

大規模データ収集システムにおける高信頼性マルチキャスト通信 Reliable Multicast Communication for a Large-Scale Data Acquisition System

梶山 真治†
Shinji Kajiyama

長坂 康史‡
Yasushi Nagasaka

1. はじめに

素粒子物理学実験における測定器の大型化に伴い、そこで用いられるデータ収集システムの開発が必要不可欠となってきた[1]。この大規模なシステムで問題となるのは、端末間の通信手段である。その中でも特に問題となるのが同一のメッセージを複数の端末に送信する1対多の通信手段である。本研究ではこの1対多の通信手段としてマルチキャスト通信に着目し、大規模データ収集システムのための高信頼性マルチキャスト通信を開発した。

2. システム概要

図1にシステムの構成とデータの流れを示す。このシステムは TR(Trigger)、TS(TriggerServer)、DH(DataHolder)、EB(EventBuilder)、DR(DataRecorder)の5つのコンポーネントから構成されている。

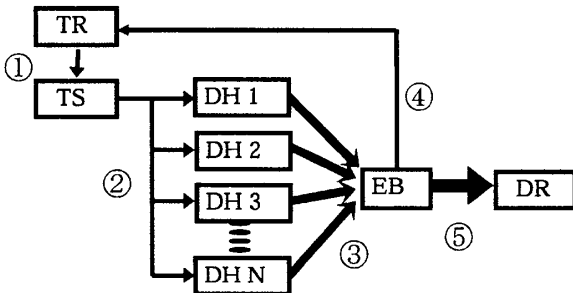


図1. システム構成とデータの流れ

このシステムの目的は、ある事象(イベント)の発生により、イベントフラグメントデータと呼ばれる複数のデータの断片を1ヶ所に集め、イベントごとのデータの集まり(イベントビルドデータ)として保存することである。イベントフラグメントデータはそれぞれのDHで周期的に発生し、DH内にバッファリングされる。イベント発生が引き金(Trigger)となりイベントフラグメントデータの収集が開始される。あるイベントに対する複数のイベントフラグメントデータをひとまとめとし1つのイベントビルドデータとすることをイベントビルドと呼ぶ。

データの流れは、TRからTSへイベント発生を知らせるTriggerメッセージを送信し(①)、TSはTriggerメッセージを受信すると複数のDHへTriggerメッセージを送信する(②)。それぞれのDHはTriggerメッセージを受信するとイベントフラグメントデータをEBに送信する(③)。EBはあるイベントに対して全てのDHからのイベントフラグメントデータを受信するとイベントビルドし、TRへイベントビルド応答を送信し(④)、DRへイベントビルドデータを送信する(⑤)。DRはEBからイベントビルドデータを受信するとデータをファイルに書き込む。EBからTRへの応答はTriggerメッセージの過剰送信を防ぐために用いられ、DHでのイベントフラグメントデータは周期的に生成される。以上

† 広島工業大学大学院工学研究科情報システム工学専攻、
Graduate School of Engineering, Hiroshima Institute of Technology

‡ 広島工業大学工学部知的情報システム工学科、
Faculty of Engineering, Hiroshima Institute of Technology

が1つのイベントに対するデータの流れであり、システム稼働時は複数のイベントが並列に処理される。

このシステムで想定しているコンポーネントの台数、コンポーネント間でやり取りされるデータサイズは、DH:1000台、その他は数台、Triggerメッセージは数十バイト、イベントフラグメントデータは数キロバイトとTriggerメッセージより大きい。

3. 通信システム

このシステムにおけるTS-DH間の通信は同一のメッセージを送信する1対多の通信である。この通信には信頼性と遅延の短さが要求される。この1対多通信の実現方法として2つの通信手段が考えられる。

まず1つめはTCP通信を用いてTSとそれぞれのDHとの間にコネクションを確立し、TSからそれぞれのDHへ1つずつメッセージを送信する手段である。TCP通信は、2つの端末間でデータを送受信する1対1のユニキャスト通信であり、信頼性を高めるための機能などを提供している。その反面、機能提供に伴いヘッダ処理にかかる時間が大きくなってしまいうという欠点を持っている。この方法では、DHの増加に伴いネットワークへ流れるデータ量も増加してしまうことと、それぞれのDHに1つずつメッセージを送信しなければならないため、それぞれのDHでのメッセージ受信時刻が大きくなってしまいTR-DH間の遅延が大きくなってしまふ。

2つめはマルチキャスト通信を用いる手段である。この通信はトランスポート層のプロトコルとしてUDPを利用し、TSから1回メッセージを送信するだけで複数のDHへ一斉にメッセージを届けることができるという1対多の通信手段である。DHが増加してもネットワークへ流れるデータ量は一定で、それぞれのDHでのメッセージ受信時刻のずれも小さい。さらに、マルチキャスト通信はUDP通信を用いているため、TCP通信のような再送制御やフロー制御といった信頼性を高める機能をサポートしていない。このためヘッダ処理にかかる時間が軽減される利点があるが、もしメッセージが何らかの原因でDHへ届かない場合、その通知もせず、復元も不可能となるので信頼性が低いという欠点がある。

両手段の問題点を比較した場合、TCP通信を用いた通信手段の問題点はTCP通信の仕様に大きく影響を受けており、変更が比較的困難であることから、マルチキャスト通信を用いた通信手段の開発を行った。

4. 高信頼性マルチキャスト通信

マルチキャスト通信の信頼性の低さから、TS-DH間の通信にこの通信をそのまま用いることはできない。そこでマルチキャスト通信にアプリケーションレベルの再送制御を付加することにより高信頼化を図った。その再送制御の形態には大きく分けて2つある。両者の違いは再送制御を送信側(TS)で行うか、受信側(DH)で行うかである。図2に両者のデータの流れを示す。

前者の再送制御を送信側で行う場合の流れは、まずTSからマルチキャストメッセージを送信する(①)。それぞれのDHでメッセージが正確に受信されると、TSへACK(肯定応答)を返信する(②)。図2のDH3のように、何らかの原因でメッセージがDHに届かず、ある時間が過ぎてもACK

が返ってこない場合、その DH に対して再送する (③)。このように再送のきっかけとなるメッセージの損失を送信側で発見するのが送信側再送制御方式 (Transmission-side Resend control System: TRS) である。

一方、再送制御を受信側で行う場合の流れは、まず TS からマルチキャストメッセージを送信する (①)。そして DH でメッセージを受信し、メッセージの内容から損失が発生したか調べ、もし損失が発生したとみなしたら、NACK (否定応答) を返信する (②)。TS は NACK を受信すると、再送を行う (③)。損失発生を調べる方法は、メッセージには識別番号が割り振られており TS から番号順に送信されてくるので、受信したメッセージのシーケンス番号が番号順でない場合、損失が発生したと認識する。このようにしてメッセージの損失を受信側で発見し送信側に知らせるのが受信側再送制御方式 (Reception-side Resend control system: RRS) である。

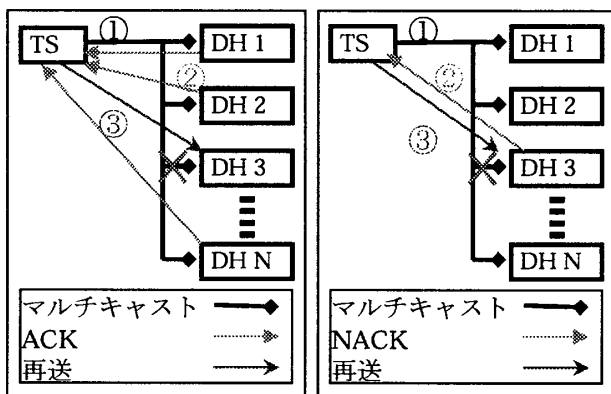


図2. 送信側再送制御方式 図3. 受信側再送制御方式

TRSはDHが増加するにつれ、DHからのACKの総数が増加し、TSへの負荷が増加してしまうという点からスケーラビリティに関して問題があるといえる。一方RRSではDHが増加した場合でも、TSへの負荷の増加が抑えられる。本研究は大規模なシステムを対象としているためスケーラビリティも重要となる。そこで、マルチキャスト通信にRRSを付加した高信頼性マルチキャスト通信システムの開発を行った。

一般的に知られている高信頼性マルチキャストプロトコルのなかではSRM(Scalable Reliable Multicast)が最もRRSに類似している[2]。この両者の共通点は、受信側でシーケンス番号の抜けからパケットの損失を発見しNACKで再送を要求する点である。相違点としてNACKと再送の通信手段が挙げられる。SRMではNACKをマルチキャストで送信し、このNACKにマルチキャストセッションに参加している誰もがマルチキャストで再送できるのに対して、RRSではNACKをユニキャストでTSのみに送信し、TSはユニキャストでNACKを送信してきたDHに再送する。これは、このシステムではマルチキャストメッセージを受信するDHはEBヘイイベントフラグメントデータを送信しなければならないDHに余分な負荷をかけないためである。DHの増加によりTSへのNACKが増加した場合は、TSの台数を増やすことで負荷分散を施す。

5. 性能評価

開発したシステムの性能を評価するため、システムを動作させ性能試験を行った。TS-DH間の通信にマルチキャスト通信にRRSを付加した高信頼性マルチキャスト通信、TCP通信によるユニキャスト通信をそれぞれ用いたシステムを動作させ性能比較を行った。測定はDual Athlon 1.6GHzのCPU、256Mbyteのメモリ、Intel e100のNICを装着し、OSに

Linux(カーネルバージョン 2.4.18)が動作しているPCを使用した。ネットワーク環境はPCI GX-08SXにCISCO Catalyst3500 series XLとREPOTEC RP-G3240Uが接続されている。システムパラメータとして、DHの数を8、16、32、48、58台、イベントフラグメントデータサイズを256、512、1024、1536 Byteとそれぞれ増加させ、Trigger送信レートを100 Hzに固定した。

図4に両システムのTRにおけるTriggerメッセージを送信してから、イベントビルド応答を受信するまでの経過時間を示す。DHの数が16台以下の場合、目立った違いはないが、16台以降ではDHの増加に伴い、両システムとも経過時間が大きくなっているが、RRSのほうがTCPに比べ増加率が抑えられている。これはマルチキャスト通信を用いることによりTS-DH間の重複したメッセージ送信処理を削減できたためである。この結果から1対多の通信にマルチキャスト通信にRRSを付加した高信頼性マルチキャスト通信を用いることにより、システムの効率が向上し、スケーラビリティに関してもよいといえる。

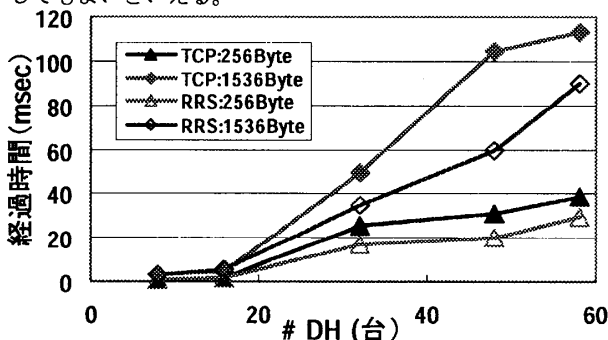


図4. TRにおけるTriggerメッセージを送信してイベントビルド応答を受信するまでの経過時間

6. 再送制御の下位層への移行

これまでの再送制御はアプリケーションレベルでの制御であった。この再送制御をより下位層のトランスポート層で行うことによりさらなる効率化を図った。トランスポート層で再送制御を行うことにより、カーネル空間とユーザ空間との間で無駄なデータのコピーがなくなり、処理が軽減される。そして、メッセージの損失をすばやく発見できるなど、処理の軽減や再送制御の効率化が期待できる。

7. まとめ

大規模データ収集システムにおける1対多の通信に着目し、この通信手段としてマルチキャスト通信を用いたデータ収集システムの開発を行った。その中でマルチキャスト通信の欠点である信頼性の低さを補うため、独自に開発した受信側再送制御方式を付加した高信頼性マルチキャスト通信を開発した。性能試験の結果、この高信頼性マルチキャスト通信を用いることにより、ネットワークに流れる重複メッセージを削除でき、システムの効率が向上することを確認した。さらに再送制御をアプリケーション層から、より下位層のトランスポート層で行うことによりシステムの性能が向上した。

8. 参考文献

- [1] ATLAS Collaboration, "ATLAS High Level Trigger, Data Acquisition and Controls", ATLAS Technical Design Report, ATL-DAQ-2003-052, 2003
- [2] S. Floyd, V. Jacobson, C-G. Liu, S. McCanne, L. Zhang: "A Reliable Multicast Framework for Light-weight Sessions and Application Level Framing", IEEE/ACM Transactions on Networking, Dec. 1997.