

## 3H-2 Ethernet 上での 全順序放送通信プロトコルの設計と実現

中村 章人、滝沢 誠  
(東京電機大学理工学部)

### 1.はじめに

現在、複数の通信実体の協調動作が分散型データベースシステムで同時実行制御、コミットメント制御、分散型問合せ処理等を実現する上で必要とされており、このためには、複数の実体間での高信頼放送通信が求められている。本論文では、EthernetのMAC層で提供されている低信頼放送サービスを利用して複数の実体に対して全順序放送(TO)サービスを提供するプロトコルについて述べる。TOサービスを利用すると、全実体が同一のPDUを同一の順序で受信する。また、本プロトコル(TOプロトコル)の実現方法についても述べる。

### 2. 信頼性

従来のプロトコルでは、各実体はPDU pが正しく受信されたことを、pが到着し、かつpがあるある条件を満たすこと、例えば通番のチェック等によって決定する。しかし、放送通信サービスに対する正しい受信の概念はこれより複雑である。正しい受信を決定するための方法として、ある1つの主制御実体がこれを行なう集中型制御と、各実体が互いに通信し合うことによって自分自信で正しい受信の決定を行なう分散型制御がある。TOプロトコルでは、後者を用いている。

[定義] 各実体E<sub>k</sub>は、以下の条件を満たすとき、PDU pを正しく受信したという。(1) E<sub>k</sub>はpを受信する。(2) E<sub>k</sub>は「pの宛先内の全実体がpを受信した」ことを知る。(3) E<sub>k</sub>は「pの宛先の各実体は、全宛先実体がpを受信したことを知っている」ことを知る。□

(1)、(2)、(3)をそれぞれpはE<sub>k</sub>で受信された、前確認された、確認されたという。複数の実体がPDUを放送する状況で、各実体が同一のPDUを同一の順序で受信するような高信頼放送サービスを定義する。

[定義] 各PDU pとqが共通の宛先のSAPで同一の順序で受信されるとき、かつそのときに限りこの(N-1)高信頼放送サービスを全順序(TO)サービスという。□

従来の2つのSAP間のコネクションの概念をn(>2)個に拡張したものが、群(cluster)である。

[定義] 1チャネル(1C)サービスとは、PDUの紛失の可能性はあるが、受信されたPDUは全実体で同一の順序である低信頼放送サービスである。□

LANのMAC層と無線網は1Cサービスを提供している。

### 3. 1Cサービス上での全順序(TO)プロトコル

群はn個の実体E<sub>1</sub>, ..., E<sub>n</sub>で構成される。各実体E<sub>k</sub>が送信したPDUをE<sub>k</sub>は受信できるとする。各実体E<sub>k</sub>について、送信ログSL<sub>k</sub>=(SP<sub>k</sub>, <<<sub>k</sub>)をE<sub>k</sub>が既に送信したPDUの系列、受信ログRL<sub>k</sub>=(RP<sub>k</sub>, <<sub>k</sub>)をE<sub>k</sub>が既に受信したPDUの系列とする。ここでSP<sub>k</sub>とRP<sub>k</sub>は各々、E<sub>k</sub>が送信したPDU

と受信したPDUの集合である。E<sub>k</sub>がq以前にpを受信したときRL<sub>k</sub>内でp <<sub>k</sub> qであり、q以前にpを送信したときSL<sub>k</sub>内でp <<<sub>k</sub> qである。p <<<sub>k</sub> qでp <<<sub>k</sub> gかつg <<<sub>k</sub> qなるgが存在しないときp →<sub>k</sub> qである。

[定義] RL<sub>k</sub>内の全てのpとqに対して、もし pとqがE<sub>j</sub>によって放送され、SL<sub>j</sub>内でp →<sub>j</sub> qで、RL<sub>k</sub>内でp <<sub>k</sub> qならば受信ログRL<sub>k</sub>は正しいという。□

RL<sub>k</sub>を、RL<sub>k</sub>の先頭からf番目までの部分系列とする。

[定義] 受信ログRL<sub>1</sub>, ..., RL<sub>n</sub>に対して、以下の条件が成り立つとき、かつそのときに限りPDU番号fを障害点FPとする。(1)全てのjとkについて、RL<sub>j</sub><sup>f+1</sup> = RL<sub>k</sub><sup>f+1</sup>で、これらの受信ログは正しい。(2)あるjについて、RL<sub>j</sub>は正しくない。□

FPの存在は、ある実体があるPDUを受信できなかつたことを意味する。APL<sub>k</sub>、PPL<sub>k</sub>、RPL<sub>k</sub>をそれぞれ確認されたPDU、前確認されたPDU、受信されたPDUからなるS<sub>k</sub>の部分ログとする。E<sub>k</sub>が放送するPDU pはE<sub>k</sub>が既に受信したPDUの確認通知を含んでいる。E<sub>k</sub>が放送する各PDU pは以下の情報を持っている。

p.A<sub>j</sub> = E<sub>k</sub>が次にE<sub>j</sub>から受信予定のPDUの通番  
(j=1, ..., n)。

p.SEQ = pの通番。

p.SRC = pを放送する実体、つまりE<sub>k</sub>。

p.BUF = E<sub>k</sub>内で利用可能なバッファ数。

各E<sub>k</sub>は、変数SEQ<sub>k</sub>とPDUの通番をチェックするためのマトリックスAL<sub>qa</sub>(p, q=1, ..., n)を持っている。

SEQ<sub>k</sub> = E<sub>k</sub>が次に放送する予定のPDUの通番。

AL<sub>jj</sub> = E<sub>k</sub>が次にE<sub>j</sub>から受信予定のPDUの通番。

AL<sub>ns</sub> = E<sub>s</sub>が次にE<sub>n</sub>から受信予定のPDUの通番。を、AL<sub>11</sub>, ..., AL<sub>nn</sub>の中の最小値minAL<sub>s</sub>は、全実体がg. SEQ < minAL<sub>s</sub>なるPDU gをE<sub>s</sub>から既に受信していることを意味している。群内の各実体は、群開設手続き[TAK87a,b]によって各E<sub>j</sub>の初期通番ISS<sub>j</sub>と、空バッファ数IBF<sub>j</sub>を知っている。最初はSEQ<sub>k</sub>=ISS<sub>k</sub>でAL<sub>jj</sub>=ISS<sub>j</sub>(h, j=1, ..., n)。各実体は変数F<sub>1</sub>, ..., F<sub>n</sub>を持ち、F<sub>j</sub>はE<sub>j</sub>内の空バッファ数であり、最初F<sub>j</sub>=IBF<sub>j</sub>(j=1, ..., n)である。minFをF<sub>1</sub>, ..., F<sub>n</sub>の中の最小値とする。

PDU pの送受信は、以下の条件を満たすとき行われる。

[受信条件] (1)p.SEQ = AL<sub>jj</sub>. (2)p.A<sub>n</sub> < AL<sub>nn</sub>.  
(h=1, ..., n). □

[受信動作] (1)AL<sub>jj</sub> = AL<sub>jj</sub> + 1, AL<sub>ns</sub> = p.A<sub>n</sub>  
(h=1, ..., n). (2)F<sub>j</sub> = p.BUF. (3)pをRL<sub>k</sub>の最後尾、つまりRPL<sub>k</sub>の最後尾に追加する。□

[フロー条件] (1)minAL<sub>k</sub> ≤ SEQ<sub>k</sub> < minAL<sub>k</sub> + min(W, minF/n<sup>2</sup>). (2)各実体がPDUを受信可能である。□

[送信動作] (1)p.BUF=F<sub>k</sub>, p.A<sub>j</sub>=AL<sub>jj</sub>(j=1, ..., n).  
(2)p.SEQ=SEQ<sub>k</sub>, p.SRC=E<sub>k</sub>, SEQ<sub>k</sub>=SEQ<sub>k</sub>+1. (3)E<sub>k</sub>はpを放送する。□

フロー条件が満たされなければ、E<sub>k</sub>は十分なバッファが利用できるようになるまで放送を待つ。

分散型制御において各実体E<sub>k</sub>がpの正しい受信をど

のように決定するかが問題となる。ここで、pをE<sub>j</sub>が放送したPDUとする。

- [前確認(PACK)条件] p.SEQ < minAL<sub>j</sub>. □
- [前確認動作] RPL<sub>k</sub>内の先頭からpまでのPDUをRPL<sub>k</sub>から取り出し、PPL<sub>k</sub>の最後尾に追加する。□

$$\begin{array}{ccc} \text{PPL}_k & \text{RPL}_k \\ \langle \dots a \rangle < b \dots p q \dots ] & ==> \\ \langle \dots a b \dots p ] < q \dots ] \end{array}$$

[定義] 各E<sub>j</sub>とRL<sub>k</sub>内のPDU pとqに対して、qがE<sub>h</sub>でpを前確認し、E<sub>h</sub>でpを前確認するg <<sub>k</sub> qなるPDU gが存在しないときqはE<sub>h</sub>についてpを最初前確認するPDUという。各E<sub>j</sub>について、pを最初前確認するPDUの中で、最後に受信したものをpを最後前確認するPDUという。□

[確認(ACK)条件] PPL<sub>k</sub>内のPDU pを最後前確認するPDUが前確認される。□

[確認動作] PS<sub>k</sub>内の先頭からpまでのPDUをPPL<sub>k</sub>から取り出し、APL<sub>k</sub>の最後尾に追加する。□

$$\begin{array}{ccc} \text{APL}_k & \text{PPL}_k & \text{RPL}_k \\ \langle \dots a \rangle < b \dots p q \dots ] < .g.. ] & ==> \\ \langle \dots a b \dots q ] < q \dots g ] < .. ] \end{array}$$

PDUが粉失した点、つまり障害点FPはE<sub>k</sub>がPDU pを受信する度に発見し、リセット手続きが適用される。

[粉失条件] (1) E<sub>j</sub>からpを受信したとき、p.SEQ > AL<sub>jj</sub>ならば、AL<sub>jj</sub> ≤ g.SEQ < p.SEQなるPDU gをE<sub>j</sub>から受信していない(j=1, ..., n)。(2) E<sub>j</sub>からqを受信したとき、あるj(j ≠ h)について、q.A<sub>j</sub> > AL<sub>jj</sub>ならば、E<sub>j</sub>はAL<sub>jj</sub> ≤ g.SEQ < q.A<sub>j</sub>なるPDU gをE<sub>j</sub>から受信していない(h=1, ..., n)。□

[リセット手続き] (1) E<sub>k</sub>はPDUの受信を停止し、r.A<sub>h</sub>=AL<sub>hh</sub>(h=1, ..., n)なるRST(RESET) PDU rを放送し、SEQ<sub>k</sub>=AL<sub>kk</sub>とする。(2) E<sub>j</sub>がE<sub>k</sub>からRST rを受信したら、各hについてRL<sub>k</sub>がr.A<sub>h</sub> ≤ p.SEQなるE<sub>h</sub>からのPDU pを含むならば、p <<sub>k</sub> gである全てのPDU gをRL<sub>k</sub>から取り除く。同時に、ALリセット手続きを適用してALを更新する。E<sub>j</sub>はPDUの受信を停止し、ra.A<sub>h</sub>=AL<sub>hh</sub>(h=1, ..., n)

なるRST\_ACK PDU raを放送し、SEQ<sub>j</sub>=AL<sub>jj</sub>とする。(3)全実体からRST\_ACKを受信し、各hについてr.A<sub>h</sub>=AL<sub>hh</sub>であれば、PPL<sub>k</sub>とRPL<sub>k</sub>内のPDUをAPL<sub>k</sub>に移す。□

[ALリセット手続き] 各E<sub>j</sub>とp=last(RL<sub>kj</sub>)に対して、AL<sub>jh</sub> = p.SEQとする(h=1, ..., n)。□

PDUの粉失は、タイムアウト処理によっても発見される。ある実体E<sub>k</sub>がPDU pを放送したにも関わらず、E<sub>k</sub>を含む全実体がpを受信していないとする。E<sub>k</sub>がp → gなるPDU gを放送しなければ、各実体はpの粉失を発見出来ない。E<sub>k</sub>はpを放送した後、タイムアウトが起きたらpの粉失を発見する。この場合、どの実体もpを受信していないので、E<sub>k</sub>はpを再放送できる。

#### 4. 評価

本プロトコルの性能を、PDU数と必要時間で評価する。PDU数は、PDU pを確認するために送信されるPDU数とする。ラウンドとは、あるSAPから宛先までの最大遅延時間である。必要時間は、ラウンド数とする。nを実体数としたとき、性能は以下のようになる。

PDU数	最良 $1 + (n-1) + n = 2n$ 個
	最悪 $1 + (n-1) + (n-1)^2 = n^2 - n + 1$ 個

必要時間 最良 3 ラウンド

Ethernet MACサービスの様にPDUを並行放送できない場合 最良  $1 + (n-1) + n = 2n$  ラウンド

最悪  $1 + (n-1) + (n-1)^2 = n^2 - n + 1$  ラウンド

#### 5. T0プロトコルの実現

データ転送手続きは、3つのキューPRQUE、ARQUE、PTQUEにより実現されている。各E<sub>k</sub>の部分ログPPL<sub>k</sub>、RPL<sub>k</sub>は各々ARQUE、PRQUEに記憶される。

[送信] フロー条件が満たされたら、PDU pをPTQUEに記憶し、pを放送する。□

[受信] E<sub>j</sub>からpを受信したら、(1)もしpが受信条件を満たせば、pを受け入れPRQUEに記憶する。(2)PRQUEの先頭のPDU qが前確認条件を満たす間、qをPTQUEから取り出しARQUEに記憶する。p.SRC = E<sub>k</sub>ならば、pをPTQUE、つまりPTQUEの先頭から取り出す。q.RPT=gであれば、ARQUEの先頭からgまでのPDUをARQUEから取り出す。つまり、E<sub>k</sub>で確認される。もしPRQUEから最後に取り出したPDUがqであればp.RPT=qとする。(3)もし放送するデータがあれば送信動作を行い、そうでなければ確認通知を含むPDUを放送する。□

[障害] (1)もし粉失条件が成り立てば、リセット手続きを行う。□

#### 6. まとめ

現在、T0プロトコルをEthernet MACサービスを用いて、Unix 4.2BSDとSystem V上でC言語で実現している。プロセス間の通信には、Unixのsocketインターフェイスを用いている。Ethernet上には、現在m380q、a400(Facom)、Sun3ワークステーション3台が接続されている。現在、このシステムを使ってT0プロトコルの性能評価を行っている。問題はバッファサイズで、本プロトコルでは、従来の一対一通信に比べて最悪の場合、n<sup>2</sup>倍以上のバッファが必要になる。しかし、今日のハードウェア技術により記憶装置のコストは非常に低くなっているので、この問題は数M byteの記憶装置を用いることで解決できると思われる。

#### 参考文献

- [BRA] Bracha, G. and Toueg, S., "Asynchronous Consensus and Broadcast Protocols," JACM, Vol. 32, No. 4, 1985, pp. 824-840.
- [CHA] Chang, J.-M. and Maxemchuck, M.F., "Reliable Broadcast Protocols," ACM TOCS, Vo. 2, No. 3, 1984, pp. 251-273.
- [NAK] 中村 章人, 滝沢 誠, "多チャネル上の全順序放送通信プロトコル", 情報処理学会MDP研究会39-1, 1988
- [SCH] Schneider, F., Gries, D., and Schlichting, R.D., "Fault-Tolerant Broadcasts," Science of Computer Programming, Vol. 4, 1984, pp. 1-15.
- [TAK87a] Takizawa, M., "Design of Highly Reliable Broadcast Communication Protocol," Proc. of IEEE COMPSAC87, 1987, pp. 731-740.
- [TAK87b] Takizawa, M., "Cluster Control Protocol for Highly Reliable Broadcast Communication," Proc. of the IFIP Conf. on Distributed Processing, Amsterdam, 1987.