計算した受信ウィンドウのサイズである計算ウィンドウサイズとに基づいて、設定用の受信ウィンドウサイズである設定ウィンドウサイズを生成し

、
前記設定ウィンドウサイズが設定された応答パケットを前記送信装置に送信し、
前記制御装置が、

前記複数の受信装置から受信した平均往復遅延時間に基づいて、前記制御装置からの設定用の受信ウィンドウサイズである制御設定ウィンドウサイズを算出し、

前記制御設定ウィンドウサイズが設定された制御パケットを前記複数の受信装置の各々に送信する

ことを特徴とする情報処理システムの制御方法。

10

【請求項 10】

パケットを送信する送信装置と複数のスイッチ装置と前記送信装置を制御する制御装置と他の受信装置とともに情報処理システムを構築し、前記送信装置が送信したパケットを前記複数のスイッチ装置を介して受信し、かつ、前記制御装置により制御される受信装置において、

前記送信装置が有する輻輳ウィンドウのサイズの増加を監視する監視部と、

前記送信装置に対する応答パケットを前記送信装置に送信してから、前記応答パケットに対する次のデータを前記送信装置から受信するまでの時間である往復遅延時間を測定する測定部と、

前記監視部が、前記輻輳ウィンドウのサイズの増加が定常状態になったことを監視した場合、パケットを受信するフローのうちボリュームが所定サイズ以上の対象フローの最小スループットを計算する計算部と、

20

過去に応答パケットに設定した受信ウィンドウのサイズである設定済ウィンドウサイズと、前記最小スループットから計算した受信ウィンドウのサイズである計算ウィンドウサイズとに基づいて、設定用の受信ウィンドウサイズである設定ウィンドウサイズを生成する生成部と、

前記往復遅延時間の平均値を前記制御装置に送信する送信部と、

前記制御装置が前記送信部が送信した往復遅延時間及び他の受信装置が送信した往復遅延時間に基づいて設定用に算出した受信ウィンドウサイズである制御設定ウィンドウサイズを受信する受信部と、

30

前記制御設定ウィンドウサイズ又は前記設定ウィンドウサイズが設定された応答パケットを前記送信装置に送信する設定部と

を有することを特徴とする受信装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理システム、情報処理システムの制御方法及び受信装置に関する。

【背景技術】

【0002】

近年、複数の情報処理装置を複数のスイッチ装置で接続して構築される情報処理システムがデータセンター等で利用されている。図34は、データセンターにおける情報処理システムの一例を示す図である。図34において、 $S_0 \sim S_n$ は送信側エンドノードを示し、 $R_0 \sim R_n$ は受信側エンドノードを示す。ここで、エンドノードは情報処理装置であり、送信側エンドノードはパケットを送信する情報処理装置であり、受信側エンドノードは送信側エンドノードが送信したパケットを受信する情報処理装置である。

40

【0003】

$SW\#00 \sim SW\#05$ は、ネットワークを構築するスイッチ装置である。 $SW\#00 \sim SW\#02$ は、エンドノードに接続されるスイッチ装置であり、リーフスイッチと呼ばれる。一方、 $SW\#03 \sim SW\#05$ は、リーフスイッチに接続されるスイッチ装置であり、スパインスイッチと呼ばれる。送信側エンドノードが送信したパケットは、リーフス

50

イチ及びスパインスイッチを経由して受信側エンドノードに送信される。

【 0 0 0 4 】

図 3 4 に示す情報処理システムでは、輻輳によるパケットの再送等を防ぐために、レート制御 (rate control) が行われる。図 3 5 は、レート制御を説明するための図である。レート制御には、受信側エンドノードによるフロー制御 (flow control) と送信側エンドノードによる輻輳制御 (congestion control) がある。

【 0 0 0 5 】

フロー制御では、受信側エンドノードで受け取ることのできるパケットの数を示す R W I N (Receive W I N d o w) を受信側エンドノードがフロー毎に A C K で明示的に送信側エンドノードに通知する。図 3 5 では、フロー # 1 については、受信バッファ 8 9 に 2 パケット格納されているので、R W I N は 3 パケットである。また、フロー # 2 については、受信バッファ 8 9 に 1 パケット格納されているので、R W I N は 4 パケットである。なお、図 3 5 では、パケットは「p t k」と表される。

10

【 0 0 0 6 】

ネットワークのスループットあるいは R T T (Round Trip Time : 伝送遅延時間) からフロー毎に理論的な最大 R W I N を計算することが可能である。なお、ここでは、R W I N をパケット数で示したが、パケット数にパケットサイズを乗じることでデータ量が得られることから、R W I N はバイト数等のデータ量で表される場合もある。

【 0 0 0 7 】

輻輳制御は、送信側エンドノードが輻輳などネットワークの状態をパケットロスやタイムアウトから推定して受信側エンドノードに送るパケットの数を示す輻輳ウィンドウ (congestion window) をフロー毎に制御する。受信側エンドノードは、輻輳ウィンドウが知らせられないので、単位時間当たりのパケット数などで推定する必要がある。

20

【 0 0 0 8 】

図 3 5 では、送信側エンドノードは、フロー # 1 については、輻輳ウィンドウに基づいて 2 個のパケットを送信し、フロー # 2 については、輻輳ウィンドウに基づいて 3 個のパケットを送信する。

【 0 0 0 9 】

T C P (Transmission Control Protocol) を用いた通信では、フロー制御と輻輳制御が同時に働いて送信側エンドノードで送信するパケットの数は R W I N と輻輳ウィンドウのうちの小さい方のウィンドウで決まる。一般的には、R W I N と輻輳ウィンドウであるが、もし、R W I N < 輻輳ウィンドウであると、受信側エンドノードによる明示的なレート制御が可能となる。

30

【 0 0 1 0 】

なお、第 1 の通信システムを用いた通信中の通信品質の低下に応じて、より低速な第 2 の通信システムにハンドダウンした場合に、通信相手に通知する受信ウィンドウサイズを縮小することで、スループットを向上する従来技術がある。

【 0 0 1 1 】

また、第 1 の伝送区間から伝送遅延がより大きい第 2 の伝送区間にパケットを送信する場合に、第 2 の伝送区間の往復遅延時間に基づいて最大送信可能データ量を増加又は減少させることで、第 2 の伝送区間の輻輳による伝送効率の低下を防止する従来技術がある。

40

【 0 0 1 2 】

また、往復遅延時間、パス最大転送単位、回線速度等のパラメータを収集し、パラメータに基づいてファイル転送の最適ウィンドウサイズを計算することで、スループットを向上する従来技術がある。

【 0 0 1 3 】

また、時間を 2 つのロットに分けて、第 1 ロットではスループットを推定し、推定したスループットと受信側だけで制限を受ける場合の期待スループットとの比率に基づいて第 2 ロットでフロー毎にレート調整を行う従来技術がある。

【 先行技術文献 】

50

【特許文献】

【0014】

【特許文献1】特開2011-176540号公報

【特許文献2】特開2003-32295号公報

【特許文献3】特開2001-195326号公報

【非特許文献】

【0015】

【非特許文献1】Haitao Wu, Zhenqian Feng, Chuanxiong Guo, Yongguang Zhang, "ICTCP: Incast Congestion Control for TCP in Data Center Networks," ACM CoNEXT 2010, November 30 - December 3 2010, Philadelphia, USA.

10

【発明の概要】

【発明が解決しようとする課題】

【0016】

しかしながら、図34に示した情報処理システムにおいてインキャスト（incast）が検出されると、受信側エンドノード間でエレファント（elephant）フロー間の競合が発生するという問題がある。ここで、インキャストとは、例えば、データセンターにおいて同期する複数の送信側エンドノードが並列に受信側エンドノードにデータを送信した場合に発生する輻輳である。

【0017】

また、エレファントフローとは、スループットに敏感でライフタイムが長くボリュームが大きいフローである。ライフタイムが短くボリュームが小さいフローはマウス（mice）フローと呼ばれる。マウスフローは、例えば、ボリュームが1Mバイト以下、ライフタイムが10秒以下と定義される。データセンターでは、ボリュームについては90%がエレファントフローで10%がマウスフローであり、フローの数については10%がエレファントフローで90%がマウスフローである。

20

【0018】

図34では、 $n + 1$ 台の送信側エンドノードが接続するSW#00の送信側でインキャストが発生し、 $n + 1$ 台の受信側エンドノードが接続するSW#02の受信側でエレファントフロー間の競合が発生している。

【0019】

本発明は、1つの側面では、受信側エンドノード間でエレファントフロー間の競合の発生を防ぎ、エレファントフローを効率的に処理することを目的とする。

30

【課題を解決するための手段】

【0020】

本願の開示する情報処理システムは、1つの態様において、パケットを送信する送信装置と、複数のスイッチ装置と、前記送信装置が送信したパケットを前記複数のスイッチ装置を介して受信する複数の受信装置とを有する。また、本願の開示する情報処理システムは、前記送信装置と前記複数の受信装置を制御する制御装置とを有する。各受信装置は、監視部と、測定部と、計算部と、生成部と、設定部とを有する。前記監視部は、前記送信装置が有する輻輳ウィンドウのサイズの増加を監視する。前記測定部は、前記送信装置に対する応答パケットを前記送信装置に送信してから、前記応答パケットに対する次のデータを前記送信装置から受信するまでの時間である往復遅延時間を測定する。前記計算部は、前記監視部が、前記輻輳ウィンドウのサイズの増加が定常状態になったことを監視した場合、パケットを受信するフローのうちボリュームが所定サイズ以上の対象フローの最小スループットを計算する。前記生成部は、過去に応答パケットに設定した受信ウィンドウのサイズである設定済ウィンドウサイズと、前記最小スループットから計算した受信ウィンドウのサイズである計算ウィンドウサイズとに基づいて、設定ウィンドウサイズを生成する。前記設定部は、前記設定ウィンドウサイズが設定された応答パケットを前記送信装置に送信する。また、前記制御装置は、算出部と通信部とを有する。前記算出部は、前記複数の受信装置から受信した平均往復遅延時間に基づいて、前記制御装置からの設定用の

40

50

受信ウィンドウサイズである制御設定ウィンドウサイズを算出する。前記通信部は、前記制御設定ウィンドウサイズが設定された制御パケットを前記複数の受信装置の各々に送信する。

【発明の効果】

【0021】

1実施態様によれば、エレファントフローを効率的に処理することができる。

【図面の簡単な説明】

【0022】

【図1】図1は、実施例1に係る情報処理システムの構成を示す図である。

【図2】図2は、RWIN制御部の実現方法を説明するための図である。

10

【図3】図3は、通信階層におけるRWIN制御部の位置づけを示す図である。

【図4】図4は、RWIN制御部の構成を示す図である。

【図5】図5は、ウィンドウテーブルがエレファントフロー毎に記憶する項目の一例を示す図である。

【図6】図6は、5タプル及びDSCPの格納場所を示す図である。

【図7】図7は、RTTの測定範囲を示す図である。

【図8】図8は、輻輳回避状態を説明するための図である。

【図9】図9は、RTTテーブルがRTT毎に記憶する項目の一例を示す図である。

【図10】図10は、スループット測定スロット及びRWIN制御スロットを説明するための図である。

20

【図11】図11は、RWINテーブルがエレファントフロー毎に記憶する項目の一例を示す図である。

【図12】図12は、setRWINの設定場所を示す図である。

【図13】図13は、輻輳ウィンドウの監視処理のフローを示すフローチャートである。

【図14】図14は、RTTの測定処理のフローを示すフローチャートである。

【図15】図15は、RWINの設定処理のフローを示すフローチャートである。

【図16】図16は、調整レート算出例を示す図である。

【図17】図17は、スループット測定スロットの途中にフローが出現した場合の調整レート算出例を示す図である。

【図18】図18は、前提条件及び事前準備を説明するための図である。

30

【図19】図19は、レート制御のシーケンスを示す第1の図である。

【図20】図20は、レート制御のシーケンスを示す第2の図である。

【図21】図21は、実施例1に係るレート制御の効果を示す図である。

【図22】図22は、実施例2に係る情報処理システムの構成を示す図である。

【図23】図23は、Cosを説明するための図である。

【図24】図24は、実施例2に係るレート制御を説明するための図である。

【図25】図25は、RWIN送信用のPDUのフォーマットを示す図である。

【図26】図26は、平均RTT送信用のPDUのフォーマットを示す図である。

【図27】図27は、コントローラRWIN制御部の構成を示す図である。

【図28】図28は、CTRL__RWINテーブルがエレファントフロー毎に記憶する項目の一例を示す図である。

40

【図29】図29は、利用可能バンド幅を説明するための図である。

【図30】図30は、マトリクスの一例を示す図である。

【図31】図31は、コントローラによる処理のフローを示すフローチャートである。

【図32】図32は、レート制御のシーケンスを示す図である。

【図33】図33は、実施例1及び2に係るRWIN制御プログラムを実行するコンピュータの構成を示す図である。

【図34】図34は、データセンターにおける情報処理システムの一例を示す図である。

【図35】図35は、レート制御を説明するための図である。

【発明を実施するための形態】

50

【 0 0 2 3 】

以下に、本願の開示する情報処理システム、情報処理システムの制御方法及び受信装置の2つの実施例を図面に基づいて詳細に説明する。実施例1では、1つの受信側エンドノードでのレート制御を説明し、実施例2では、複数の受信側エンドノード間でのレート制御を説明する。なお、これらの実施例は開示の技術を限定するものではない。

【 実施例 1 】

【 0 0 2 4 】

まず、実施例1に係る情報処理システムの構成について説明する。図1は、実施例1に係る情報処理システムの構成を示す図である。図1に示すように、情報処理システム1は、 $S_0 \sim S_n$ で表される $n + 1$ 台の送信側エンドノード2と、 $R_0 \sim R_n$ で表される $n + 1$ 台の受信側エンドノード3と、 $SW \# 00 \sim SW \# 05$ で表される6台のスイッチ装置4とを有する。なお、送信側エンドノード2と受信側エンドノード3の数は異なってもよい。また、スイッチ装置4の数は、より多くでも少なくてもよい。

10

【 0 0 2 5 】

送信側エンドノード2は、情報処理装置であり、スイッチ装置4を介してパケットを受信側エンドノード3に送信する。受信側エンドノード3は、情報処理装置であり、送信側エンドノード2から送信されたパケットをスイッチ装置4を介して受信する。スイッチ装置4は、パケットを中継する装置である。

【 0 0 2 6 】

$SW \# 00 \sim SW \# 02$ は、リーフスイッチであり、 $SW \# 03 \sim SW \# 05$ はスパインスイッチである。 $SW \# 00$ は $S_0 \sim S_n$ に接続され、 $SW \# 02$ は $R_0 \sim R_n$ に接続される。 $SW \# 00 \sim SW \# 02$ は、それぞれ $SW \# 03 \sim SW \# 05$ に接続される。

20

【 0 0 2 7 】

受信側エンドノード3は、 $RWIN$ 制御部5を有する。 $RWIN$ 制御部5は、各エレメントフローのスループットを定期的に測って $RWIN$ を制御してレートを調整する。図2は、 $RWIN$ 制御部5の実現方法を説明するための図である。図2では、1台の受信側エンドノード3が1台のサーバとして動作する場合と、1台の受信側エンドノード3が m 台の仮想サーバとして動作する場合の $RWIN$ 制御部5の実現方法が示されている。なお、図2では、 $m = 4$ である。

【 0 0 2 8 】

1台の受信側エンドノード3が1台のサーバとして動作する場合には、 $RWIN$ 制御部5は $OS8$ (Operating System)の一部としてソフトウェアにより実現される。 $OS8$ は、 NIC (Network Interface Card) 7を用いて送信側エンドノード2等の他の装置と通信を行う。また、送信側エンドノード2等の他の装置で実行される $APP6$ (Application: アプリケーション)と通信を行う $APP6$ が $OS8$ の制御の下で実行される。

30

【 0 0 2 9 】

1台の受信側エンドノード3が m 台の仮想サーバとして動作する場合には、 m 個の $OS8$ と、仮想サーバの通信を制御する仮想スイッチ9が受信側エンドノード3で動作し、 $RWIN$ 制御部5は仮想スイッチ9の一部としてソフトウェアにより実現される。仮想スイッチ9は、 $NIC7$ を用いて送信側エンドノード2等の他の装置と通信を行う。

40

【 0 0 3 0 】

図3は、通信階層における $RWIN$ 制御部5の位置づけを示す図である。図3に示すように、 $RWIN$ 制御部5は、 TCP 、 UDP (User Datagram Protocol) が属するトランスポート層 ($L4$) の上位層において動作する。

【 0 0 3 1 】

次に、 $RWIN$ 制御部5の構成について説明する。図4は、 $RWIN$ 制御部5の構成を示す図である。図4に示すように、 $RWIN$ 制御部5は、ウィンドウテーブル51と、ウィンドウ監視部52と、 RTT テーブル53と、 RTT 測定部54と、 RTT タイマ55と、 $AvgRTT$ 計算部56と、 NT タイマ57と、スループット測定部58とを有する。また、 $RWIN$ 制御部5は、スループット計算部59と、最小スループット計算部60

50

と、RWINテーブル61と、RWIN計算部62と、RWIN比較部63と、RWINタイマ64と、setRWIN計算部65と、RWIN設定部66とを有する。また、RWIN制御部5は、タイムスロット設定部67とコントローラ通信部68とを有する。

【0032】

ウィンドウテーブル51は、輻輳ウィンドウの監視に用いられる情報を記憶する。図5は、ウィンドウテーブル51がエレファントフロー毎に記憶する項目の一例を示す図である。図5に示すように、ウィンドウテーブル51は、「インデックス」、「5タプル」、「前パケット数」、「現パケット数」、「RTT」及び「初期ACK」をエレファントフロー毎に記憶する。

【0033】

「インデックス」は、エレファントフローのインデックスである。「5タプル」は、フローを特定するための情報であり、送信元IP (Internet Protocol)、送信先IP、送信元ポート、送信先ポート及びプロトコルである。送信元IPは、パケットの送信元のIPアドレスである。送信先IPは、パケットの送信先のIPアドレスである。送信元ポートは、パケットの送信元のポート番号である。送信先ポートは、パケットの送信先のポート番号である。プロトコルは、フローの通信プロトコルである。

【0034】

図6は、5タプル及びDSCP (Differentiated Services Code Point) の格納場所を示す図である。図6に示すように、送信元ポート及び送信先ポートはTCPヘッダに含まれる。また、送信元IP、送信先IP、プロトコル及びDSCPはIPヘッダに含まれる。なお、DSCPは、フローがエレファントフローであるか否かを示す情報の格納場所として使用される。

【0035】

「前パケット数」は、前回のデータ転送で受信されたパケットの数である。「現パケット数」は、現在のデータ転送で受信されたパケットの数である。「RTT」は、ACKが送信されてから次のデータが受信されるまでの時間である。図7は、RTTの測定範囲を示す図である。図7に示すように、RTTは、 i 番目のデータに対してACKが送信された時刻から $(i + 1)$ 番目のデータが受信された時刻までの時間として測定される。「初期ACK」は、RTTを図るための最初のACKの時刻である。

【0036】

図4に戻って、ウィンドウ監視部52は、エレファントフローの輻輳ウィンドウが輻輳回避状態 (Congestion Avoidance) になるまでをウィンドウテーブル51を用いて監視する。ここで、輻輳回避状態とは、輻輳ウィンドウのサイズが前回から例えば1だけ増加するように緩やかに増加する定常状態である。

【0037】

図8は、輻輳回避状態を説明するための図である。図8において、横軸はデータの転送回数を示すRTTインデックスであり、縦軸は輻輳ウィンドウをパケット数で示す。図8に示すように、データ転送の初期は、輻輳が発生しないように輻輳ウィンドウが例えば指数関数的に激しく増加する。そして、輻輳ウィンドウがある程度の大きさになると増加が緩やかになる。このような輻輳ウィンドウの増加が緩やかになった定常状態が輻輳回避状態である。受信側エンドノード3は、各フローのスループットを輻輳回避状態で測定する。

【0038】

RTTテーブル53は、RTTの測定に用いられる情報を記憶する。RWIN制御部5は、RTTテーブル53をエレファントフロー毎に有する。図9は、RTTテーブル53がRTT毎に記憶する項目の一例を示す図である。図9に示すように、RTTテーブル53は、「インデックス」、「RTT」、「ACKタイム」、「ACK番号」及び「失敗数」をRTT毎に記憶する。

【0039】

「インデックス」は、RTTのインデックスである。「RTT」は、ACKが送信され

10

20

30

40

50

てから次のデータが受信されるまでの時間である。「ACKタイム」は、ACKが送信された時刻である。「ACK番号」は、ACKのシーケンス番号にMSS(Maximum Segment Size)が加えられた値である。MSSは、例えば1500である。「失敗数」は、ACK番号とデータのシーケンス番号であるデータ番号との比較で一致しなかった数である。

【0040】

RTT測定部54は、ACKが送信されてから次のデータが受信されるまでの時間をRTTテーブル53を用いてRTTとして測定する。RTTタイマ55は、RTTの平均値を計算するスロットを測るタイマである。ここで、スロットとは一定の大きさの時間区間である。RTT測定部54は、RTTタイマ55が0でない間はRTTを測定する。Avg RTT計算部56は、RTTタイマ55が0になるとRTTテーブル53を用いて各エレファントフローのRTTの平均値を計算する。

10

【0041】

NTタイマ57は、スループットを測定するスロットを測るタイマである。スループット測定部58は、輻輳回避状態になったエレファントフローのスループットをスループット測定スロットで測定する。ここで、スループット測定スロットとは、スループットを測定する時間帯である。RWIN制御部5は、時間をスループット測定スロットとRWIN制御スロットに2分する。

【0042】

図10は、スループット測定スロット及びRWIN制御スロットを説明するための図である。図10において、横軸のTimeは時間を表す。図10に示すように、時間は一定の時間間隔のタイムスロットに分割され、タイムスロットは交互にスループット測定スロットとRWIN制御スロットとなる。すなわち、スループット測定スロットがenableされている場合にはRWIN制御スロットがdisableされており、スループット測定スロットがdisableされている場合にはRWIN制御スロットがenableされている。例えば、タイムスロット(0)はスループット測定スロットであり、タイムスロット(1)はRWIN制御スロットであり、タイムスロット(2)はスループット測定スロットであり、タイムスロット(3)はRWIN制御スロットである。

20

【0043】

スループット計算部59は、各エレファントフローのスループットを計算する。最小スループット計算部60は、エレファントフローのスループットの中で最小スループットを計算する。

30

【0044】

RWINテーブル61は、RWINの計算に用いられる情報を記憶する。図11は、RWINテーブル61がエレファントフロー毎に記憶する項目の一例を示す図である。図11に示すように、RWINテーブル61は、「インデックス」、「バイト数」、「RTT」、「スループット」及び「calRWIN」をエレファントフロー毎に記憶する。

【0045】

「インデックス」は、エレファントフローのインデックスである。「バイト数」は、スループット測定スロットで測定されたバイト数である。「RTT」は、Avg RTT計算部56により計算された平均RTTである。「スループット」は、エレファントフローのスループットである。「calRWIN」は、最小スループットから計算されたRWINである。

40

【0046】

RWIN計算部62は、最小スループットとRTTからRWINを計算する。RWINの単位はバイトである。また、ここで使用されるRTTは、各エレファントフローのRTTを平均した値である。

【0047】

RWIN比較部63は、RWIN計算部62が計算したRWINと1パケット分を比較して大きい方の値をcalRWINとする。RWIN計算部62が計算したRWINがパ

50

ケットのサイズより小さい場合には、RWINが1パケット分ない場合であるので、RWIN比較部63は、1パケットの大きさをcalRWINとする。

【0048】

RWINタイマ64は、RWIN制御スロットを測るタイマである。setRWIN計算部65は、RWIN制御スロットでcalRWINとから調整レートsetRWINを計算する。具体的には、setRWIN計算部65は、 $setRWIN = calRWIN + calRWIN \times \alpha$ によりsetRWINを計算する。ここで、 α は以下の式(1)で定義される。

【数1】

$$\alpha = \frac{\text{MAX}_{i=0..I} |f_i - \text{Rate}|}{\text{Rate}} \quad \dots(1)$$

10

【0049】

式(1)で、Rateは1つ前のRWIN制御スロットで計算された調整レートであり、 f_i はi番目のエレファントフローの現在のスループットである。 $\text{MAX}_{i=0..I}$ は、0番目のエレファントフローからI番目のエレファントフローについて、 $|f_i - \text{Rate}|$ の最大値をとる。I+1がエレファントフローの数である。また、 f_i には、最小スループットとして選択されたエレファントフローのスループットは含まれない。

【0050】

式(1)に示すように、 α は、「エレファントフローの現在のスループットと1つ前の調整レートの差の最大値」と「1つ前の調整レート」との比率である。ただし、 α が1より大きい場合は1とする。 α が1に近い場合は、1つ前の調整レートと現在のエレファントフローのスループットとの差が大きい場合であり、 α が0に近い場合は、1つ前の調整レートと現在のエレファントフローのスループットとの差が小さい場合である。また、setRWIN計算部65は、複数のエレファントフローがあって、 $\alpha = 0$ である場合には、例えば、強制的に $\alpha = 0.1$ とする。setRWIN計算部65は、強制的に α の値を0.1以外の値にしてもよい。

20

【0051】

RWIN設定部66は、setRWIN計算部65により計算されたsetRWINをACKに設定する。図12は、setRWINの設定場所を示す図である。図12に示すように、setRWINは、TCPヘッダの中のRWINフィールドに設定される。

30

【0052】

タイムスロット設定部67は、スループット測定スロット、RWIN制御スロット及びRTTスロットの時間間隔をそれぞれNTタイマ57、RWINタイマ64及びRTTタイマ55に設定する。時間間隔は、例えば、1ms(ミリ秒)、100ms、RTT等である。スループット測定スロットとRWIN制御スロットの時間間隔は異なってもよい。

【0053】

コントローラ通信部68は、RTTの平均値をコントローラに送信し、RWINの設定値を受信してRWIN設定部66に渡す。なお、コントローラについては、実施例2で説明する。

40

【0054】

次に、輻輳ウィンドウの監視処理のフローについて説明する。図13は、輻輳ウィンドウの監視処理のフローを示すフローチャートである。図13に示すように、ウィンドウ監視部52は、エレファントフローのSYNパケットを受信し(ステップS1)、SYNパケットから5タプルを取り出す(ステップS2)。

【0055】

そして、ウィンドウ監視部52は、5タプルの情報を含むエレファントフローの情報をウィンドウテーブル51に記録し(ステップS3)、ACKが送信された時刻をウィンドウテーブル51に初期ACKとして記録する(ステップS4)。そして、ウィンドウ監視部52は、パケットを受信したか否かを判定し(ステップS5)、受信していない場合に

50

は、パケットの受信を待つ。

【0056】

一方、パケットを受信した場合には、ウィンドウ監視部52は、受信したパケットの5タプルでウィンドウテーブル51を検索してエレファントフローを特定する(ステップS6)。また、ウィンドウ監視部52は、RTTフラグがfalseであるか否かを判定する(ステップS7)。ここで、RTTフラグは、RTTの測定が行われたか否かを示すフラグであり、falseは測定が行われていないことを示す。また、ステップS7における「=」は、「等しい」を表す記号である。

【0057】

そして、RTTフラグがfalseでない場合には、RTTが測定されているので、ウィンドウ監視部52は、ステップS10へ進む。一方、RTTフラグがfalseの場合には、RTT測定部54が、ACK データ間のRTTを測定し(ステップS8)、RTTフラグをtrueに設定する(ステップS9)。

10

【0058】

そして、ウィンドウ監視部52は、ステップS6で特定したエレファントフローに関するウィンドウテーブル51の情報に基づいて、パケット受信間隔がRTTより小さいか否かを判定する(ステップS10)。その結果、パケット受信間隔がRTTより小さい場合には、データ受信中なので、ウィンドウ監視部52は、ウィンドウテーブル51の現パケット数に1を加え(ステップS11)、ステップS5に戻る。

【0059】

20

一方、パケット受信間隔がRTTより小さくない場合には、ウィンドウ監視部52は、ウィンドウテーブル51の現パケット数が前パケット数より大きいか否かを判定する(ステップS12)。その結果、現パケット数が前パケット数より大きくない場合には、輻輳回避状態ではないので、ウィンドウ監視部52は、現パケット数を前パケット数とし(ステップS13)、現パケット数を0とする(ステップS14)。そして、ウィンドウ監視部52は、ステップS5に戻る。

【0060】

一方、現パケット数が前パケット数より大きい場合には、ウィンドウ監視部52は、前パケット数から現パケット数への増加値が1であるか否かを判定し(ステップS15)、1でない場合には、輻輳回避状態ではないので、ステップS13へ移動する。一方、増加値が1である場合には、ウィンドウ監視部52は、ウィンドウ数を1増加し(ステップS16)、ウィンドウ数が2より大きいか否かを判定する(ステップS17)。ここで、ウィンドウ数は、前パケット数から現パケット数への増加値が1である回数を表す。すなわち、ウィンドウ数は、ウィンドウの状態が輻輳回避になった回数を表す。

30

【0061】

そして、ウィンドウ数が2より大きい場合には、ウィンドウ監視部52は、輻輳回避状態になったと判定し、処理を終了する。一方、ウィンドウ数が2より大きくない場合には、輻輳回避状態になったとはまだ判定できないので、ウィンドウ監視部52は、ステップS13へ移動する。

【0062】

40

このように、ウィンドウ監視部52は、ウィンドウテーブル51を用いて1回のデータ転送におけるパケット数を監視することで、エレファントフローが輻輳回避状態になったか否かを判定することができる。

【0063】

次に、RTTの測定処理のフローについて説明する。図14は、RTTの測定処理のフローを示すフローチャートである。図14に示すように、AvgRTT計算部56は、対象のエレファントフローについてウィンドウテーブル51から5タプルを取得してRTTテーブル53を作成する(ステップS21)。

【0064】

そして、AvgRTT計算部56は、RTTタイマ55が0でないか否かを判定し(ス

50

テップS 2 2)、0でない場合には、R T Tを測定する時間帯なので、パケットがA C Kであるか否かを判定する(ステップS 2 3)。その結果、パケットがA C Kでない場合には、A v g R T T計算部5 6は、A C Kが来るまで待つ。

【0 0 6 5】

一方、パケットがA C Kである場合には、A v g R T T計算部5 6は、A C Kの時刻をR T Tテーブル5 3のA C Kタイムに記録し(ステップS 2 4)、A C K番号にM S Sを加える(ステップS 2 5)。そして、A v g R T T計算部5 6は、データパケットを受信したか否かを判定し(ステップS 2 6)、データパケットを受信しない場合には、受信するまで待つ。

【0 0 6 6】

一方、データパケットを受信した場合には、A v g R T T計算部5 6は、A C K番号とデータ番号が等しいか否かを判定する(ステップS 2 7)。その結果、等しくない場合には、A v g R T T計算部5 6は、失敗数に1を加え(ステップS 2 8)、失敗数が3より大きいか否かを判定する(ステップS 2 9)。そして、A v g R T T計算部5 6は、失敗数が3より大きくない場合には、ステップS 2 6に戻り、失敗数が3より大きい場合には、ステップS 2 2に戻る。

【0 0 6 7】

一方、A C K番号とデータ番号が等しい場合には、R T T測定部5 4が、A C K時刻とデータパケットを受信した時刻からR T Tを測定し(ステップS 3 0)、R T TをR T Tテーブル5 3に記録する(ステップS 3 1)。そして、A v g R T T計算部5 6は、ス

【0 0 6 8】

一方、R T Tタイマ5 5が0である場合には、A v g R T T計算部5 6は、R T Tテーブル5 3に記録された複数のR T Tから平均R T Tを計算し(ステップS 3 2)、コントローラ通信部6 8がコントローラに平均R T TをB P D U (Bridge Protocol Data Unit)化して送信する(ステップS 3 3)。なお、コントローラについては、実施例2で説明する。

【0 0 6 9】

そして、A v g R T T計算部5 6は、R W I Nテーブル6 1に平均R T Tを記録し(ステップS 3 4)、R T Tテーブル5 3をクリアし(ステップS 3 5)、ステップS 2 2に戻る。

【0 0 7 0】

このように、A v g R T T計算部5 6がR T Tの平均値を算出してR W I Nテーブル6 1に記録することによって、スループット計算部5 9はR T Tの平均値を用いてエレファントフローのスループットを算出することができる。

【0 0 7 1】

次に、R W I Nの設定処理のフローについて説明する。図1 5は、R W I Nの設定処理のフローを示すフローチャートである。図1 5に示すように、R W I N制御部5は、スロットフラグがt r u eであるか否かを判定する(ステップS 4 1)。ここで、スロットフラグは、スループット測定スロットであるかR W I N制御スロットであることを示すフラグであり、t r u eである場合にはスループット測定スロットであることを示し、f a l s eである場合にはR W I N制御スロットであることを示す。

【0 0 7 2】

そして、スロットフラグがt r u eでない場合には、R W I N制御部5は、R W I Nタイマ6 4が0でないか否かを判定し(ステップS 4 2)、0でない場合には、R W I Nフラグがf a l s eであるか否かを判定する(ステップS 4 3)。ここで、R W I Nフラグは、現在のR W I N制御スロットでR W I N制御を行ったか否かを示すフラグであり、t r u eである場合にはR W I N制御を行ったことを示し、f a l s eである場合にはR W I N制御を行っていないことを示す。そして、R W I Nフラグがf a l s eでない場合には、R W I N制御部5は、ステップS 4 2に戻る。

10

20

30

40

50

【0073】

一方、RWINフラグがfalseである場合には、スループット計算部59が各エレファントフローのスループットを計算し(ステップS44)、RWIN計算部62が最スループットからRWINを計算する(ステップS45)。そして、RWIN比較部63がcalRWINを計算する。そして、setRWIN計算部65がsetRWINを計算し(ステップS46)、RWIN設定部66が各エレファントフローに対してsetRWINをACKに設定する(ステップS47)。そして、RWIN制御部5は、RWINフラグをtrueに設定し(ステップS48)、ステップS42に戻る。

【0074】

また、ステップS42においてRWINタイマ64が0である場合には、RWIN制御スロットからスループット測定スロットに変更するために、スロットフラグにtrueを設定し(ステップS49)、ステップS41へ戻る。

【0075】

また、ステップS41においてスロットフラグがtrueである場合には、RWIN制御部5は、NTタイマ57が0でないか否かを判定する(ステップS50)。その結果、0でない場合には、スループット測定部58がスループットを測定し(ステップS51)、RWIN制御部5は、ステップS50へ戻る。一方、0である場合には、スループット測定スロットからRWIN制御スロットに変更するために、RWIN制御部5は、スロットフラグをfalseに設定し(ステップS52)、ステップS41に戻る。

【0076】

このように、RWIN制御スロットにおいてRWIN設定部66が各エレファントフローに対してsetRWINをACKに設定することで、RWIN制御部5は各エレファントフローに対してスループットが均等になるようにフロー制御を行うことができる。

【0077】

次に、調整レート算出例について図16及び図17を用いて説明する。図16は、調整レート算出例を示す図である。図16において、E#0~E#2は、輻輳回避状態にあるエレファントフローである。図16に示すように、タイムスロット(0)では2つのエレファントフローE#0及びE#1があり、タイムスロット(1)でエレファントフローE#2が出現したとする。

【0078】

スループット測定スロットであるタイムスロット(0)でスループット測定部58がスループットを測定し、RWIN制御スロットであるタイムスロット(1)でスループット計算部59がスループットを計算する。ここでは、スループット計算部59が計算したE#0及びE#1のスループットがそれぞれ1Gbps(ギガビット/秒)及び1.2Gbpsであったとする。

【0079】

すると、最小スループットは1Gbpsであり、前回の調整レートはないので、をデフォルト値の0として、setRWIN=calRWIN=1Gbpsである。ここで、RTT=0.5msとすると、setRWIN=1Gbps×0.5×10³s=500Kビット=62.5Kバイトとなる。なお、RTTの値は、E#0とE#1のRTTの平均値である。

【0080】

その後、スループット測定スロットであるタイムスロット(2)でスループット測定部58がスループットを測定し、RWIN制御スロットであるタイムスロット(3)でスループット計算部59がスループットを計算する。ここでは、スループット計算部59が計算したE#0、E#1及びE#2のスループットがそれぞれ1Gbps、1.1Gbps及び1.2Gbpsであったとする。

【0081】

すると、最小スループットは1Gbpsであり、=0.2Gbps/1Gbps=0.2であり、calRWIN=1Gbpsであるので、setRWIN=1Gbps+1

10

20

30

40

50

$Gbps \times 0.2 / 2 = 1.1 Gbps$ である。ここで、 $RTT = 0.5 ms$ とすると、 $setRWIN = 1.1 Gbps \times 0.5 \times 10^3 s = 550 K$ ビット $= 68.75 K$ バイトとなる。なお、 RTT の値は、 $E \# 0 \sim E \# 2$ の RTT の平均値である。このように、調整レートより測定スループットが大きいと、 $RWIN$ 制御部5は、調整レートを上げる。

【0082】

なお、タイムスロット(3)でスループット計算部59が計算した $E \# 0$ 、 $E \# 1$ 及び $E \# 2$ のスループットがそれぞれ $1 Gbps$ 、 $1.1 Gbps$ 及び $0.8 Gbps$ であったとすると、最小スループットは $0.8 Gbps$ であり、 $\alpha = 0.1 Gbps / 1 Gbps = 0.1$ であり、 $calRWIN = 0.8 Gbps$ であるので、 $setRWIN = 0.8 Gbps + 0.8 Gbps \times 0.1 / 2 = 0.84 Gbps$ である。ここで、 $RTT = 0.5 ms$ とすると、 $setRWIN = 0.84 Gbps \times 0.5 \times 10^3 s = 420 K$ ビット $= 52.5 K$ バイトとなる。この場合、調整レートは変動する。

10

【0083】

図17は、スループット測定スロットの途中にフローが出現した場合の調整レート算出例を示す図である。図17に示すように、タイムスロット(0)ではエレファントフロー $E \# 0$ があり、タイムスロット(2)でエレファントフロー $E \# 1$ が出現したとする。

【0084】

スループット測定スロットであるタイムスロット(0)でスループット測定部58がスループットを測定し、 $RWIN$ 制御スロットであるタイムスロット(1)でスループット計算部59がスループットを計算する。ここでは、スループット計算部59が計算した $E \# 0$ のスループットが $1 Gbps$ であったとする。

20

【0085】

すると、最小スループットは $1 Gbps$ であり、前回の調整レートはないので、 α をデフォルト値の0として、 $setRWIN = calRWIN = 1 Gbps$ である。ここで、 $RTT = 0.5 ms$ とすると、 $setRWIN = 1 Gbps \times 0.5 \times 10^3 s = 500 K$ ビット $= 62.5 K$ バイトとなる。

【0086】

その後、スループット測定スロットであるタイムスロット(2)でスループット測定部58がスループットを測定すると、途中で $E \# 1$ が出現する。そして、 $RWIN$ 制御スロットであるタイムスロット(3)でスループット計算部59がスループットを計算する。ここでは、スループット計算部59が計算した $E \# 0$ 及び $E \# 1$ のスループットがそれぞれ $1 Gbps$ 及び $0.6 Gbps$ であったとする。

30

【0087】

すると、途中で出現した $E \# 1$ のスループットは使用せず、最小スループットは $1 Gbps$ であり、 $\alpha = 0$ であるので、 $setRWIN = calRWIN = 1 Gbps$ である。 $RTT = 0.5 ms$ とすると、 $setRWIN = 1 Gbps \times 0.5 \times 10^3 s = 500 K$ ビット $= 62.5 K$ バイトとなる。そして、 $RWIN$ 制御部5は、途中で出現した $E \# 1$ の $RWIN$ も $62.5 K$ バイトに設定する。

【0088】

そして、 $RWIN$ 制御スロットであるタイムスロット(5)でスループット計算部59がスループットを計算する。ここでは、スループット計算部59が計算した $E \# 0$ 及び $E \# 1$ のスループットがそれぞれ $1 Gbps$ 及び $0.8 Gbps$ であったとする。

40

【0089】

すると、最小スループットは $0.8 Gbps$ であり、 $\alpha = 0$ である。ここで、複数のエレファントフローがあり、 $\alpha = 0$ であるので、 $setRWIN$ 計算部65は、強制的に $\alpha = 0.1$ とする。 $setRWIN$ 計算部65は、強制的に α の値を0.1以外の値にしてもよい。すると、 $setRWIN = 0.8 Gbps + 0.8 Gbps \times 0.1 / 2 = 0.84 Gbps$ である。ここで、 $RTT = 0.5 ms$ とすると、 $setRWIN = 0.84 Gbps \times 0.5 \times 10^3 s = 420 K$ ビット $= 52.5 K$ バイトとなる。このように、

50

スループット測定スロットの途中でフローが出現した場合には、途中で出現したフローのスループットは、次の次のRWIN制御スロットで使われる。

【0090】

次に、レート制御のシーケンスについて図18～図20を用いて説明する。なお、図19は図16に示した例に対応し、図20は図17に示した例に対応する。図18は、前提条件及び事前準備を説明するための図である。図18において、実線の矢印はエレファントフローのデータを示し、破線の矢印はエレファントフローのACKを示し、点線の矢印はミスフローを示す。

【0091】

図18に示すように、前提条件として、送信側エンドノード $S_0 \sim S_2$ は、エレファントフローの packets 及びミスフローの packets を受信側エンドノード R_0 に送信する。そして、 $SW\#00$ でインキャストが発生し、 $SW\#02$ でエレファントフロー間の競合が発生する。

【0092】

また、ウィンドウ監視部52は、エレファントフローの輻輳ウィンドウが輻輳回避状態になることを監視する(ステップt1)。そして、RTT測定部54は、ACKとデータの間のRTTを測定する(ステップt2, ステップt3)。RTTの測定は、エレファントフローの通信が終わるまで実施される。

【0093】

図19は、レート制御のシーケンスを示す第1の図である。図19は、RWIN制御スロットでエレファントフローE#2が出現する場合を示す。なお、図19及び図20において、 T_n はタイムスロット(n)を表し、 T_0 、 T_2 及び T_4 はスループット測定スロットを表し、 T_1 、 T_3 及び T_5 はRWIN制御スロットを表す。また、実線の矢印はデータ packets を示し、破線の矢印はRWIN制御があるACKを示し、点線の矢印はRWIN制御がないACKを示す。

【0094】

図19に示すように、まず、図18に示した前提条件及び事前準備が整う(ステップt11)。すなわち、インキャストが発生し、エレファントフローの輻輳ウィンドウが輻輳回避状態になる。そして、 T_0 において、スループット測定部58がエレファントフローE#0及びE#1のスループットを測定する(ステップt12)。

【0095】

そして、 T_1 において、setRWIN計算部65がsetRWIN(62.5KB)を計算し、RWIN設定部66がE#0とE#1のACKにsetRWINを設定して送信する(ステップt13)。setRWINの計算には、E#0及びE#1のスループットとしてそれぞれ1Gbps及び1.1Gbpsが用いられる。また、 T_1 において、E#2が輻輳回避状態になって出現する(ステップt14)。ただし、E#1のACKにはsetRWINは設定されない。

【0096】

そして、 T_2 において、スループット測定部58がエレファントフローE#0、E#1及びE#2のスループットを測定する(ステップt15)。そして、 T_3 において、setRWIN計算部65がsetRWIN(68.75KB)を計算し、RWIN設定部66がE#0とE#1とE#2のACKにsetRWINを設定して送信する(ステップt16)。setRWINの計算には、E#0、E#1及びE#2のスループットとしてそれぞれ1Gbps、1.1Gbps及び1.2Gbpsが用いられる。

【0097】

なお、setRWINの計算に、E#0、E#1及びE#2のスループットとしてそれぞれ1Gbps、1.1Gbps及び0.8Gbpsが用いられた場合には、setRWINの値は52.5KBとなる。

【0098】

図20は、レート制御のシーケンスを示す第2の図である。図20は、スループットの

10

20

30

40

50

測定途中でエレファントフローが出現する場合を示す。図20に示すように、まず、図18に示した前提条件及び事前準備が整う(ステップt21)。そして、 T_0 において、スループット測定部58がエレファントフローE#0のスループットを測定する(ステップt22)。

【0099】

そして、 T_1 において、setRWIN計算部65がsetRWIN(62.5KB)を計算し、RWIN設定部66がE#0のACKにsetRWINを設定して送信する(ステップt23)。setRWINの計算には、E#0のスループットとして1Gbpsが用いられる。

【0100】

そして、 T_2 において、スループット測定部58がエレファントフローE#0のスループットを測定する(ステップt24)。また、 T_2 において、エレファントフローE#1が輻輳回避状態になって出現する(ステップt25)。そして、 T_3 において、setRWIN計算部65がsetRWIN(62.5KB)を計算し、RWIN設定部66がE#0のACKにsetRWINを設定して送信する(ステップt26)。setRWINの計算には、E#0のスループットとして1Gbpsが用いられる。ただし、setRWINの計算にはE#1のスループットは用いられず、E#1のACKにはsetRWINは設定されない。

【0101】

そして、 T_4 において、スループット測定部58がエレファントフローE#0及びE#1のスループットを測定する(ステップt27)。そして、 T_5 において、setRWIN計算部65がsetRWIN(52.5KB)を計算し、RWIN設定部66がE#0とE#1のACKにsetRWINを設定して送信する(ステップt28)。setRWINの計算には、E#0及びE#1のスループットとしてそれぞれ1Gbps及び0.8Gbpsが用いられる。また、setRWINを計算する際に $\alpha = 0$ となるが、複数のエレファントフローが存在するため、setRWIN計算部65は強制的に $\alpha = 0.1$ として、setRWINを少し大きい数値に設定する。

【0102】

図21は、実施例1に係るレート制御の効果を示す図である。図21は、7本の25.6MBのエレファントフローと50本の256KBのマイスフローを用いて10Gのリンクでインキャストを発生させ、各エレファントフローのスループットを従来技術A~Dと実施例1に係るレート制御で比較したものである。縦軸はスループットであり、横軸はエレファントフローのインデックスを示すフローインデックスである。図21に示すように、実施例1に係るレート制御は、従来技術と比較して、スループットが高い値で平均化されている。

【0103】

上述してきたように、実施例1では、最小スループット計算部60が、輻輳回避状態にある複数のエレファントフローのスループットから最小のスループットを計算する。また、RWIN計算部62が、最小のスループットに基づいてRWINを計算し、RWIN比較部63がcalRWINを計算する。そして、setRWIN計算部65は、前回の調整レートと現在のスループットの差の絶対値を最小スループット以外の全てのエレファントフローについて計算し、最も大きな値と前回の調整レートに基づいて β を計算する。そして、setRWIN計算部65は、calRWINと β に基づいてsetRWINを計算し、RWIN設定部66がsetRWINを各エレファントフローのACKに設定する。したがって、RWIN制御部5は、競合するエレファントフロー間でスループットを平均化することができる。

【0104】

また、実施例1では、RWIN制御部5は、仮想スイッチ9が存在する場合には仮想スイッチの一部として実現され、仮想スイッチ9が存在しない場合にはOS8の一部として実現される。したがって、受信側エンドノード3は、仮想スイッチ9が存在する仮想化シ

10

20

30

40

50

システムであるか否かにかかわらず、競合するエレファントフロー間でスループットを平均化することができる。

【0105】

また、実施例1では、RWIN比較部63は、RWIN計算部62が計算したRWINとパケットサイズとを比較し、大きい方をcalRWINとする。したがって、RWIN制御部5は、setRWINの大きさを1パケット以上とすることができる。

【0106】

また、実施例1では、複数のエレファントフローが存在して = 0 になった場合に、setRWIN計算部65は強制的に = 0.1 とする。したがって、RWIN制御部5は、setRWINを大きくすることができ、スループットを向上することができる。

10

【0107】

なお、実施例1では、をcalRWINに乘じる際に、を2で割ることとしたが、ネットワークの状態、特徴等を考慮して、を2で割らない、の平方根を用いる、を他の自然数で割る、を所定の指数で割ることとしてもよい。また、calRWINにを乘じる代わりにcalRWINをの平方根で割ってもよい。

【実施例2】

【0108】

ところで、上記実施例1では、1つの受信エンドノード3で複数の競合するエレファントフローのレートを制御する場合について説明したが、エレファントフローは複数の受信エンドノード3間で競合する場合もある。そこで、実施例2では、複数の受信エンドノード3間で複数のエレファントフローが競合する場合のレート制御について説明する。

20

【0109】

図22は、実施例2に係る情報処理システムの構成を示す図である。なお、ここでは説明の便宜上、図1に示した装置と同様の役割を果たす装置については同一符号を付すこととしてその詳細な説明を省略する。図22に示すように、情報処理システム1aは、n+1台の送信側エンドノード2と、n+1台の受信側エンドノード3と、6台のスイッチ装置4と、コントローラ10とを有する。

【0110】

スパインスイッチSW#03~SW#05は、パケットのヘッダを見てエレファントフローのパケットであれば、SYNパケット、FINパケット及びRSTパケットをスヌープしてコントローラ10に送信する。スパインスイッチSW#03~SW#05は、エレファントフローであるか否かをパケットのヘッダのDSCP、CoS(Class of Service)等を用いて判定する。図23は、CoSを説明するための図である。図23に示すように、CoSは、パケットのヘッダのVLANタグに含まれる。

30

【0111】

また、各受信エンドノード3は、エレファントフローのRTTの平均値を算出すると、算出した平均RTTをコントローラ10に送信する。

【0112】

コントローラ10は、スパインスイッチSW#03~SW#05から送信されたパケット及び各受信エンドノード3から送信された平均RTTに基づいて、複数の受信エンドノード3間で競合するエレファントフローのレート制御を行う。コントローラ10はRWINを計算して各受信側エンドノード3に送信し、各受信側エンドノード3はコントローラ10から送信されたRWINを用いてエレファントフローのレート制御を行う。

40

【0113】

図24は、実施例2に係るレート制御を説明するための図である。図24に示すように、コントローラ10はコントローラRWIN制御部5aを有する。コントローラRWIN制御部5aはネットワーク11からエレファントフローに関する情報を収集してRWINを計算し、計算したRWINを各受信側エンドノード3のRWIN制御部5に送信し、RWIN制御部5がレート制御を行う。各受信側エンドノード3のRWIN制御部5は、自身が計算したsetRWINよりコントローラ10から受信したRWINを優先してレ

50

ト制御を行う。

【0114】

このように、コントローラ10がRWINを計算することによって、実施例2に係る情報処理システム1aは受信側エンドノード3間で競合するエレファントフローのレート制御を公平に行うことができる。

【0115】

図25は、RWIN送信用のPDUのフォーマットを示す図である。図25に示すように、RWIN送信用のPDUでは、送信先アドレスとしてあらかじめ予約された「01-80-C2-00-00-XX」が用いられる。また、図25に示す2バイトのBPDUにRWINの値が設定される。

10

【0116】

図26は、平均RTT送信用のPDUのフォーマットを示す図である。図26に示すように、平均RTT送信用のPDUでは、送信先アドレスとしてあらかじめ予約された「01-80-C2-00-00-XX」が用いられる。また、図26に示す38バイトのBPDUに5タプルと平均RTTが設定される。

【0117】

図27は、コントローラRWIN制御部5aの構成を示す図である。図27に示すように、コントローラRWIN制御部5aは、5タプル検出部71と、CTRL_RWINテーブル72と、パス計算部73と、パス検索部74と、RTT測定部75と、スループット計算部76と、ABW計算部77とを有する。また、コントローラRWIN制御部5a

20

【0118】

5タプル検出部71は、SYNパケットのヘッダから5タプルを取り出す。CTRL_RWINテーブル72は、RWINの算出に用いられる情報を記憶する。図28は、CTRL_RWINテーブル72がエレファントフロー毎に記憶する項目の一例を示す図である。図28に示すように、CTRL_RWINテーブル72は、「インデックス」、「5タプル」、「パス」、「RTT」、「ABW」、「スループット」、「RWIN」及び「トポロジー」をエレファントフロー毎に記憶する。

【0119】

「インデックス」は、エレファントフローを識別する識別子である。「5タプル」は、エレファントフローを特定するための5タプルである。「パス」は、エレファントフローが通る端から端までのパスである。「RTT」は、ACKとデータ間の平均RTTである。

30

【0120】

「ABW」は、パス上で最小の利用可能バンド幅である。図29は、利用可能バンド幅を説明するための図である。図29に示すように、利用可能バンド幅は、リンク容量から使用中の帯域幅を引いた帯域幅である。例えば、リンク容量が10Gbpsで使用中の帯域幅が5Gbpsのとき、利用可能バンド幅は5Gbpsである。

【0121】

「スループット」は、エレファントフローのスループットである。「RWIN」は、ABWとスループットで計算されたRWINである。「トポロジー」は、スイッチ装置間のリンクの有無を示し、マトリクスで定義される。図30は、マトリクスの一例を示す図である。図30に示すように、スイッチ装置間にリンクがある場合には、マトリクスの対応する要素は「」であり、スイッチ装置間にリンクがない場合には、マトリクスの対応する要素は「-」である。

40

【0122】

図27に戻って、パス計算部73は、5タプル、マトリクス等のネットワーク情報を用いてエレファントフローの端から端までのパスを計算し、CTRL_RWINテーブル72に記録する。パス検索部74は、CTRL_RWINテーブル72を検索して同じリンクでエレファントフローが2以上であるかを判定する。

50

【 0 1 2 3 】

R T T測定部 7 5 は、受信側エンドノード 3 が送信した平均 R T Tを通信部 8 0 から受け取り、C T R L _ R W I Nテーブル 7 2 に記録する。スループット計算部 7 6 は、平均 R T Tとフロー数から各エレファントフローのスループットを計算する。

【 0 1 2 4 】

A B W計算部 7 7 は、リンク数とリンク容量から各リンクの利用可能バンド幅を計算し、計算した利用可能バンド幅の最小値 A B Wを C T R L _ R W I Nテーブル 7 2 に記録する。具体的には、A B W計算部 7 7 は、以下の式 (2) で A B Wを計算する。

【 数 2 】

$$ABW = \min_{l=0,L} \{Linkcap(l) - Linkused(l)\} \quad \dots(2)$$

10

式 (2) で、L i n k c a p はリンク容量であり、L i n k u s e d は使用中の帯域幅である。

【 0 1 2 5 】

R W I N計算部 7 8 は、A B W、各エレファントフローのスループット、複数の受信側エンドノード 3 から受信した平均 R T Tの平均値に基づいて R W I Nを計算する。具体的には、R W I N計算部 7 8 は、以下の式 (3) で R W I Nを計算する。

【 数 3 】

$$RWIN = \left\{ (ABW + \sum_{i=0}^I BW(f_i)) \text{Avg}(RTT) \right\} / n \quad \dots(3)$$

20

式 (3) で B W (f_i) はエレファントフロー i が使用するバンド幅 (スループット) であり、A v g (R T T) は、複数の受信側エンドノード 3 から受信した平均 R T Tの平均値である。

【 0 1 2 6 】

R W I N設定部 7 9 は、R W I N計算部 7 8 により計算された R W I Nを競合が起きている受信側エンドノード 3 に通信部 8 0 を介して送信する。通信部 8 0 は、受信側エンドノード 3 から送信される平均 R T Tの受信、受信側エンドノード 3 への R W I Nの送信、スイッチ装置 4 からのパケットの受信等の通信を行う。

【 0 1 2 7 】

30

次に、コントローラ 1 0 による処理のフローについて説明する。図 3 1 は、コントローラ 1 0 による処理のフローを示すフローチャートである。図 3 1 に示すように、通信部 8 0 は、パケットを受信し、受信したパケットが S Y Nパケットであるか否かを判定する (ステップ S 6 1)。

【 0 1 2 8 】

その結果、S Y Nパケットでない場合には、通信部 8 0 は、パケットが F I Nパケット又は R S Tパケットであるか否かを判定する (ステップ S 6 2)。その結果、F I Nパケットでも R S Tパケットでもない場合には、通信部 8 0 は、パケットが平均 R T Tの B P D Uかを判定する (ステップ S 6 3)。その結果、パケットが平均 R T Tの B P D Uでない場合には、コントローラ 1 0 は、ステップ S 6 1 に戻る。

40

【 0 1 2 9 】

一方、受信したパケットが平均 R T Tの B P D Uである場合には、R T T測定部 7 5 は平均 R T Tを C T R L _ R W I Nテーブル 7 2 に記録し (ステップ S 6 4)、R T Tフラグを t r u e に設定する (ステップ S 6 5)。また、ステップ S 6 2 において、受信したパケットが F I Nパケット又は R S Tパケットである場合には、コントローラ 1 0 は、ステップ S 7 0 に進む。

【 0 1 3 0 】

また、ステップ S 6 1 において、受信したパケットが S Y Nパケットである場合には、5 タプル検出部 7 1 は、S Y Nパケットから 5 タプルを取り出す (ステップ S 6 6)。そして、パス計算部 7 3 は、トポロジーを読み込み (ステップ S 6 7)、エレファントフロ

50

ーのパスを計算して(ステップS68)、CTRL__RWINテーブル72にパスを保存する(ステップS69)。

【0131】

そして、パス検索部74は、CTRL__RWINテーブル72を検索し、同じパス上でエレファントフローの数が2より大きいかなかを判定し(ステップS70)、2より大きくない場合には、コントローラ10は、ステップS61に戻る。一方、2より大きい場合には、コントローラ10は、RTTフラグがtrueであるかなかを判定し(ステップS71)、trueでない場合には、受信側エンドノード3から平均RTTを受信していないので、ステップS61に戻る。

【0132】

一方、RTTフラグがtrueである場合には、通信部80は、パス上のスイッチ装置4からリンク数を収集する(ステップS72)。そして、ABW計算部77は、各リンクの利用可能バンド幅を計算し(ステップS73)、利用可能バンド幅の最小値ABWを計算する(ステップS74)。

【0133】

そして、通信部80は、最小値ABWを有するスイッチ装置4からエレファントフロー数を収集する(ステップS75)。そして、スループット計算部76は、Avg(RTT)とエレファントフロー数でスループットを計算する(ステップS76)。そして、RWIN計算部78は、ABWとスループットでRWINを計算する(ステップS77)。そして、RWIN設定部79は、エレファントフロー間競合が生じた受信側エンドノード3に通信部80を介してRWINを設定する(ステップS78)。

【0134】

このように、コントローラ10がRWINを計算し、計算したRWINを受信側エンドノード3に送信することによって、エレファントフローが競合する受信側エンドノード3間でスループットを均等化することができる。

【0135】

次に、レート制御のシーケンスについて説明する。図32は、レート制御のシーケンスを示す図である。図32に示すように、まず、前提条件及び事前準備が整う(ステップt31)。すなわち、インキャストが発生する。

【0136】

そして、 S_0 は R_0 へエレファントフローE#0のSYNパケットを送信し、 S_n は R_n へエレファントフローE#1のSYNパケットを送信する。そして、コントローラ10の5タプル検出部71はSYNパケットから5タプルを取り出し、パス計算部73はパスを計算する(ステップt32)。また、 R_0 は S_0 へエレファントフローE#0のACKを送信し、 R_n は S_n へエレファントフローE#1のACKを送信する。

【0137】

そして、エレファントフローの輻輳ウィンドウが輻輳回避状態になる(ステップt33)。そして、受信側エンドノード3のRTT測定部54がRTTを測定し、AvgRTT計算部56が平均RTTを計算する(ステップt34)。そして、 R_0 のコントローラ通信部68が平均RTTをBPDU化して送信し(ステップt35)、 R_n のコントローラ通信部68が平均RTTをBPDU化して送信する(ステップt36)。

【0138】

そして、コントローラ10の通信部80はリンク数を収集して、ABW計算部77は最小値ABWを計算する(ステップt37)。また、コントローラ10の通信部80はエレファントフロー数を収集して、スループット計算部76はスループットを計算する(ステップt38)。そして、コントローラ10のRWIN計算部78はRWINを計算し、通信部80はBPDU化したRWINを各受信側エンドノード3に送信する(ステップt39)。

【0139】

そして、コントローラ10からのRWINを各受信側エンドノード3のRWIN設定部

10

20

30

40

50

66がACKに設定する(ステップt40)。そして、 R_0 は S_0 へエレファントフローE#0のACKを送信し、 R_n は S_n へエレファントフローE#1のACKを送信する。

【0140】

上述してきたように、実施例2では、スループット計算部76が各エレファントフローのスループットを計算し、ABW計算部77が利用可能バンド幅の最小値ABWを計算する。そして、RWIN計算部78が、スループットから計算される各エレファントフローの利用可能バンド幅とABWとAvg(RTT)を用いてRWINを計算し、RWIN設定部79がRWINを各受信側エンドノード3に送信する。そして、コントローラ10からのRWINを各受信側エンドノード3のRWIN設定部66がACKに設定する。したがって、コントローラRWIN制御部5aは受信側エンドノード3間でエレファントフロー間の競合の発生を防ぎ、情報処理システム1aはエレファントフローを効率的に処理することができる。

10

【0141】

また、実施例1及び2では、RWIN制御部5及びコントローラRWIN制御部5aについて説明した。しかし、RWIN制御部5及びコントローラRWIN制御部5aが有する構成をソフトウェアによって実現することで、同様の機能を有するRWIN制御プログラムを得ることができる。そこで、RWIN制御プログラムを実行するコンピュータについて説明する。

【0142】

図33は、実施例1及び2に係るRWIN制御プログラムを実行するコンピュータの構成を示す図である。図33に示すように、コンピュータ90は、メインメモリ91と、CPU(Central Processing Unit)92と、LAN(Local Area Network)インタフェース93と、HDD(Hard Disk Drive)94とを有する。また、コンピュータ90は、スーパーIO(Input Output)95と、DVI(Digital Visual Interface)96と、ODD(Optical Disk Drive)97とを有する。

20

【0143】

メインメモリ91は、プログラムやプログラムの実行途中結果などを記憶するメモリである。CPU92は、メインメモリ91からプログラムを読み出して実行する中央処理装置である。CPU92は、メモリコントローラを有するチップセットを含む。

【0144】

LANインタフェース93は、コンピュータ90をLAN経由で他のコンピュータに接続するためのインタフェースである。HDD94は、プログラムやデータを格納するディスク装置であり、スーパーIO95は、マウスやキーボードなどの入力装置を接続するためのインタフェースである。DVI96は、液晶表示装置を接続するインタフェースであり、ODD97は、DVDの読み書きを行う装置である。

30

【0145】

LANインタフェース93は、PCIエクスプレス(PCIe)によりCPU92に接続され、HDD94及びODD97は、SATA(Serial Advanced Technology Attachment)によりCPU92に接続される。スーパーIO95は、LPC(Low Pin Count)によりCPU92に接続される。

40

【0146】

そして、コンピュータ90において実行されるRWIN制御部プログラムは、DVDに記憶され、ODD97によってDVDから読み出されてコンピュータ90にインストールされる。あるいは、RWIN制御部プログラムは、LANインタフェース93を介して接続された他のコンピュータシステムのデータベースなどに記憶され、これらのデータベースから読み出されてコンピュータ90にインストールされる。そして、インストールされたRWIN制御部プログラムは、HDD94に記憶され、メインメモリ91に読み出されてCPU92によって実行される。

【0147】

なお、実施例1及び2では、インキャストが発生している場合について説明したが、本

50

発明はこれに限定されるものではなく、インキャストが発生していない場合のエレファントフロー間の競合にも同様に適用することができる。

【0148】

また、実施例1及び2では、エレファントフロー間で競合が発生している場合について説明したが、本発明はこれに限定されるものではなく、例えばミスフローなど他の種類のフロー間で競合が発生している場合にも同様に適用することができる。

【0149】

また、実施例1及び2では、ネットワークがスパインスイッチとリーフスイッチで構成され、リーフスイッチの下にエンドノードが接続される場合について説明した。しかしながら、本発明はこれに限定されるものではなく、例えば、リーフスイッチの下にTOR (Top Of Rack) スイッチが接続し、TORにエンドノードが接続される場合など他の構成のネットワークにも同様に適用することができる。

10

【符号の説明】

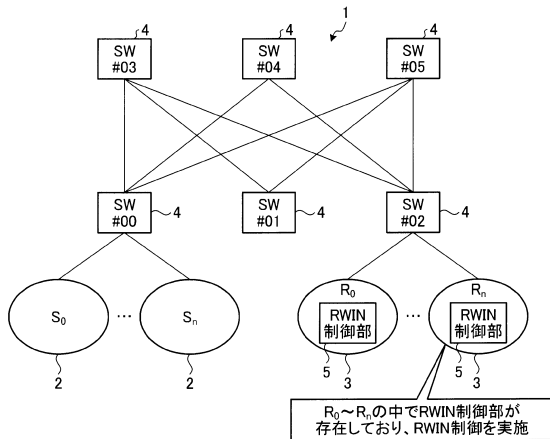
【0150】

1	, 1 a	情報処理システム	
2		送信側エンドノード	
3		受信側エンドノード	
4		スイッチ装置	
5		RWIN制御部	
5 a		コントローラRWIN制御部	20
6		APP	
7		NIC	
8		OS	
9		仮想スイッチ	
10		コントローラ	
11		ネットワーク	
51		ウィンドウテーブル	
52		ウィンドウ監視部	
53		RTTテーブル	
54		RTT測定部	30
55		RTTタイマ	
56		AvgRTT計算部	
57		NTタイマ	
58		スループット測定部	
59		スループット計算部	
60		最小スループット計算部	
61		RWINテーブル	
62		RWIN計算部	
63		RWIN比較部	
64		RWINタイマ	40
65		setRWIN計算部	
66		RWIN設定部	
67		タイムスロット設定部	
68		コントローラ通信部	
71		5タブル検出部	
72		CTRL_RWINテーブル	
73		パス計算部	
74		パス検索部	
75		RTT測定部	
76		スループット計算部	50

- 7 7 A B W 計算部
- 7 8 R W I N 計算部
- 7 9 R W I N 設定部
- 8 0 通信部
- 8 9 受信バッファ
- 9 0 コンピュータ
- 9 1 メインメモリ
- 9 2 C P U
- 9 3 L A N インタフェース
- 9 4 H D D
- 9 5 スーパー I O
- 9 6 D V I
- 9 7 O D D

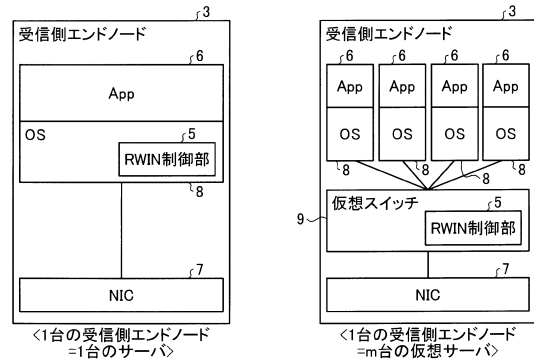
【 図 1 】

実施例1に係る情報処理システムの構成を示す図



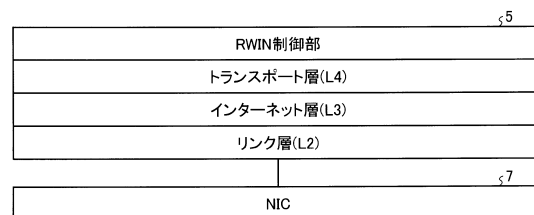
【 図 2 】

RWIN制御部の実現方法を説明するための図



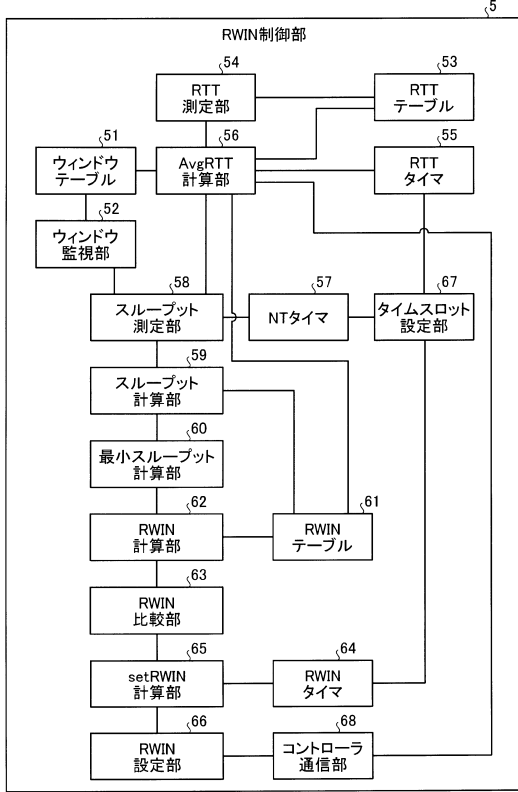
【 図 3 】

通信階層におけるRWIN制御部の位置づけを示す図



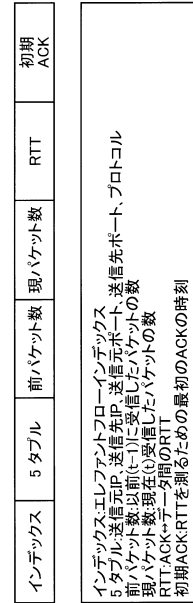
【図4】

RWIN制御部の構成を示す図



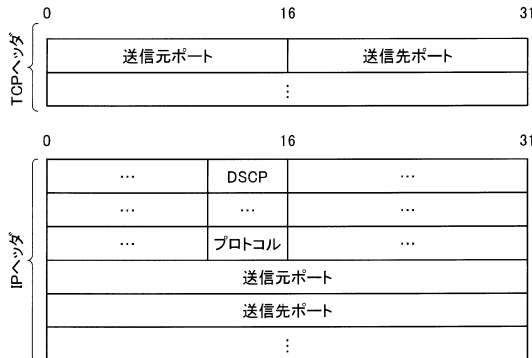
【図5】

ウィンドウテーブルがエレファントフロー毎に記憶する項目の一例を示す図



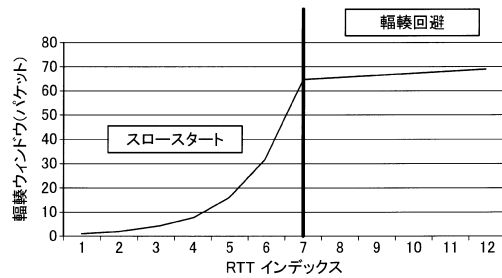
【図6】

5タブル及びDSCPの格納場所を示す図



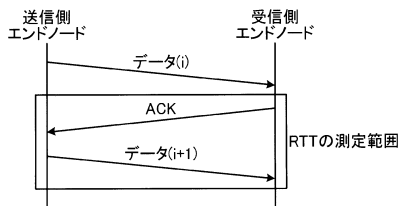
【図8】

輻輳回避状態を説明するための図



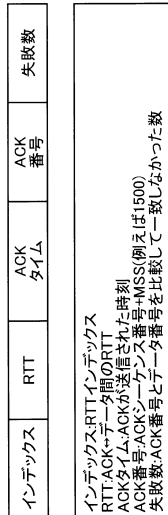
【図7】

RTTの測定範囲を示す図



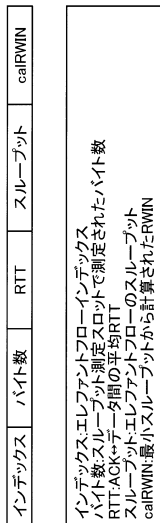
【 図 9 】

RTTテーブルがRTT毎に記憶する項目の一例を示す図



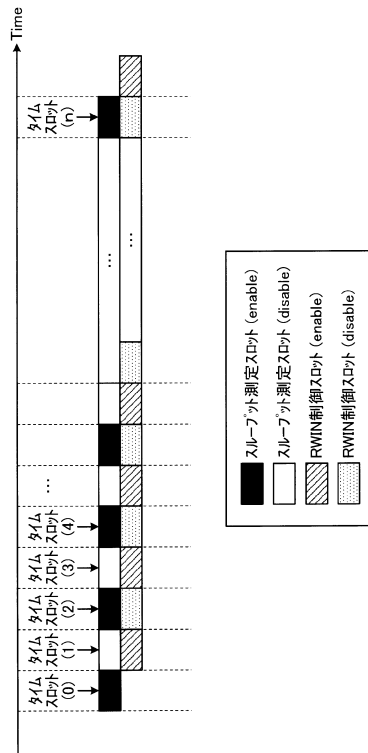
【 図 11 】

RWINテーブルがエレファントフロー毎に記憶する項目の一例を示す図



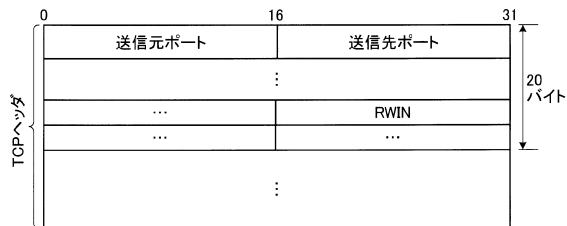
【 図 10 】

スルーバット測定スロット及びRWIN制御スロットを説明するための図

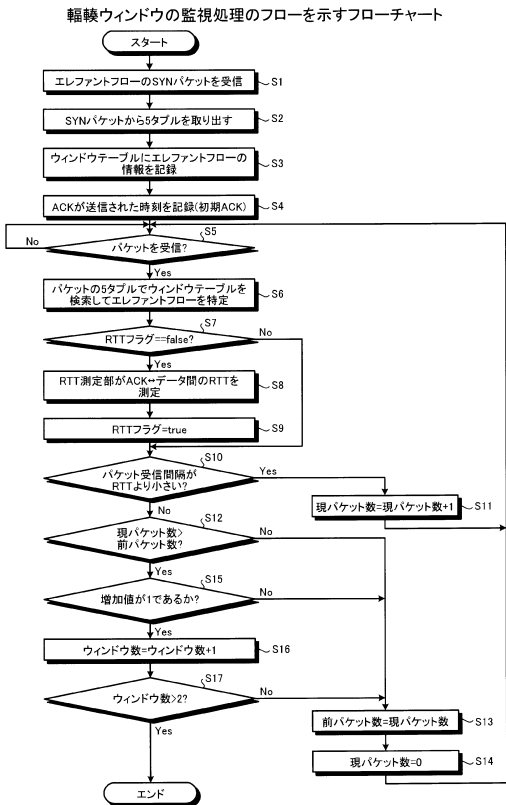


【 図 12 】

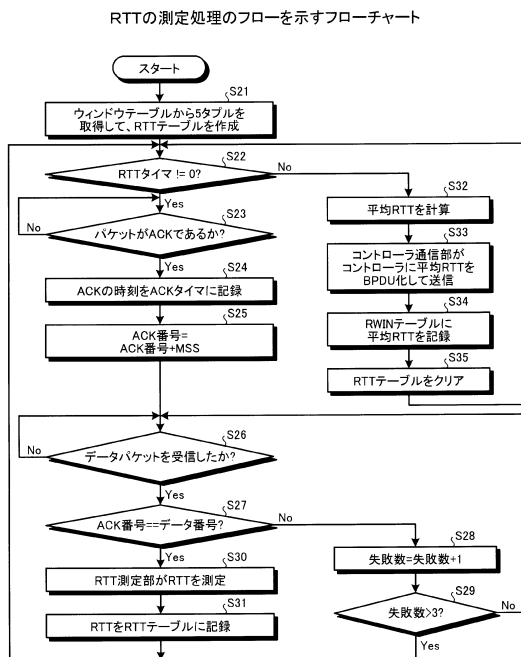
setRWINの設定場所を示す図



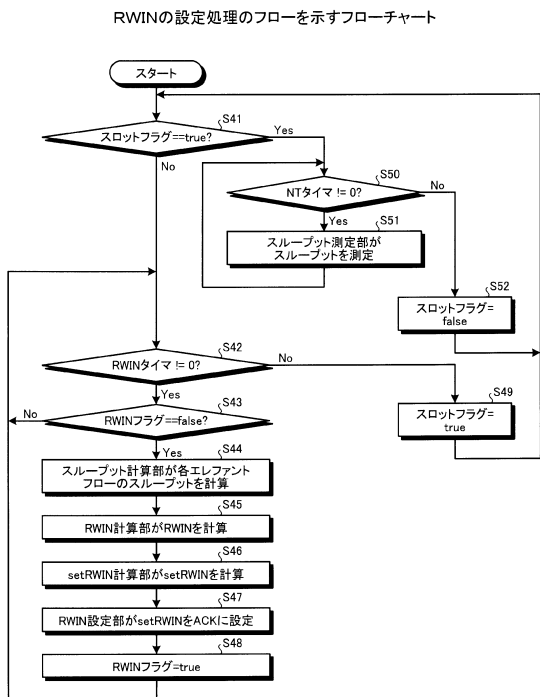
【図13】



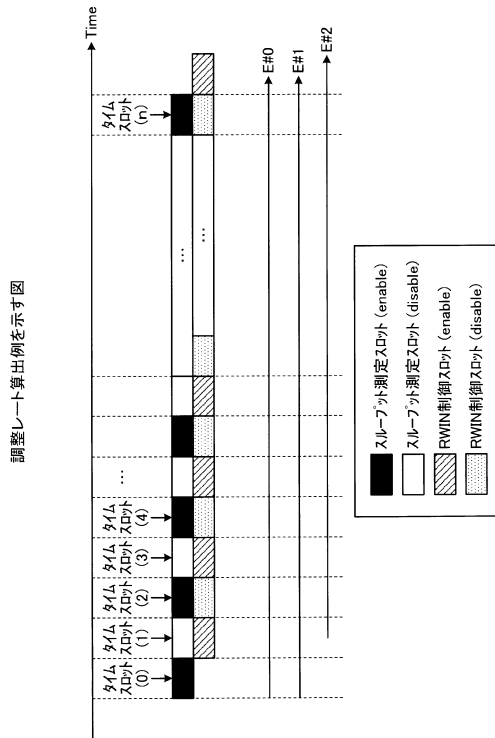
【図14】



【図15】

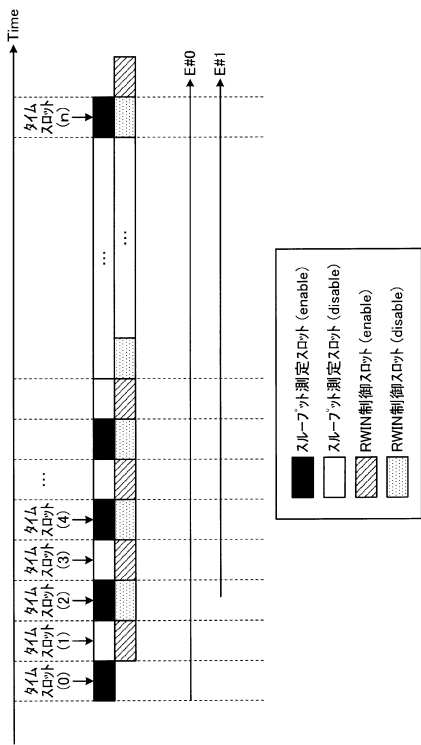


【図16】



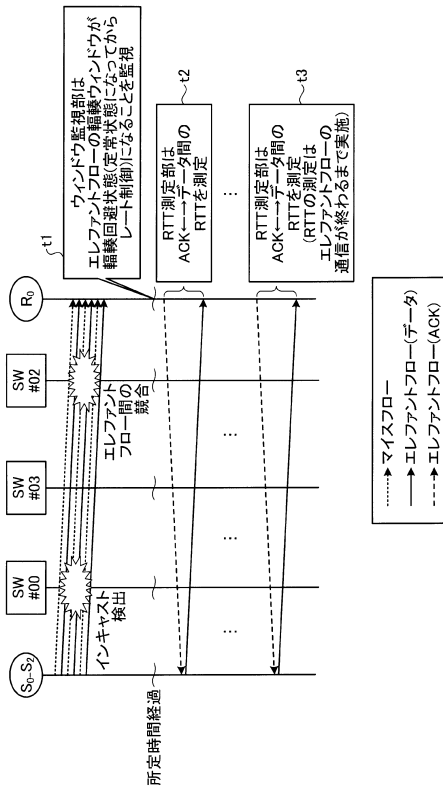
【図 17】

スループット測定スロットの途中にフローが出現した場合の調整スロット抽出例を示す図



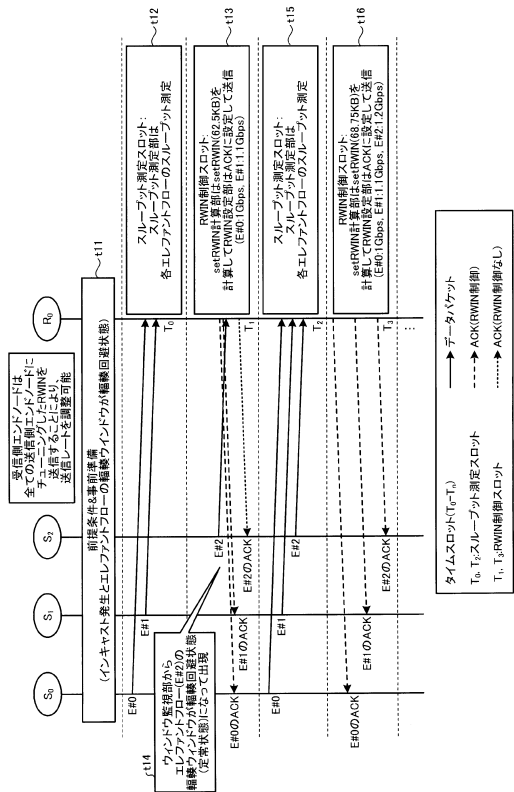
【図 18】

前提条件及び事前準備を説明するための図



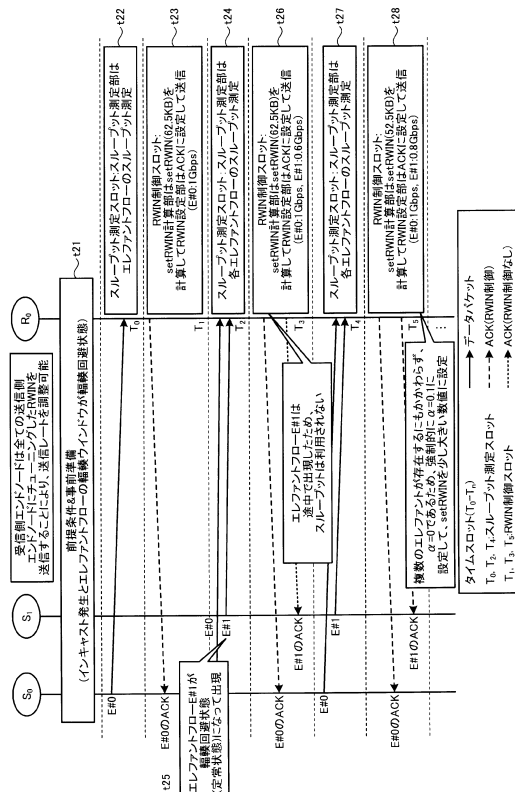
【図 19】

レート制御のシーケンスを示す第1の図



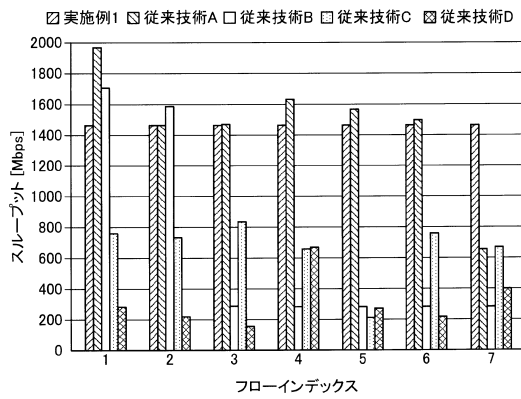
【図 20】

レート制御のシーケンスを示す第2の図



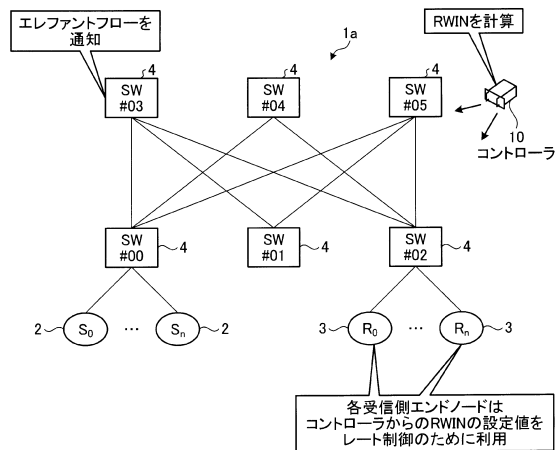
【図 2 1】

実施例1に係るレート制御の効果を示す図



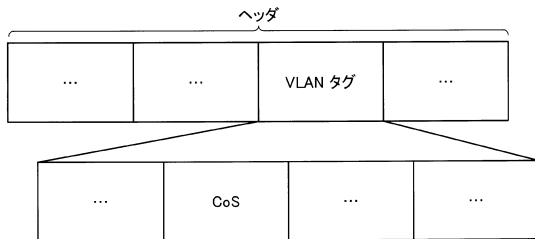
【図 2 2】

実施例2に係る情報処理システムの構成を示す図



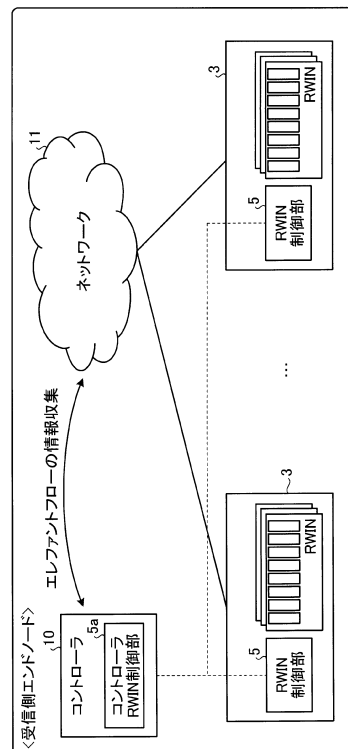
【図 2 3】

CoSを説明するための図



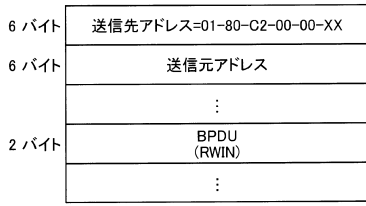
【図 2 4】

実施例2に係るレート制御を説明するための図



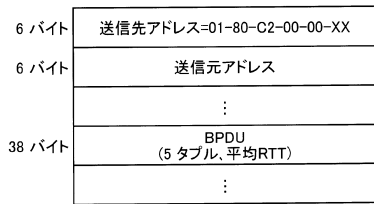
【図 25】

RWIN送信用のPDUのフォーマットを示す図



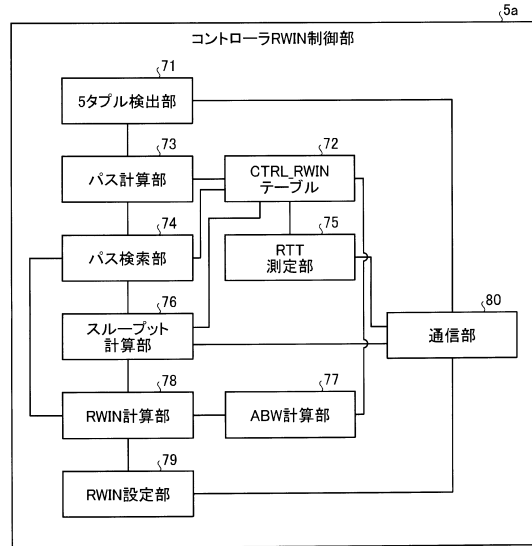
【図 26】

平均RTT送信用のPDUのフォーマットを示す図



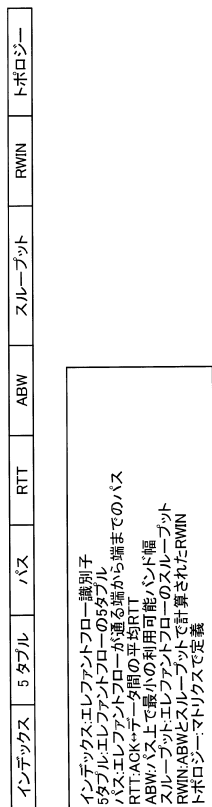
【図 27】

コントローラRWIN制御部の構成を示す図



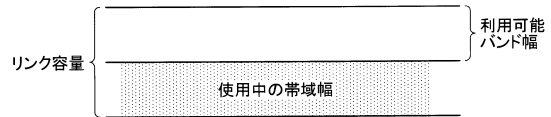
【図 28】

CTRL_RWINテーブルがエレファントフロー毎に記憶する項目の一例を示す図



【図 29】

利用可能バンド幅を説明するための図



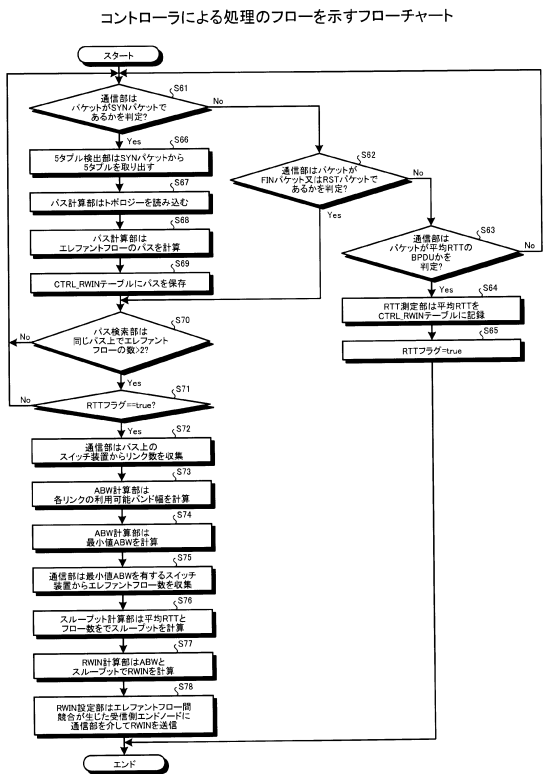
【図 30】

マトリクスの一例を示す図

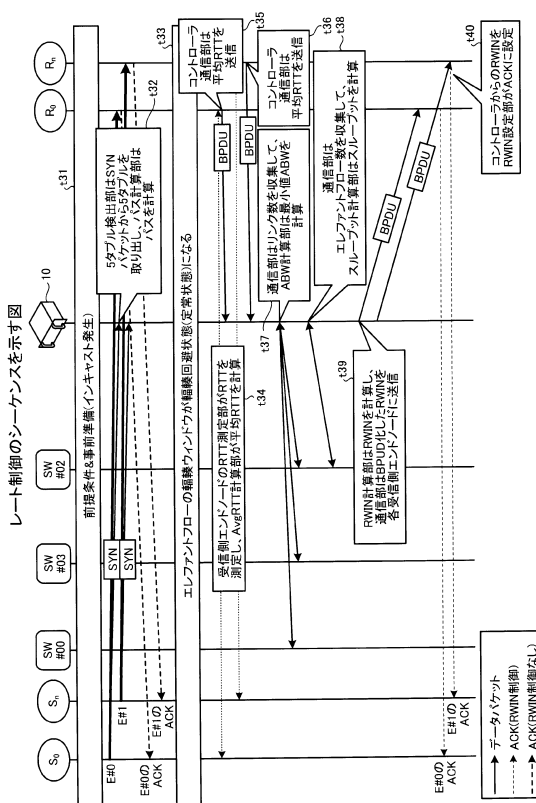
トポロジー	SW#00	...	SW#yy
SW#00	-	...	○
SW#01	○	...	○
⋮	⋮		⋮
SW#xx	○	...	-

-:スイッチ装置の間でリンクが存在しない
 ○:スイッチ装置の間でリンクが存在する

【図 3 1】

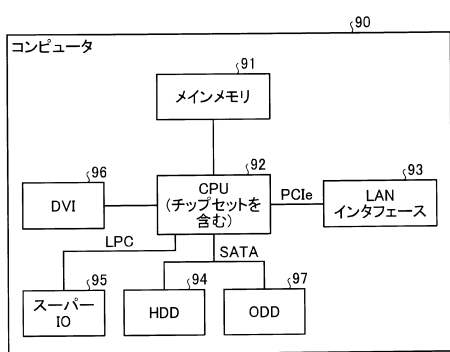


【図 3 2】

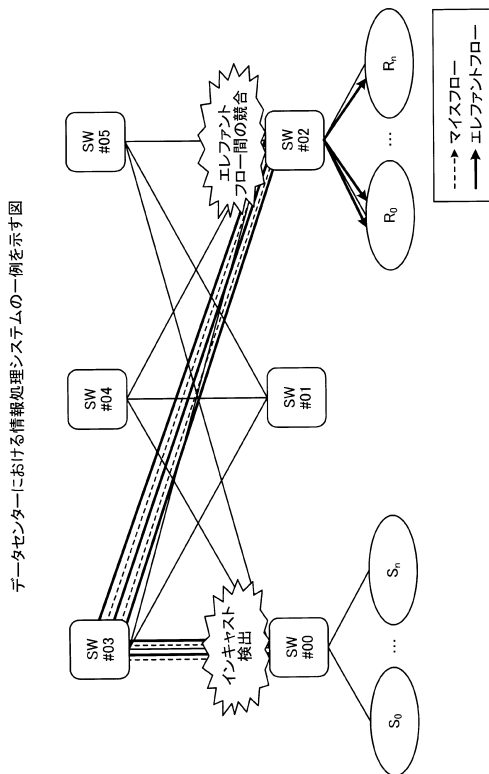


【図 3 3】

実施例1及び2に係るRWIN制御プログラムを実行するコンピュータの構成を示す図

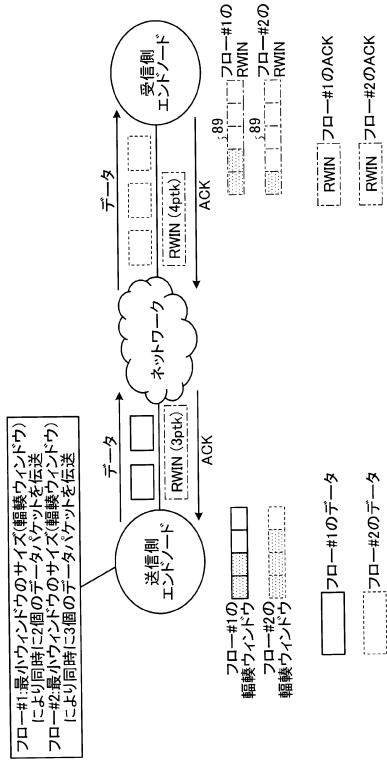


【図 3 4】



【 図 35 】

レート制御を説明するための図



フロントページの続き

(56)参考文献 米国特許出願公開第2007/0076726(US, A1)

特開2005-218069(JP, A)

特開2011-035608(JP, A)

特開2009-206733(JP, A)

吉水 卓 Suguru YOSHIMIZU, 高速TCPの公平性を改善する輻輳制御アルゴリズムの検討 Congestion Control Algorithms to Improve Fairness of High-Speed TCPs, 電子情報通信学会 技術研究報告 Vol. 109 No. 327 IEICE Technical Report, 日本, 社団法人電子情報通信学会 The Institute of Electronics, Information and Communication Engineers, 2009年12月 3日, 第109巻, p.25-30

(58)調査した分野(Int.Cl., DB名)

H04L 12/807

H04L 29/08