論文

拡散現象を指導原理とする自律分散フロー制御の実装に向けた アクティブフロー数計測技術の検討

高野 知佐^{†a)} 山内 正志^{††} 会田 雅樹^{†††}

On Implementation Issues of Diffusion-Type Flow Control for Multiple Flows

Chisa TAKANO^{†a)}, Tadashi YAMAUCHI^{††}, and Masaki AIDA^{†††}

あらまし 筆者らはこれまで,高速ネットワークの制御に求められる短い制御遅延の要求を満足するために, 自律分散で動作する拡散型フロー制御技術を提案してきた.この技術は,局所的なネットワーク情報のみに基づ くノードの自律動作によって,間接的にネットワーク全体を適切な状態に導く制御技術である.本論文では,自 律分散システムとしてのモジュール構成を考慮しながらネットワークシミュレータ ns2 への実装を行い,動作特 性の評価・分析を行った.この結果,複数フローが混在する環境での適切な制御動作を実現するためには,アク ティブなフロー数に関する情報を局所的に推定する技術が必要であることが分かった.この課題を解決するため, 拡散型フロー制御技術の改良方式を提案し,その効果を確認した.また,この改良により,複数フロー間のパ ケットスケジューリングを簡略化することが可能となり,より簡易な実装技術に結び付けることができた. キーワード トラヒック,フロー制御,拡散方程式,自律分散制御

1. まえがき

近年のインターネットの急速な普及・需要の拡大に より,将来をにらんだ高速なネットワークの構築が進 められている.超高速ネットワークでは,光速度で決 まる固定伝搬遅延に比べてノードの動作速度が相対 的に高くなり,ノードの高速動作を特徴づける時間ス ケールにおいて,各ノードが収集可能な情報は非常に 限られたものになる.遠くのノードから収集した情報 は過去の情報であり,ネットワーク全体に関する現時 点の状態を知ることはできないからである.この意味 で,ネットワークの高速化は各ノードが情報的に孤立 する状況を招く.また,高速ネットワークでは制御遅 延の影響が大きいため,非常に高速で動作する制御機 構が必要となる.このように,超高速ネットワークで 適切に動作する制御を実現するためには以下の課題を 解決する必要がある.

局所的な状態情報のみで動作する制御機構

超高速ネットワークでは,リンク伝搬遅延(光速度 により決まる)によりネットワークの広範囲な状態情 報をリアルタイムに把握することができない.このた め,リアルタイムに知ることができる局所的な状態情 報のみで動作可能な制御が必要になる.

超高速で動作する制御機構

超高速ネットワークでは,ネットワーク状態が短時 間で大きく変化するため,非常に短い制御遅延で動作 するネットワーク制御技術が必要になる.

フロー制御の課題に対する多くの研究は,線形計画 問題として定式化される[1]~[3].これらの定式化で は,ネットワークの状態情報を収集が可能であること を前提にしている.しかし,高速ネットワークでは広 範な状態情報をリアルタイムに収集することができな い.また,線形計画問題の求解には十分な時間が必要 なので,高速ネットワークに要求される非常に短い制 御遅延を満足することが難しい.このように,ネット ワーク全体の情報を1箇所に収集することを前提と した集中制御の枠組みは,超高速ネットワークでのフ ロー制御として適切でない.このため,自律分散制御

 [†]広島市立大学情報科学研究科,広島市 Graduate School of Information Sciences, Hiroshima City University, Hiroshima-shi, 731-3194 Japan
 ^{††}東京都立科学技術大学工学部,日野市 Faculty of Engineering, Tokyo Metropolitan Institute of Technology, Hino-shi, 191-0065 Japan

^{†††} 首都大学東京大学院システムデザイン研究科,日野市 Graduate School of System Design, Tokyo Metropolitan University, Hino-shi, 191-0065 Japan

a) E-mail: takano@hiroshima-cu.ac.jp

に基づくフロー制御機構が必要とされる.

現在,自律分散型のフロー制御として,例えば TCP のようなエンドホストによる制御が広く用いられ,そ の特性について活発に研究されている[4]~[6].エンド ホストによるフロー制御では、タイムスケールはラウ ンドトリップタイム (RTT) 程度となる. 低速なネッ トワークでは, RTT 程度の制御遅延による影響を無 視することが可能だが,ネットワークが高速化すると RTT 程度の制御遅延が大きな影響を与えるようにな る可能性がある.ネットワークの高速化に対して RTT 自身の物理的な値は不変であっても、ノードの動作速 度の時間スケールを時間の基準として見た RTT は大 きくなるからである.これは,RTT 程度の時間内に ノードが処理するパケット数が増大することを意味し, より制御遅延の影響を強く受けるようになる.RTT よりも短い時間スケールで俊敏に動作する制御を実現 するためには、ノードが局所情報に基づいて自律動作 するタイプの制御機構が必要になる.

我々は,高速ネットワーク環境でも適切に動作す るフロー制御方式として,拡散型フロー制御技術 (diffusion-type flow control:DFC)を提案してきた. このフロー制御方式は,各ノードが自分自身で利用可 能な局所的情報のみに基づいて自律動作を行い,そ の結果としてネットワーク全体を安定させるものであ る.図1はネットワーク全体を安定させるものであ る.図1はネットワークの各種制御を,要求される 制御遅延時間の大きさにより階層化したものである. routing 制御,session initiation protocol(SIP)な どのシグナリング処理,TCPやexplicit congestion notification(ECN)などの流量制御は,それぞれ数十 分,秒,RTT程度の制御遅延をもつ.DFCは,RTT よりも更に厳しい制御遅延の要求を満たすための制御 技術として位置づけられる.DFC が対象とする制御



図 1 各種ネットワーク制御の要求される制御遅延による 階層化

遅延の時間スケールは,ネットワークが高速化するの に伴い重要性が増す領域になっている.

これまで,シミュレーションによって DFC による ネットワーク状態の安定性,状態変化に対する適応性, TCP との親和性を評価し,DFC がねらいどおりの性 能を発揮することを確かめてきた[7]~[10].これらの 結果に加え,DFC が実ネットワークでも正しく動作 することを確かめるためには,DFC の実装に向けた 技術課題として以下の課題を確認する必要がある.

DFC が自律分散制御としての装置構成によって実装が可能であること.

• 自律分散制御として実装した DFC がねらいど おりの性能を発揮すること.

これまでの評価では, TCP と DFC との共存可能 性や相互補完性を調べるために,ネットワークシミュ レータである ns2 [13] を利用し, DFC 機能を ns2 に 実装することで評価を行ってきた.ns2 は本来,装置 のモジュール構成を反映したシミュレータとなってい るため,自律分散制御としての装置構成を意識しなが ら DFC 機能を ns2 に実装することによって,DFC 機 能の実装可能性を確認することができるはずである. しかし,これまでの評価では DFC の構成を完全に自 律分散制御の装置構成で実装したわけではないため, DFC の実装可能性が確認できていなかった.

ns2 に組み込んだ DFC 機能のこれまでの実装では, シミュレーションモデルとして「各ノードは自分を通 過するフロー数を知っている」という前提のもとで評 価を行い, DFC が適切に動作することを確認してき た.自律分散的な装置構成では,すべてのフロー情報 を知っている管理装置の存在を仮定していないので, 各ノードは自分自身で局所的に通過フロー数の情報 を取得しなければならない.もしネットワーク内のフ ローが常にトラヒックを発生させていれば,到着パケッ トを観測することでノードを通過するフロー数を知る ことができる.しかし,一般にはフローが一時的にパ ケットを発生させない状態になることがあり得るし, 発生させていてもパケット間隔が広く、適当な時間区 間ではパケットの発生がない場合も考えられる.した がって,自律分散的な装置構成をとる限り,各ノード は自分を通過する正確なフロー数を知ることができ ない.

本論文では,モジュール構成を考慮した ns2 の特性 を利用し,自律分散で動作するモジュールとして DFC の実装を行う.更に,実装した DFC の振舞いとこれ

Fig. 1 Hierarchical structure of network controls with respect to their required control delay.

までのシミュレーション評価結果との違いから,通過 フロー数の情報取得に関連する問題点を明らかにし, その解決策を示す.更に改良した DFC がねらいどお りの特性を示すことを明らかにすることで,DFCの 実装可能性を確認する.

本論文の構成は以下のとおりである.2.で,これ まで評価してきた DFC の概要を示し,3.で,既存の DFC を ns2 に実装する方法を示す.また 4. で既存の DFC を ns2 に実装した場合の問題点を明らかにした 後,5. で DFC の実装に向けた改良方法を示す.更に, 6. で,改良した DFC を評価し,その効果を実証する. 最後に 7. でまとめる.

2. 拡散型フロー制御の概要

2.1 制御の基本コンセプト

インターネット型のネットワークにおいても,エン ドツーエンドでの QoS 保証を要求するフローについ ては,RSVP のように特定の経路を用いた通信を行 う.本論文では,このような特定の経路をもつフロー を対象とし,ノードではフローごとのパケットキュー イングが行われることを仮定する.

DFC は,各ノードが自分自身の知り得る局所情報 のみに基づいて自律的に局所フローを制御し,間接的 にネットワーク全体の状態を望ましい方向に誘導する フロー制御技術の枠組みである.具体的には,ネット ワークがふくそうに陥ることをを回避したり,ふくそ う状態から回復するように動作する.各ノードが自律 的に動作することで,状況の変化に即応した高速な制 御動作が可能となる.

DFC は, ノードの局所的自律動作によってネット ワーク全体の状態を望ましい方向に誘導するために, 物理の近接作用の考え方に基づいた制御の仕組みを採 用している.これを熱拡散現象を例にして説明する. 図 2 のように鉄棒の一点を熱すると,温度分布は正規 分布となり拡散する.このとき,鉄棒の一点は鉄棒全 体の温度分布に関する知識はなく,左右の温度差に比 例した熱量を高温側から低温側に流しているにすぎな



Fig. 2 Example of thermal diffusion phenomenon.

い.しかしながら,温度分布は拡散方程式の解として 秩序ある振舞いを示す.このように,システムの各部 分が局所情報のみに基づいて自律動作をしつつも,シ ステム全体の状態に秩序を生むことが可能である.こ の仕組みをネットワークに適用すると,各ノードの動 作ルールを適切に決めておけば,ノードの自律動作に より間接的にネットワークの全体性能を望ましい方向 に誘導することが可能であると考えられる.DFC で は各ノードが拡散現象に似たルールで局所的なパケッ トフローを制御することにより,一部のノードへのパ ケットの集中を回避し,パケットロスの起きにくい状 態を維持するように動作する.

2.2 拡散型フロー制御方式

あるフローに沿って一次元的に連結したノードi(i = 1, 2, ..., N)を考え(図3), ノードiとノードii+1の間のリンクの帯域を B_i とする.また, ノードiとノードi+1の間のリンクを共有するフロー数を M_i とする.ノードiは,時刻tにおける下流ノードi+1へのフローj($j = 1, 2, ..., M_i$)のパケット転送レー ト $J_i^j(t)$ を自分が知り得る情報のみに基づいて決定し, パケット転送を行う.また,上流ノードi-1に向けて ノードiの情報 $\mathbf{F}_i^j(t)$ をフィードバックする.パケッ ト転送レート $J_i^j(t)$ の決定は,下流ノードi+1から伝 搬遅延 d_i を伴って通知される情報 $\mathbf{F}_{i+1}^j(t-d_i)$ の到 着ごとに実行する.上流ノードi-1に向けたフィー ドバック情報 $\mathbf{F}_i^j(t)$ の通知は,ノードi-1とi間の 伝搬遅延 d_{i-1} に比例した時間間隔で行う.

以降では、パケット転送レート $J_i^j(t)$ とフィードバッ ク情報 $\mathbf{F}_i^j(t)$ の算出方法について述べながら、DFC の全体像を説明するが、それに先立ち、DFC の説明 に利用されるアクティブフローの概念について明らか にしておく、ノード i とノード i+1 の間のリンクを 共有する、または共有する可能性のあるすべてのフ ロー数を M_i としたとき、ある時点で必ずしもすべて のフローにパケットが流れているとは限らない、本論 文では、何らかの計測によってパケットが流れている



Fig. 3 Node interactions in our flow control model.

と判断されたフローをアクティブフローと呼ぶことに する^(注1).特定の時間帯においてあるノードがパケッ トの流れを観測したときに,実際にパケットが観測さ れるかどうかについては,音声の無音区間やパケット レートの高低などのトラヒックの特性に影響を受ける. このため,どのフローがアクティブと判断されるかは, 時刻ごとやノードごとに一般に異なる結果となること があり得る.

時刻 t , ノード i におけるフロー j の転送レート $J_{i}^{j}(t)$ は

$$J_{i}^{j}(t) = \max(0, \min(L_{i}^{j}(t), \tilde{J}_{i}^{j}(t)))$$
(1)
$$\tilde{J}_{i}^{j}(t) = r_{i}^{j}(t - d_{i}) - D_{i} \left(n_{i+1}^{j}(t - d_{i}) - n_{i}^{j}(t)\right)$$
(2)

によって決定する.ここで, $n_i^j(t)$ は時刻 t でノー ドiに存在するフローjのパケット数, $r_i^j(t-d_i)$, $n_{i+1}^{j}(t-d_{i})$ は下流ノードi+1からのフィードバッ クにより,(リンク伝搬遅延 d_i を伴って)通知される 情報で,それぞれ下流からの指示レート,下流ノー ド内のフロー j のパケット数を表す L_i^j はノード iからノードi+1に向かうリンクにおけるフローjの利用可能帯域である^(注2). また D_i は D(>0) を 定数として $D_i = D/d_i$ である. D は拡散係数と呼 ばれ, 0 < D < 1/2 の範囲の値をとる[11], [12].フ ロー j に関するパケット転送レート $J_i^j(t)$ は, すべて の $j = 1, 2, ... M_i$ について計算するのではなく,ア クティブフローについてのみ行う.ここで得られたパ ケット転送レート $J_{i}^{j}(t)$ により, アクティブフローの パケット転送レートを更新する.アクティブフローで ないフローのパケット転送レートは0とする.

一方,ノードiが生成するフローjに関するフィードバック情報 $\mathbf{F}_{i}^{j}(t)$ は

$$\mathbf{F}_i^j(t) = (r_{i-1}^j(t), n_i^j(t)) \tag{3}$$

からなり, これを上流ノード *i* - 1 に通知する.こ こで,

$$r_{i-1}^{j}(t) = \max(0, \min(\tilde{J}_{i}^{j}(t), B_{i}/M_{i}^{\text{act}}(t))) \quad (4)$$

$$M_{i}$$

$$M_i^{\text{act}}(t) = \sum_{j=1}^{m_i} \mathbb{1}_{\{j = \text{active}\}}$$
(5)

とする.ただし, $1_{\{j=\text{active}\}}$ はフロー jがアクティブ フローのとき1,そうでなければ0となる指示関数で あり, $M_i^{\text{act}}(t)$ はその時刻tでのアクティブフロー数 を表す. つまり, あるフロー j の指定レート $r_{i-1}^{j}(t)$ の上限は, 帯域 B_i をアクティブフロー数で割った値になる.

次に,DFC の直観的意味を説明する.議論を簡略化 するため,ノード ID を連続化し, $i \rightarrow x$ で表す.ま た,連続極限をとってリンクの伝搬遅延を $d_i \rightarrow 0$ とす る.このとき, $\tilde{J}_i^j(t) \rightarrow \tilde{J}(x,t)$, $r_i^j(t-d_i) \rightarrow r(x,t)$, $n_i^j(t) \rightarrow n(x,t)$ とすれば,式 (2) は

$$\tilde{J}(x,t) = r(x,t) - \kappa \frac{\partial n(x,t)}{\partial x}$$
 (6)

の形で書ける.ここで $\kappa > 0$ は離散のときの D に対応する定数である.連続の式

$$\frac{\partial n(x,t)}{\partial t} = -\frac{\partial \tilde{J}(x,t)}{\partial x}$$

を用いると,パケット密度n(x,t)の時間発展は

$$\frac{\partial n(x,t)}{\partial t} = -\frac{\partial r(x,t)}{\partial x} + \kappa \frac{\partial^2 n(x,t)}{\partial x^2} \tag{7}$$

と書ける.これは拡散型の方程式である.したがって, DFC はノードが局所情報に基づいて自律動作をする ことで,ネットワーク全体としてノード内パケット数 を均一化する効果が期待できる.

2.3 ネットワーク入口でのパケットレート制限法 DFC がサポートされているネットワークを DFC がサポートされていないネットワークと接続するため に,ネットワーク境界部分での境界条件を考える必要 がある.

ネットワーク外のノードまたはエンドホストは拡散 型フロー制御方式をサポートしていないことを前提 とし,ネットワーク入口での入力パケットレートは, フィードバック情報 $r_0^j(t-d_0)$ で指定されたレートに シェーピング (バッファで遅延させることによるレー トの調整)がされるものとする.ノードi = 1が通知 する情報 $r_0^j(t)$ は,ネットワーク外のノードまたはエ ンドホストのバッファに存在するパケット数を0と仮 定し,式(2)によってレート $\tilde{J}_0(t)$ を計算する.つま り, $r_0^j(t)$ は式(4)の代わりに,

$$r_0^j(t) = J_1^j(t) - D_0 \, n_1^j(t) \tag{8}$$

を用いて決定する. $r_0^j(t)$ はノードi = 1で取得可能

 ⁽注1): アクティブフローの判定方法は後述する.
 (注2): L^j_{*}の決定方法は 2.4 で述べる.

な量のみによって計算することができる.これによ り,ネットワーク外からの入力パケットレートとネッ トワーク内部の DFC によるレートとの親和性を高め ることができる.

しかし, DFC は局所情報に基づくノードの自律動 作を用いた制御なので,ネットワーク内の状態を高速 に平滑化させる効果には優れるが,ネットワーク内に 外部から流入する余分なパケットを防ぐ効果は少な い.そこで,必要に応じてネットワーク入口でパケッ トレートを制限する機能を併用し,ネットワーク入口 でのパケット流入量を適切に制限することが必要と なってくる.これを実現するためには,ネットワーク 入口部分でフローに沿った経路上のネットワーク状態 を把握し,その状態に合わせてネットワーク入口での パケットシェーピングを施すことが有効である.つま り,ノード i が生成するフィードバック情報 $\mathbf{F}_{i}^{i}(t)$ に リンクの最大利用可能帯域に関する情報 $\ell_{i}^{i}(t)$ を加え,

$$\mathbf{F}_{i}^{j}(t) = (r_{i-1}^{j}(t), n_{i}^{j}(t), \ell_{i}^{j}(t))$$
(9)

を上流ノードi-1に通知する. $\ell_i^j(t)$ は,

$$\ell_i^j(t) = \min(\ell_{i+1}^j, B_i / M_i^{\text{act}}(t))$$
(10)

のように生成する . $\ell_i^j(t)$ の計算は , 下流のリンクのうちで一番小さい (とノード i が認識する) 帯域の値を用いる . この値で , 式 (8) の値を制限することにより , ネットワーク内に外部から流入する余分なパケットを防ぐことができる .

2.4 各フローの利用可能帯域の決定方法

リンク帯域は複数のフローで共有されているため, フロー間で適切に Lⁱ_i(t) を調節する必要がある.

ノードiからノードi+1へのリンクの帯域を B_i としたとき, $L_i^j(t)$ ($j=1, 2, \ldots, M_i$)は,

$$\sum_{j=1}^{M_i} 1_{\{j=\text{active}\}} \times L_i^j(t) = B_i$$

$$L_i^j(t) = 0 \quad (j \neq \text{active})$$
(12)

を満たすとする.これは,アクティブフローのみの利 用可能帯域で全帯域を分け合うことを意味する.

DFC では,転送レートの理想値は $\tilde{J}_{i}^{j}(t)$ であり, $J_{i}^{j}(t)$ は理想値 $\tilde{J}_{i}^{j}(t)$ を物理的に利用可能な帯域を上限として制限したものである.多くのフローがある場合,理想の転送レート $\tilde{J}_{i}^{j}(t)$ は, $\sum_{j=1}^{M_{i}}\tilde{J}_{i}^{j}(t) > B_{i}$ となり,すべてのフローがそれぞれの理想の転送レート $ilde{J}_{i}^{j}(t)$ を使用することができず,パケット密度のスムーズな平滑化を妨げる原因となり得る. $L_{i}^{j}(t)$ の最も単純な決定方法は,アクティブフローごとに重み $ilde{J}_{i}^{j}(t)$ を付けて帯域 B_{i} を分け合うことである.つまり,それぞれのフローの重み $W_{i}^{j}(t)$ を

$$W_{i}^{j}(t) = \frac{\tilde{J}_{i}^{j}(t)}{\sum_{k=1}^{M_{i}} 1_{\{k=\text{active}\}} \times \tilde{J}_{i}^{k}(t)}$$
(13)

とし,それぞれのフローの利用可能帯域を

$$L_i^j(t) = B_i \times W_i^j(t) \tag{14}$$

とすることで,より大きい重み $\tilde{J}_{i}^{j}(t)$ をもったフロー がより大きい転送レートを確保でき,下流ノードによ り多くのトラヒックを流すことができることになる. このとき,重みの小さなフローは相対的に他のフロー の転送レートはより小さく制限されることになる.こ れによりフロー間で帯域が分配され,共通経路上にあ るノード内パケット数のフロー間のばらつきが平滑化 される.また,フローごとに利用可能帯域が適切に調 整されるため,ボトルネックに対してフローの下流方 向にノード内パケット数を平滑化する効果も得られる.

3. ns2 に実装された DFC の概要

本章では, 2. で示した既存の DFC 機能を実装した ns2 の装置構成を示す.

3.1 各モジュールの概要

ns2 に実装された DFC 機能は主に四つのモジュー ルから構成されている(図4).以下にそれぞれの機能 の概要を示す.

• DFCFlowQ

受信したパケットをフローごとにキューイングする. また,あるフローのパケットが DFCQueue に存在しな ければ,そのフローに関する DFCF1owQ の先頭パケッ トを DFCDBAgent により定められた転送レート $\tilde{J}_i^j(t)$ で DFCQueue へ転送する (図 5 参照).

• DFCQueue

DFCF1owQ から受信したパケットをリンクへ送出す る.パケットを送出するフローは,各フローに定めら れた重みにより weighted round robin (WRR)で決 定される(図5参照).

• DFCDBAgent

下流ノードからのフィードバック情報が記された DFC パケットを受信し,受信した値と自ノードの情報をもとに自ノードのフィードバック情報を算出する.



図 4 DFC のモジュール構成 Fig. 4 Modules of DFC function.





• DFCAgent

DFCDBAgent により算出されたフィードバック情報 をパケット化し, DFCFlowQ を通して上流ノードへ送 出する.

3.2 DFCQueue からリンクへのパケット送出方法 DFCF1owQ から DFCQueue へのパケット転送レート は $\tilde{J}_{i}^{j}(t)$ である.これは単一のフローjに着目して算 出された値であり,DFCQueue からリンクへのパケッ ト送出時に,式(1)のようにレート $\tilde{J}_{i}^{j}(t)$ はフローjの 利用可能帯域 $L_{i}^{j}(t)$ によって $J_{i}^{j}(t)$ に制限される.これ を実現するために,これまでの実装ではDFCQueue か らリンクへのパケット送出方法を重み $J_{i}^{j}(t)$ のWRR 方式としている.

3.3 アクティブフロー数の測定方法

式 (4) と式 (10) で用いられている $M_i^{\text{act}}(t)$ の測定 方法について説明する.

DFCF1owQ から DFCQueue に送出が行われるごとに, パケットが到着したフロー数をアクティブフローとし て測定し,フィードバック情報生成に必要なアクティ ブフロー数 $M_i^{act}(t)$ を計算するための値を保持する. 具体的には,



activeNum

= (DFCFlowQ から DFCQueue への送出回数)

activeCount

= (パケット到着のあったフロー数の累計)

とし,図6のようにして,activeCount/activeNum によりアクティブフロー数 $M_i^{act}(t)$ を測定している. また,activeNum,activeCountはフィードバック情 報が生成されるごとにリセットされる.この方法では, フィードバック情報が生成される間の平均アクティブ フロー数を求めていることになる.

4. DFC 実装上の問題点と原因の分析

本章では,自律分散制御としての装置構成で既存の DFC の動作特性を評価し,問題点抽出と原因の分析 を行う.具体的には,既存のDFC が実装された ns2 を用いて,単純なモデルでシミュレーションを実行し, DFC の動作特性を分析する.また,本論文の評価で は,DFC のパラメータをD = 0.4,フィードバック 情報 $\mathbf{F}_{i}^{i}(t)$ の送信間隔を d_{i-1} とした.



4.1 シミュレーションモデルとシナリオ

シミュレーションに用いるモデルは,あるフローに 沿った一次元ネットワークを考える(図7参照).ノー ド数を30,各リンクの伝搬遅延を0.1 msでリンク長 は固定,各ノードのバッファ容量1800パケット,1パ ケットは固定長1500 Byteで,1リンク上の最大パケッ ト数は100パケットとした.ネットワーク内のノード はDFCの機能をもつ.また,パケットの発生源とし てソースノード1と2を配置した.これらはDFC方 式適用外(シェーピング機能はあり)のノードである. 各ソースノードのパケット生成レートは,フロー制御 がかかっていない状態で単位時間当り100パケットと した.

シミュレーションシナリオは次のとおりである.ネッ トワーク上にはフローj = 1とフローj = 2の二つ の TCP フローを流す. ソースノード 1 で発生するフ $\mathbf{D} - j = 1$ は性能評価対象のフローで, ノード1から ノード 30 へと流れる. ソースノード 2 で発生するフ $\mathbf{D} - i = 2$ はフローi = 1の妨げとなる背景フローで, 対象フローの経路上のノード 15 に流入し, ノード 30 へと流れる. TCP の window size は両フローともに RTT の bandwidth-delay product より十分大きい値 2 bld = 1 bld = 1 bld = 0.0 s $\mathbf{D} - 2$ はt = 0.06sにそれぞれ送信を始め, t = 0.12s でシミュレーション終了とする.フローi = 2が送信 され始めると、ノード15からノード16へのリンクは ボトルネックとなり,二つのフローは DFC のルール によって転送レートが規制されることになる.シミュ レーション結果より背景フローが流入直後からのネッ トワーク状態の時間変化を調べることで, DFC の従 来の実装について性能を評価する.

4.2 シミュレーション評価結果

フロー j = 2 が流入した直後からの,フロー j = 1のネットワーク内のノード上パケット数の振舞いを調 べた結果が図 8 であり,また各時間のノード上パケッ ト数の総和を示したものが図 9 である.図 8 の横軸







Fig. 9 Temporal evolution of the total number of stored packets in nodes.

はフローの経路に沿ったノード ID を示しており,縦 軸はそのノード内パケット数を表している.グラフ ごとに異なる時刻の状態を表示している.図 9 の横 軸は時刻を示しており,縦軸はその時刻でのノード 内パケット数の総和を表している.図 8 から,時刻 $t = 0.06 \, \text{s}$ にフロー 2 が流入したことによるふくそう が, $t = 0.073 \, \text{s}$ でいったん解消していることが分か る.しかし,その後 $t = 0.082 \, \text{s}$ で再び小さなふくそ うを起こしている.図 9 によれば,1度目のふくそう が解消された後に,2度目,3度目の小さなふくそう が繰り返し起こっていることが分かる.つまり,ネッ トワーク内のノードに存在するパケット数が平滑化し, いったんふくそうが解消したにもかかわらず,何らか の原因で再度ふくそう状態が引き起こされていること になる.

次に,フローj = 1のリンク上に存在するパケット 数の総和の時間推移を図 10 に示す.この評価は,ネッ トワークがフロー 1 のパケットをどれだけ運んでいる かを示している.シミュレーション条件より1本のリ



Fig. 10 Temporal evolution of the total number of packets on links.

ンク上に乗る最大パケット数が 100 なので, リンクの 帯域がフローj = 1 とフローJ = 2 で公平に分けら れたとするならば,フローj = 2が流入した後の,フ ローj = 1のリンク上のパケット数の総和の適切な値 は,約 1450 である.図 10 から,2 度目,3 度目のふ くそうが起こっていた付近で,リンク上のパケット数 の総和の値が適切な値を超えていることが分かる.

これまでの評価 (ノードは通過フロー数を知ってい るという状況下)では,1度ふくそうが解消されると 再度ふくそうが現れるということはなく,ネットワー ク性能は安定していた.今回の評価において複数回ふ くそうが起こる原因として,ネットワーク内への過剰 なパケット流入と,フローj = 1のノード 15 への過 剰なパケット送出の2点が考えられる.

ネットワーク内へのパケット流入量を決定するとき に用いられるのは, $\ell_i^j(t)$ であり, 各ノードの転送レー トの決定するときに用いられるのは r^j(t) である. 双方 ともふくそうが複数回起きる原因となっているので,そ れぞれの値の決定に共通して用いられている $M_i^{\rm act}(t)$ の決定方法に問題がある可能性が高い.これを確かめ るため $B_i/M_i^{\text{act}}(t)$ の値がどのように変化しているか を検証する.図11は各ノードでのフロー当りの利用 可能帯域 $B_i/M_i^{act}(t)$ をノードごとに表示し, グラフ ごとに異なる時間の状態を表示したものである. 各グ ラフの縦軸は利用可能帯域 $B_i/M_i^{\rm act}(t)$, 横軸はノー ド ID を表す. ノード 1 から 14 まではフロー j = 1 の 1本しか存在しないため, $B_i/M_i^{\rm act}(t) \simeq 12 \, {\rm Gbit/s} \, \mathbb{R}$ 度となり, ノード 15 以降は 2 本のフローで帯域を分 け合うため,その半分程度の値となるのが適切である. この図より,1度目のパケット平滑化が行われた後,適 切な値から大きく外れていくことが確認できる.これ





らのことから, $\ell_i^j(t)$ と $r_i^j(t)$ の決定に用いた $M_i^{act}(t)$ の振舞いが,DFCの適切な実装を妨げていることが分かる.

5. DFC の実装に向けた方式の改良

前章で, DFC を実装するためには $\ell_i^j(t)$ と $r_i^j(t)$ の 決定に用いられていた $B_i/M_i^{
m act}(t)$ の値が DFC のね らいどおりに変化していないことが分かった.これま で,シミュレーションモデルを用いた DFC の基本動 作特性の評価では,フィードバック情報計算の対象と なるフロー数,つまりノードを通過するアクティブフ ロー数 $M_i^{\text{act}}(t)$ の情報はシミュレーションの初期設定 のパラメータとして決定され,ノードにとって既知の 情報としていた.このようなシミュレーション条件で は,アクティブフロー数 $M_i^{\text{act}}(t)$ の決定が不適切で あっても DFC の動作に深刻な影響を与えることはな かった.しかし,モジュール構成を意識した DFC の 実装では,ノードがアクティブフロー数 $M_i^{\text{act}}(t)$ の情 報を知ることができるとは限らない.そのため DFC を実装するモジュール内に、ノード自身がアクティブ フローに関する情報を取得する仕組みを組み込まなけ ればならない.

本章では $\ell_i^j(t) \ge r_i^j(t)$ の決定方法の考察とノード 自身がアクティブフローに関する情報を取得する仕組 みの考察を行い, DFC の方式改良を行う.

5.1 $\ell_i^j(t) \geq r_i^j(t)$ の制限方法

式 (2) の転送レート $\tilde{J}_{i}^{j}(t)$ はネットワーク内のノード内パケット数平滑化のために計算された値であるが,フロー j のみに着目して決定した値であり,他のフローの影響を考慮に入れていない.そのため,実際のパケット送出のレート $J_{i}^{j}(t)$ は,式 (1) に示すように $L_{i}^{j}(t)$ によって $\tilde{J}_{i}^{j}(t)$ を制限した値となっている.

式 (14) から分かるように , $L_i^j(t)$ はノード内の他の フローの転送レートの値を考慮して算出されたもので ある .

これまで, $\ell_i^j(t) \ge r_i^j(t)$ の決定には $B_i/M_i^{act}(t)$ の 値で上限を決めていたが,転送レートの決定式 (1) と 同様に $L_i^j(t)$ の値を上限値として用いれば,他のフ ローの影響を考慮に入れた DFC の転送レート決定法 と整合がとれる.つまり, $\ell_i^j(t) \ge r_i^j(t)$ の値をそれ ぞれ,

$$\ell_i^j(t) = \min(\ell_{i+1}^j(t), L_i^j(t))$$
(15)

$$r_{i}^{j}(t) = \max(0, \min(\tilde{J}_{i}^{j}(t), L_{i}^{j}(t)))$$
(16)

とすれば, ネットワーク全体のパケット平滑化という ねらいが考慮された値になると考えられる. $L_i^j(t)$ を 適切に決定するにはアクティブフローの判定が必要で あり, 逆にアクティブフローの判定が適切に実行でき れば,自然にアクティブフロー数 $M_i^{act}(t)$ の値も明ら かになる.そのため,以降ではアクティブフローの判 定法と $L_i^j(t)$ の決定法を考察する.

5.2 アクティブフロー数の測定法と $L_i^j(t)$ の決定法

フィードバック情報が生成される間にどのフローが アクティブフローであるかを測定する方法は,従来の アクティブフロー数の測定方法を変更することで実 現できる(図12).具体的には,まず DFCQueue がパ ケットを更新するごとに,それぞれのフローについて, DFCF1owQ と DFCQueue にパケットが存在しているか 否かを調べる.調べたフローにパケットが存在するな らば,そのフローのフロー情報にフラグを立てる.そ の後,DFCDBAgent がフィードバック情報を生成する ごとにフラグをリセットする.この操作を行うことで, フィードバック情報が生成される間にどのフローがア クティブであったかを測定することができる.

DFCDBAgent はこの得られた情報を参考に,パケット転送レートやフィードバック情報の計算対象となる



図 12 アクティブフロー数の測定

Fig. 12 Measurement of the number of active flows.

アクティブフローを特定することができる.また,式 (13) により $W_i^j(t)$ を計算し,式 (14) により $L_i^j(t)$ を 決定することができる.

5.3 パケット送出方法に関する考察

DFCの実ネットワークへの実装を考えるとき,DFC はデータリンク層で実現され,DFCQueue は例えばネッ トワークカード内に実装されることが考えられる.こ のため,もしリンクへのパケット送出方法を簡略する ことができれば,既存のネットワークカードにDFC を実装することが容易になると考えられる.

従来の DFC では, DFCFlowQ から DFCQueue への パケット転送レートを $\tilde{J}_{i}^{j}(t)$ とし, DFCQueue からリン クへの転送レートを WRR を用いて J^j_i(t) に制限して いた.しかし,今回の改良によりアクティブフロー数 の情報を DFCDBAgent が利用できるようになるので, DFCFlowQ から DFCQueue へのパケット転送レートを $J_i^j(t)$ とすることが可能となった.この $J_i^j(t)$ は,下流 のリンク帯域と,他のフローの転送レートを考慮に入 れた値となっているので, DFCQueue からリンクへの パケット送出をよりシンプルな Round Robin (RR) や Single Server Queue に変更しても, 従来の性能を 保つことができると考えられる (図 13 参照). DFC 機能を実装を考えると, DFCQueue では単純な FIFO queue を用いた実装が可能となり, ハードウェアの簡 略化が実現することに結び付く.ここで,DFCFlowQ から DFCQueue へのパケット転送レート $\tilde{J}_i^j(t)$ が

$$\sum_{j=1}^{M_i} 1_{\{j = \text{active}\}} \times \tilde{J}_i^j(t) < B_i \tag{17}$$

を満たすならば,アクティブフローでないフローの パケットであっても DFCFlowQ から DFCQueue へ転 送することができる.そのときのパケット転送レート は $B_i - \sum_{j=1}^{M_i} 1_{\{j=\text{active}\}} \times \tilde{J}_i^j(t)$ となる.このため, ノード*i*においてアクティブフローでないと判断され,



図 13 FIFO (Single Server Queue)方式 Fig. 13 Single server queue-based rate control.

パケット転送レート 0 が割り当てられたパケットが ノード *i* に到着したとしても, 混雑していなければ次 のレート更新時刻を待たずに転送可能である.

6. 評 価

上記の改良を行った ns2 を用いてシミュレーション を行い,自律分散制御の装置構成をもつ DFC 機能が 適切な性能を発揮することを確認することで,改良 した DFC の妥当性を評価する.なお,今回のシミュ レーションでは,DFCQueue からリンクへのパケット 送出方法を Single Server Queue に変更した.

6.1 ふくそうの繰返しを避ける効果についてのシ ミュレーションの評価結果

まず,4.1 で用いたモデルとシナリオを用いて,ふ くそうの繰返しを避ける効果についての評価を行う. 図 14 と図 15 はそれぞれ改良前と改良後のフロー *j* = 1 の,ノード上パケット数の総和の時間推移,リ ンク上パケット数の総和の時間推移を示したものであ る.それぞれ,実線が改良後,破線が改良前の値を示



図 14 ノード上のパケット数の総和の時間変化の比較 Fig. 14 Comparison of temporal evolution of the total number of stored packets in nodes.



図 15 **リンク上のパケット数の総和の時間変化の比較** Fig. 15 Comparison of temporal evolution of the total number of packets on links.

している.これを見ると, 複数回起こっていたふくそうが起こらなくなり, それに伴ってリンク上のパケット数も改良前に比べて早い段階で安定した値になっていることが分かる.

次に,ノード上の総パケット数及びリンク上の総パ ケット数に関して,DFC機能による収束の速さを改 良前後で比較する.図16は,横軸がふくそうが開始 してからの時間,縦軸がその時間以降のノード上のパ ケット数の総和の最大値を表す.実線が改良後,破線 が改良前の値を示している.この図から改良後は改 良前に比べてノード上の総パケット数が速やかに減少 していることが分かる.例えば,ノード上のパケット 数の総和が30個以下に収束する時間は,改良前では 0.06s,改良後では0.015sであり,改良前の約25%ま で削減できる.また図17は,横軸がふくそうが開始 してからの時間,縦軸がリンク上のパケット数の総和 と理想値(1450)との差の比率を表す.この結果から,









改良後のリンク上のパケット数の総和が改良前に比べ て素早く理想値に近づいていることが分かる.例え ば,理想値との誤差が1%未満になる時間は,改良前 では 0.055 s, 改良後では 0.024 s であり, 改良前の約 40%まで削減できる.

これらのことから, $\ell_i^j(t)$ と $r_i^j(t)$ の決定に用いら れていた $B_i/M_i^{\rm act}(t)$ を $L_i^j(t)$ に変更することで,フ ローj = 1の転送レートの値が適切な値となり,また ネットワーク入口部分のシェーピング機能がうまく働 き,ネットワークへの過剰なパケット流入が抑えられ, ふくそうが複数回起こるという現象を解決できたこと が分かった.

6.2 アクティブフロー数の適切な管理についての シミュレーションモデルと評価結果

ここでは,アクティブフロー数が適切に管理されて いるかについての評価を行う.この評価では,これま で用いたモデルを以下のように変更した.ノード 23 に新たにソースノードを接続し, $t = 0.09 \, \text{s}$ からノー ド30に向けてパケット送出を開始する(このフローを フローj = 3とする).更に同じ時刻に,フローj = 2のパケット生成を停止させる.

まず,ネットワーク内からフローi = 2のパケット が消滅するまでの時間を調べるために,フローj = 2のリンク上パケット数の総和の時間推移についてのグ ラフを示す(図18).この図からフローj = 2のすべて のパケットがノード 30 に到達した時刻が約 t = 0.15 s であることが分かる.この時刻以降のフローi = 1の転送レート $J_i^1(t)$ の変化を調べたものが図 19 であ る.横軸はノード ID,縦軸はそのノードの転送レー ト $J_i^1(t)$ を表し, それぞれのグラフは異なる時刻の状 態を表示している.

フローi = 2のパケットがネットワーク内から消滅 すると, ノード 30 からノード1 に向かって転送レート



flow 2's packets on links.

 $J_i^1(t)$ の値がフローが 2 本の場合の値 (このとき存在 するフローは j = 1, 3) に回復していき, t = 0.156 s にすべてのノードの転送レートの値が回復したことが 分かる.つまり,t = 0.156s 付近からフローj = 1の リンク上パケット数の総和の値が,1450パケット付近 まで回復していれば, それぞれのノードでアクティブ フロー数が適切に管理されていることが分かる.

図 20 はフロー *j* = 1 のリンク上パケット数の総和 の時間推移を示したものである.これを見ると,それ ぞれの時刻のフロー数に応じて,リンク上パケット数 の総和が適切に変化していることが分かる.具体的 には,アクティブフローが二つである t = 0.06 s から $t = 0.09 \, \mathrm{s}$ にかけては 1450 パケット付近, アクティブ フローが三つである t = 0.09 s から t = 0.150 s にか けては 966 パケット付近, アクティブフローが一つ減









少した t = 0.150 s 以降は 1450 パケット付近となって いることが分かる.

これらのことから,ノード上でのアクティブフロー 数が適切に管理されていることが確認でき,本研究で 示したアクティブフロー数に関する情報の取得方法が 適切であることが確認された.

7. む す び

拡散型フロー制御(DFC)とは,拡散現象を指導原 理としたフロー制御で,局所的なネットワーク情報に 基づいたノードの自律動作が,間接的にネットワーク 全体の状態を良好な特性に導く制御方式である.これ まで,DFC がねらいどおりの特性をもつことがシミュ レーションにより確かめられてきたが,その評価では ノードを通過するフロー数が既知であることを用いて いた.実ネットワークにDFCを実装するためには, ノードを通過するアクティブフロー数の情報をノード が自律的に取得できる枠組みにする必要がある.

本論文では DFC の実装可能性を確かめるため ns2 に DFC 機能を組み込み,自律分散制御として実装可 能であるかどうか検討を行った.その結果,従来の DFC の枠組みではノードの通過フロー数の情報を適 切に取得することができないため,ネットワーク内に 流入するパケット数を制御する機構が正しく動作せ ず,ネットワークがふくそう状態から回復するまでに 時間がかかってしまうことが分かった.これを解決す るために,フィードバック情報として通知される転送 レートの制限方法にを改良し,更にノード自身が通過 するアクティブフロー数を計測する方法を改良した. この改良により, DFC の実装に関して従来 weighted round robin が必要であったパケットスケジューリン グ部分を,単純な FIFO で置き換える単純化が可能と なった.また,改良した DFC の特性を ns2 によりシ ミュレーションを行い,その効果を実証した.

謝辞 本研究の一部は,情報通信研究機構(NICT) の委託研究「新世代ネットワークの構成に関する設計・ 評価手法の研究開発」により実施した.

文 献

- Y. Bartal, J. Byers, and D. Raz, "Global optimization using local information with applications to flow control," Proc. 38th Ann. IEEE Symp. on Foundations of Computer Science, pp.303–312, Oct. 1997.
- [2] S.H. Low and D.E. Lapsley, "Optimization flow control-I: Basic algorithm and convergence," IEEE/ACM Trans. Netw., vol.7, no.6, pp.861–874,

1999.

- [3] K. Kar, S. Sarkar, and L. Tassiulas, "A simple rate control algorithm for maximizing total user utility," Proc. IEEE INFOCOM 2001, pp.133–141, 2001.
- [4] J. Mo and J. Walrand, "Fair end-to-end window based congestion control," IEEE/ACM Trans. Netw., vol.8, no.5, pp.556–567, Oct. 1999.
- [5] S. Kunniyur and R. Srikant, "A decentralized adaptive ECN marking algorithm," Proc. IEEE GLOBE-COM '00, pp.1719–1723, 2000.
- [6] R. Johari and D. Tan, "End-to-end congestion control for the Internet: Delays and stability," IEEE/ACM Trans. Netw., vol.9, no.6, pp.818-832, Dec. 2001.
- [7] C. Takano and M. Aida, "Stability and adaptability of autonomous decentralized flow control in highspeed networks," IEICE Trans. Commun., vol.E86-B, no.10, pp.2882–2890, Oct. 2003.
- [8] C. Takano, M. Aida, and S. Kuribayashi, "Autonomous decentralized flow control in high-speed networks with inhomogeneous configurations," IE-ICE Trans. Commun., vol.E87-B, no.6, pp.1551– 1560, June 2004.
- [9] C. Takano and M. Aida, "Diffusion-type autonomous decentralized flow control for end-to-end flow in highspeed networks," IEICE Trans. Commun., vol.E88-B, no.4, pp.1559–1567, April 2005.
- [10] C. Takano, K. Muranaka, K. Sugiyama, and M. Aida, "Mutual complementarity of diffusion-type flow control and TCP," IEICE Trans. Commun., vol.E89-B, no.10, pp.2850–2859, Oct. 2006.
- [11] C. Takano, K. Muranaka, K. Sugiyama, and M. Aida, "Parameter design for diffusion-type autonomous decentralized flow control," Lecture Notes in Computer Science, vol.4238, pp.461–470, Springer-Verlag Heidelberg, 2006.
- [12] C. Takano and M. Aida, "Diffusion-type autonomous decentralized flow control for multiple flows," IEICE Trans. Commun., vol. E90-B, no.1, pp.21–30, Jan. 2007.
- [13] The Network Simulator—ns2. http://www.isi.edu/nsnam/ns/
 (平成 20 年 1 月 15 日受付, 5 月 13 日再受付)



高野知佐(正員)

平12 阪大・工・電子通信卒.平20 首都 大東京大学院博士後期課程了.平12 NTT アドバンステクノロジ(株)入社.平20 広島市大大学院・情報科学研究科・准教授. 通信トラヒック制御,自律分散制御技術の 研究に従事.平14 本会学術奨励賞,平15

本会情報ネットワーク研究賞受賞.



山内正志(正員)

平 18 都立科技大・工・生産情報システム 工卒.在学中は自律分散制御技術の研究に 従事.現所属はトーワシップ債権回収(株).



会田 雅樹 (正員)

昭 62 立教大・理・物理卒.平元同大大 学院博士課程前期原子物理学専攻了.同年 日本電信電話(株)入社.平 17 首都大東 京・システムデザイン学部・准教授.平 19 同大大学院・システムデザイン研究科・教 授.博士(工学).通信トラヒックの設計技

術,通信品質計測技術,自律分散制御技術,通信トラヒックの べき乗則の研究に従事.日本 OR 学会,IEEE 各会員.