

フォールトトレランスな群制御プロトコルの設計*

5 V - 2

鈴木等 中村章人 滝沢誠†

東京電機大学 ‡

1 はじめに

通信技術の発展と標準化に伴って、システムは分散型となってきた。グループウェアシステム等の新しい分散型応用では、二つの通信実体間の通信に加えて、複数の通信実体間のグループ通信が必要とされている。

これまでに、EthernetのMAC副層や無線システムが提供する低信頼な放送通信サービス上で、群と呼ばれる $n(\geq 2)$ 個のSAPに対して高信頼な放送通信サービスを提供する放送通信プロトコルが検討されている[1-4]。下位層でのプロトコルデータ単位(PDU)の紛失に加えて、実体の停止等の障害を考察する必要がある。

本論文では、実体が停止、復旧する動的群に対する高信頼な放送通信サービスを提供するためのプロトコルを検討する。本論文では、第2章で基本定義について、第3章で動的群制御の手続きについて論じる。

2 基本定義

2.1 1チャンネル(1C)サービス

本論文では、Ethernet MAC副層が提供する低信頼な放送通信サービスを抽象化したものを1チャンネル(1C)サービスとする。

[定義] 1Cサービスとは、PDUの紛失の可能性はあるが、ある実体が放送したPDUは、全実体において、同一の順序で受信されるサービスである。□

2.2 群

群とは、従来の2つのSAP間のコネクションの概念を、 $n(\geq 2)$ 個のSAP間に拡張した概念である。 $\langle N-1 \rangle$ 群スキーム Q とは、 $n(\geq 2)$ 個の $\langle N-1 \rangle$ SAP S_1, \dots, S_n の組 $\langle S_1, \dots, S_n \rangle$ である。 S_j でSDUを送受信している時、 S_j は活性(A)である。そうでなければ、非活性(I)とする。 $state(S_j)$ を S_j の状態とする。 S_j は E_j により提供されているとする。このとき E_j は C 内にあるとする。群はスキームとその例で与えられる[図1]。

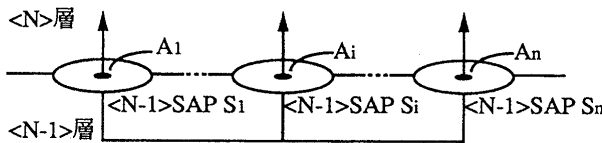


図1. $\langle N-1 \rangle$ 群

$$state(S_i) = E_i \text{ が提供する } S_j \text{ の状態 } \in \{A, I\}$$

$$state(C) = \langle state(S_1), \dots, state(S_n) \rangle$$

各群スキーム Q に対して、同時に高々1つの例 C が存在する。群には静的群と動的群の2種類がある。静的群で

*Design of Fault-tolerant Dynamic Cluster Control Protocol

†Hitoshi Suzuki, Akihito Nakamura, and Makoto Takizawa

‡Tokyo Denki University

は、群内のSAPが全て活性である。動的群では、SAPの状態が変化する。

2.3 ログ

各実体が利用するサービスを、PDUの系列であるログの集合としてモデル化する。ログ L とは、PDUの系列 $\langle p_1, \dots, p_m \rangle (m \geq 0)$ であり、各実体の通信履歴である。群内の各実体 E_k は、送信ログ SL_k と受信ログ RL_k を持つ($k = 1, \dots, n$)。L内でPDU p が q に先行するとき、 $p \rightarrow q$ とする。

[ログの同値関係]2つの受信ログ RL_i と RL_j について、以下の同値関係が存在する。

- (1) 順序同値: RL_i と RL_j の両方に含まれる任意の二つのPDU p と q について、一方で $p \rightarrow q$ ならば、他方でも $p \rightarrow q$ である。
- (2) 情報同値: RL_i 内のPDUの集合と、 RL_j 内のPDUの集合が同一である。
- (3) 同値: RL_i と RL_j が順序かつ情報同値である。

[ログの性質] RL_i について、以下の性質を定義する。

- (1) 順序保存: 各 E_j と、 RL_i 内の E_j から受信した任意の二つのPDU p と q について、 SL_j 内で $p \rightarrow q$ ならば、 RL_i でも $p \rightarrow q$ である。
- (2) 情報保存: RL_i 内のPDUの集合が、 SL_1, \dots, SL_n 内に含まれるPDU集合の和である。
- (3) 正しい: 順序保存でかつ情報保存である。

2.4 高信頼放送通信サービス

放送通信サービス内の活性な各SAP S_j を提供する E_j の RL_j が正しいならば、このサービスを信頼放送通信サービスとする。

[全順序放送通信(TO)サービス]サービス内の全受信ログが正しく、かつ互いに順序同値である。

TOサービスでは、各実体から受信したPDUの順序は、放送順序と同じであり、PDUの紛失がない。また、異なる実体から受信したPDUの受信順序も各実体で同一である。図2にTOサービスの例を示す。

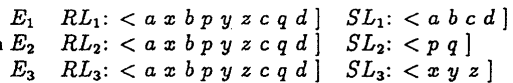


図2. 全順序放送通信(TO)サービス

3 動的群制御

動的群は、群スキームに対する動的群と群スキームの例に対する動的群の2種類について考える。

3.1 群スキームの例に対する動的群

信頼放送通信サービス[1-4]では、群内のあるSAP S_j が非活性になる、即ち E_j が停止した場合、群内の全実体に対する通信サービスを提供できない。例えば、冗長なモジュールを含んだシステムでは、あるモジュールが停

止しても残りのモジュールにより処理を続行できる。本論文では、実体停止障害をモデル化し、群内のある実体が停止した場合でも、他の実体間に TO サービスを提供するための手続きについて述べる。

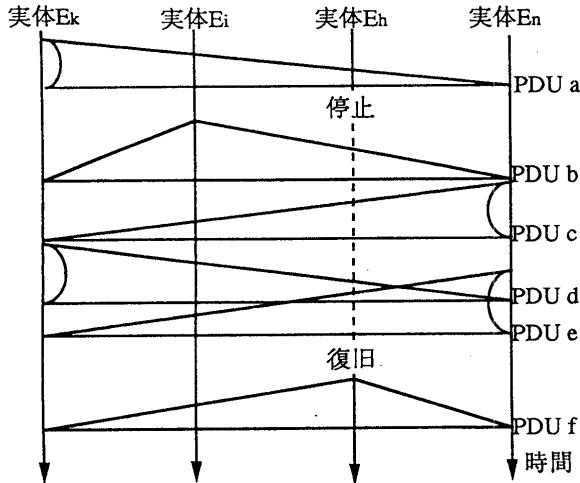


図3. 実体の停止と復旧

[実体停止] ある実体が、一定時間以上動作しないとき、停止していると看做す。□ 群内の実体は、ある実体の停止を、確認通知待ちでのタイムアウトによって検出する。

3.2 障害手続き

低信頼な 1C サービス上の TO プロトコル [1] を基本として、障害実体を検出し、正常動作している実体間で正しい送受信を行なうための手続きを示す。本方式は、データ転送を停止せずに動作できる。TO プロトコルにおける各 PDU p に以下の項目を追加する。また、各実体は停止した実体を示すビットマップ $emask$ を持つ。

$p.BMAP =$ 停止した実体を示すビットマップ。

[実体停止の検出条件] ある実体 E_k が、 E_h からの PDU を一定時間以上受信できないとき、 E_k は実体停止合意手続きを実行する。

[実体停止合意手続き] (1) E_k は、障害実体 E_h に対して、 $emask$ の h ビットを ON にし、 $s.ACK_i = REQ_i (i=1, \dots, n)$ 、 $s.BMAP = emask$ なる STOP PDU s を放送する。

(2) STOP 又は STOP_PACK s を受信した E_j は、 $emask = emask \vee s.BMAP$ 、 $AL_{ji} = s.ACK_j (j=1, \dots, n)$ として、 $sp.BMAP = emask$ 、 $sp.ACK_i = REQ_i (i=1, \dots, n)$ なる STOP_PACK sp を放送する。

(3) 各 E_i は、 $emask$ のビットが OFF の全実体から STOP_PACK 又は STOP_ACK を受信したなら、 $sa.BMAP = emask$ とした STOP_ACK sa を放送する。ここで、 $REQ_j = AL_{ji} (j=1, \dots, n)$ 、 $sp.ACK_j = REQ_j$ 。

(4) 各 E_i は、 $emask$ のビットが OFF の全実体から STOP_ACK sa を受信したとき、 $emask$ のビットが ON の障害実体 E_h から受信した PDU の中で通番が REQ_h 以上の PDU p^h を RL_i の中から削除する。

(5) 各 E_i は、 AL 行列内 [1-4] の E_h に関する確認情報に対して停止印を付け、以後データ PDU の前確認、確認条件に関して E_h の情報を無視する。

3.3 復旧手続き

停止していた E_h が、復旧したとき、 E_h を群に戻し、信頼放送通信サービスを再開するための手続きが必要である。本論文では、ある実体の復旧についての合意をとった後、データ転送を再開する手続きを紹介する。 E_h を復旧した実体とする。

[復旧手続き] (1) 障害復旧した E_h は、 $emask$ の全ビットを OFF にし、 $rc.BMAP = emask$ 、 $rc.SEQ = rc.ACK_h = 0$ なる RCV(RECOVERY) PDU rc を放送する。

(2) E_h から RCV を受信した E_j は、 $emask$ の h ビットを OFF とする。 E_j は、 E_h からの RCV 又は RCV_PACK rc を受信したならば、 $emask = emask \vee rc.BMAP$ とする。 E_j が復旧した実体ならば、 $REQ_h = rc.SEQ + 1$ 、 $AL_{hj} = rc.ACK_j$ として、 AL リセット手続きを実行する。また、 $rcp.BMAP = emask$ 、 $rcp.ACK_i = REQ_i (i=1, \dots, n)$ なる RCV_PACK rcp を放送する。

(3) 各 $E_j (j \neq h)$ は、 $emask$ のビットが OFF の全実体から RCV_PACK rcp を受信し、 $emask = rcp.BMAP$ ならば、 $REQ_h = AL_{hi} = rc.SEQ$ 、 $REQ_j = SEQ_j + 1$ 、 $AL_{hj} = rcp.ACK_j$ として、 AL リセット手続きを実行する。また、 $rca.BMAP = emask$ なる RCV_ACK rca を放送する。そうでなければ、 $emask = emask \vee s.BMAP$ なる RCV_PACK rcp を放送し、(3) を実行する。RCV_PACK 後に受信した RCV は、無視する。

(4) 各 E_i は、 $emask$ のビットが OFF の全実体から RCV_ACK rca を受信し、各 E_j について、 $rca.ACK_j = REQ_j$ ならば、 PRL_i と RRL_i 内の PDU の中で rca 以前に受信した PDU を ARL_i に移し、それ以外の PDU を RL_i 内から削除する。

(5) 各 E_i は、 AL テーブル内の E_h に対する停止印を削除し、データ PDU の送信を再開する。

[AL リセット手続き] 各 E_j と $p = last(RL_{kj})$ に対して、 $REQ_j = AL_{jh} = p.SEQ + 1 (j=1, \dots, n)$ 。

4 まとめと今後の課題

本論文では、データ PDU の送受信を停止することなく、実体停止時に他の実体で合意するための手続きを論じた。また、復旧手続きについて論じた。

今後は、データ転送を停止しない復旧手続きを考察する。また、実体の停止と復旧が同時に発生する場合や、実体停止以外の実体障害についても合意手続きを考察すると共に、群を構成する実体数が変化する群スキームに対する動的群制御プロトコルを設計し、フォールトトレランスなサービスを提供し続けるための手続きの検討を行う。

参考文献

- [1] Takizawa, M., "Cluster Control Protocol Highly Reliable Broadcast Communication," *Proc. of the IFIP Conf. on Distributed Processing*, 1987, pp.431-345.
- [2] Takizawa, M. and Nakamura, A., "Partially Ordering Broadcast Communication," *Proc. of the IEEE INFOCOM*, 1990, pp.357-364.
- [3] Nakamura, A. and Takizawa, M., "Reliable Broadcast Protocol for Selectively Ordering PDUs," *Proc. of the 11th IEEE ICDCS*, 1991, pp.239-246.
- [4] Nakamura, A. and Takizawa, M., "Priority-Based Total and Sifti-Total Ordering Broadcast Protocols," *Proc. of the 12th IEEE ICDCS*, 1992, pp.178-185.