

# 高信頼なデータストリーム処理システムにおけるリカバリ時間短縮手法の提案

An Approach to Reducing Recovery Cost for High-Availability Stream Processing Systems

長野 恭子<sup>▽</sup> 糸川 剛<sup>◇</sup>  
北須賀 輝明<sup>▲</sup> 有次 正義<sup>▲</sup>

Kyoko NAGANO Tsuyoshi ITOKAWA  
Teruaki KITASUKA Masayoshi ARITSUGI

近年、継続的に発生するデータの処理を行うストリーム処理システムの需要が高まっており、このようなシステムは分散環境でも構築されている。分散環境においては1台の計算機の故障が、システム全体に影響を及ぼさないように高信頼化の技術を組み込む必要がある。そこで、プライマリの計算機の他にバックアップの計算機を配置する高信頼化手法が盛んに研究されている。高信頼化の既存手法では、故障が発生していない通常処理の場合、バックアップの計算機はアイドル状態となっている。本研究ではバックアップの計算機にも通常の処理を振り分ける手法を提案する。これによりリカバリ時に再処理するタプル数を削減し、リカバリ時間の短縮を実現する。

Recently, a requirement for stream processing systems which deal with data continually generated has been increasing. Such systems are often developed in distributed environments. Since the failure of a computer may cause damages to all of the systems in distributed environments, high-availability mechanisms must be involved in distributed stream processing systems. One of the high-availability mechanisms is to provide a backup computer for a primary computer. When a primary computer fails, the backup computer takes over the operation of the primary computer. It is important to reduce the cost for the recovery. In this paper, we propose an approach to realizing short recovery from fails by making use of an idle backup computer.

<sup>▽</sup>学生会員 熊本大学大学院自然科学研究科博士前期課程

[kyoko@dbms.cs.kumamoto-u.ac.jp](mailto:kyoko@dbms.cs.kumamoto-u.ac.jp)

<sup>◇</sup>正会員 熊本大学大学院自然科学研究科

[itokawa@cs.kumamoto-u.ac.jp](mailto:itokawa@cs.kumamoto-u.ac.jp)

<sup>▲</sup>正会員 熊本大学大学院自然科学研究科

[kitasuka@cs.kumamoto-u.ac.jp](mailto:kitasuka@cs.kumamoto-u.ac.jp)

<sup>▲</sup>正会員 熊本大学大学院自然科学研究科

[aritsugi@cs.kumamoto-u.ac.jp](mailto:aritsugi@cs.kumamoto-u.ac.jp)

## 1. はじめに

近年、センサを用いた交通状況の監視システムや、株価のモニタリングといった金融アプリケーションなど、リアルタイムで処理を行うシステムが増加している。このように継続的に発生するストリームデータの処理を行うシステムとして、ストリーム処理システムに対する需要が高まっている。また、データソースと処理部が地理的に離れた場所にあるこれらのアプリケーションを、分散環境で実現する分散ストリーム処理システムの研究が行われている[1, 2]。

分散環境においては、1台の計算機で故障が発生しても全体へ影響を及ぼすことなく、システムが稼働し続けなければならないため、高信頼化の技術を組み込む必要がある。そこで主に処理を行う計算機と、故障が発生した場合に処理を引き継ぐバックアップの計算機の2台を用意するという高信頼化手法の研究が行われている[3, 4, 5]。

Hwangら[3]は高信頼化のアルゴリズムとして、Passive Standby, Upstream backup, Active Standbyの3手法を示している。これらの中で、Upstream backupに注目する。

Upstream backupは高信頼化のために生じるオーバヘッドがほとんど無いという利点があるが、障害回復のためのリカバリ時間が長いという欠点がある。また、故障が発生していない間、バックアップのために用意してある計算機はアイドル状態で処理を行っていない。

そこで本研究では、このアイドル状態であるバックアップの計算機にも処理を振り分ける手法を提案する。2台の計算機を用いることで、故障が発生した後のリカバリ処理において、再処理するべきタプルの数を削減することができ、リカバリ時間が短縮できることを示す。

本稿の構成は以下の通りである。まず2章で分散ストリーム処理システムに関連する研究を述べて、高信頼化手法の従来研究を詳しく説明する。3章では、提案手法を説明し、4章で提案手法の有効性を検討するため評価実験を行い、考察する。5章でまとめと今後の課題について示す。

## 2. 関連研究

Brandeis/Brown/MITの合同プロジェクトであるAurora[6]や、UC BerkeleyによるTelegraphCQ[7]などの分散ストリーム処理システムが開発されており、これらに関連した研究が幅広く行われている。

分散ストリーム処理システムは、データの処理部の各演算を複数の計算機へと割り当てている。そこで、各演算をどの計算機に配置するかを決定するためのアルゴリズムが必要となる。RepantisとKalogeraki[8]は、各計算機間の物理的な距離によって発生する通信の遅延と、各計算機の処理能力を考慮して処理に必要な演算を計算機へと配置するアルゴリズムを提案している。また、分散ストリーム処理システムにおける高信頼化の技術に関する研究も盛んに行われており、いくつかの手法が提案されている[3, 4, 5]。これらは、分散ストリーム処理システムを構成する各計算機とは別に、バックアップ用の計算機を配置することで、ある計算機の故障時にバックアップが処理を引き継ぐという共通した方式を持っている。

これらの関連研究の中でも、1つの計算機の故障がシステム全体に影響を及ぼさないための高信頼化技術は、分散環境において重要である。そこで本研究では、この高信頼化手法に焦点を当てる。

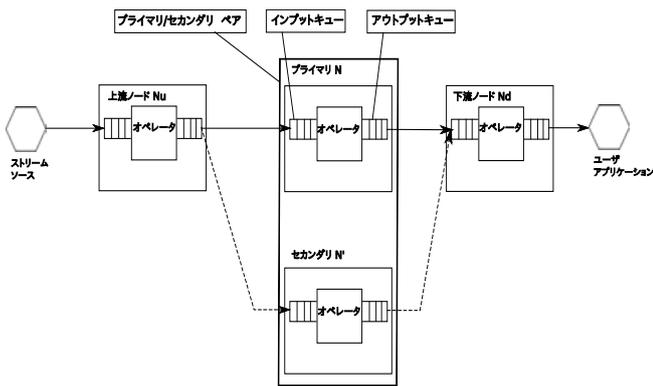


図1 分散ストリーム処理システムの例[3]

Fig.1 A high availability distributed stream processing system[3]

## 2.1 従来の高信頼化手法

Hwangら[3]によって、Passive Standby, Upstream backup, Active Standbyという3つの高信頼化手法が示されている。ここではまず、本研究で扱う分散ストリーム処理システムのモデルを示し、高信頼化のための手法として提案された3手法について詳しく説明する。

### 2.1.1 システムモデル

図1に分散ストリーム処理システムの例を示す。システム上を流れるストリームデータは、リアルタイムに処理されるタプルのシーケンスである。システムの演算処理が上流ノード  $N_u$ 、ノード  $N$ 、下流ノード  $N_d$  の3つのノードに分散している。各ノードは入力と出力のキューを持っている。また、ストリームデータの流れる方向は実線の矢印で描かれており、上流ノード  $N_u$  からノード  $N$  へとタプルが流れる。ノード  $N$  でタプルが処理され、下流ノード  $N_d$  へと送信される。

このようなシステムにおいて、処理を行うノード  $N$  をプライマリと言い、プライマリに接続されるバックアップのためのノード  $N'$  をセカンダリと言う。セカンダリは、プライマリで故障が発生したことを検知し、処理を引き継ぐ役割を持つ。故障時のタプルの流れは破線の矢印で表されている。

故障を発見するために、セカンダリは定期的に生存確認のメッセージをプライマリへ送信する。それを受信したプライマリは応答を返す。故障が発生した場合、プライマリは応答を返さない。セカンダリで生存確認がタイムアウトになったら、セカンダリはプライマリで故障が発生したと判断して上流ノードにリクエストを送信し、セカンダリが処理を引き継ぐ。なお、生存確認のメッセージを送る間隔を生存確認送信間隔と定義する。システムの基本的な処理の流れはこのようになっているが、プライマリとセカンダリの利用方法により高信頼化手法は3つに分類されている。これらの高信頼化手法について以下に説明する。

### 2.1.2 Passive Standby

Passive Standbyは、プライマリ  $N$  が完全にタプルを処理するまで上流ノード  $N_u$  の出力キュー内に、タプルを保存しておく方式である。プライマリはタプルの処理を行うと同

時に、その内部状態の変化を知らせるメッセージをセカンダリ  $N'$  へ送信する。この時の内部状態をチェックポイントと言い、このメッセージをチェックポイントメッセージと言う。セカンダリはチェックポイントメッセージを受信すると、上流ノードへのタプルまでの処理が完了したかを知らせ、上流ノードは出力キューから該当タプルを削除する。

プライマリが故障した場合、セカンダリが上流ノードにプライマリの故障を通知することで、上流ノードは出力キュー内にあるタプルをセカンダリへ再送信する。セカンダリは再送信されたタプルの処理を、最新のチェックポイントから行う。その後、プライマリの処理を引き継ぐ。

### 2.1.2 Upstream Backup

Upstream backupも上流ノードの出力キュー内に、タプルを保存する方式である。しかし、プライマリは内部状態をセカンダリへ送信せず、上流ノードの出力キューからタプルを削除するために、Level-0 ACKとLevel-1 ACKの2つのメッセージを用いる。

プライマリ  $M$  に上流ノード  $N_u$  からタプル  $t_1$  が届き、それを処理して下流ノード  $N_d$  にタプル  $t'1$  を出力した状況を考える。下流ノードは、プライマリから受信したタプル  $t'1$  の識別情報を付加したLevel-0 ACK をプライマリに送信する。なお、Level-0 ACKを送信する間隔をACK送信間隔と定義し、これは一定とする。Level-0 ACKを受け取ったプライマリは、その中の識別情報を参照してタプル  $t'1$  を生成したタプル  $t_1$  を見つけ、 $t_1$  の識別情報を付加したLevel-1 ACKを上流ノードへ送信する。このLevel-1 ACKを受信すると、上流ノードはタプル  $t_1$  を出力キューから削除する。

プライマリで故障が発生した場合、セカンダリは上流ノードにプライマリの故障を通知し、上流ノードは出力キュー内に保存されていたタプルをセカンダリへ再送信する。セカンダリはそれらを受信し、タプルの再処理を行い処理を引き継ぐ。

### 2.1.4 Active Standby

Active Standbyはプライマリ  $N$  とセカンダリ  $N'$  の両方が上流ノードからタプルを受信し、並列に処理を行う。プライマリで故障が発生していない場合はセカンダリもタプルの処理を行うが、下流ノードに出力せず、セカンダリの出力キュー内にタプルを保存しておく。Upstream backupと同様に、プライマリに上流ノード  $N_u$  からタプル  $t_1$  が届き、それを処理して下流ノード  $N_d$  にタプル  $t'1$  を出力した状況を考える。下流ノードはプライマリから受信したタプル  $t'1$  の識別情報を付加したLevel-0 ACKをプライマリに送信する。また、セカンダリではプライマリと同時に処理が行われ、タプル  $t'1$  が出力キューに保存されている。Level-0 ACKを受信したプライマリは、タプル  $t'1$  が下流ノードまで到達したことをセカンダリへ知らせる。すると、セカンダリは出力キューからタプル  $t'1$  を削除する。

プライマリで故障が発生した場合、セカンダリは出力キュー内のデータを下流ノードへ送信し、処理を引き継ぐ。

また、Active StandbyをWAN環境に応用した手法も提案されている[4]。この場合、セカンダリも下流ノードに処理を行ったタプルを送信する。下流ノードはプライマリとセカンダリからタプルを受信することとなるが、下流ノードでは最も速く到達した方のタプルを用いて処理を継続するという方式をとる。セカンダリも下流ノードにタプルを送

信しているため、出力キューから古いタブルを削除する必要がなく、Level-0 ACKは不要となる。計算機の故障やネットワークの障害がより頻繁に起こると想定される環境においては、ストリームデータだけでなくLevel-0 ACKも該当するノードに到達しない可能性が大きい。したがってWAN環境に適した手法と言える。

2.2 従来手法の特性と問題点

Hwangら[3]は3手法の特性をバンド幅オーバーヘッドとリカバリ時間という2つの指標を用いて評価している。全データ転送に用いたバンド幅のうち、高信頼化のために用いたバンド幅をバンド幅オーバーヘッドという。また、データの処理中に故障が発生した場合、データの損失が生じる。故障を検知してから損失したデータを再処理し、故障前の状態へ復旧するまでの時間をリカバリ時間という。

Active Standbyは、プライマリとセカンダリの両方にタブルを送信するため、バンド幅オーバーヘッドはほぼ100%となる。一方、セカンダリにおいて並列にタブルを処理しているため、リカバリ時間は非常に短い。Upstream backupでは、故障時のみセカンダリへデータを送信するため、バンド幅オーバーヘッドはほとんどない。一方で、上流ノードからタブルの再送信を行い、セカンダリで再処理を行うためリカバリ時間は長くなる。Passive Standbyはセカンダリがプライマリの内部状態を保持しており、故障時にはその状態から再処理を行う。したがってリカバリ時間はUpstream backupよりも短くなる。しかしプライマリがセカンダリへ内部状態を送信する際には多くのバンド幅を消費してしまうため、そのオーバーヘッドは大きくなる。

上記のことから、Upstream backupは3手法の中ではバンド幅オーバーヘッドが最も少ないという利点をもっている。しかし、故障が発生した際に、リカバリ時間が長いという問題が挙げられる。上流ノードの出力キュー内にあるタブルの数が多ければ、その分セカンダリで再処理するため、リカバリ時間はこのタブル数に比例して増加する。

ここで、故障が発生していない場合の各手法のセカンダリの動作について注目する。Active Standbyでは上流ノードから常にストリームデータが送信され、セカンダリは待機している。Passive Standbyでは、セカンダリはプライマリの内部状態を知らせるメッセージを受信しており、メッセージを受信した時点までで処理が完了しているタブルを、上流ノードへ知らせるという働きを持っている。すなわち、これらの2手法は故障時でない場合もセカンダリは動作している。一方、Upstream backupでは故障が発生していない間、セカンダリはプライマリに対して生存確認を行うのみでタブルの処理を行わないので、アイドル状態である。

ここで、Upstream backupに注目する。Upstream backupはバンド幅オーバーヘッドが少ないため、システムを構築する際、高信頼化のために用いるバンド幅を考慮する必要がないという利点がある。さらにセカンダリは故障が発生していない間、待機しているのみである。

3. 提案手法

高信頼化のための従来手法であるUpstream backupは、リカバリ時間が長いという欠点がある。また、システムにはプライマリとセカンダリの計算機を配置するが、通常はプライマリだけで処理を行い、故障時のみセカンダリを利用するという特徴を持つ。そこで、我々はUpstream backup

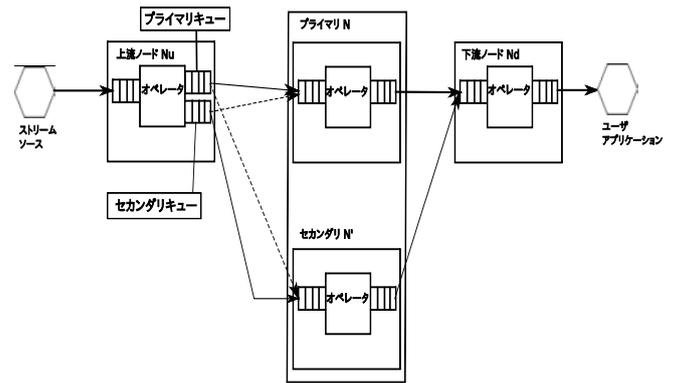


図2 提案手法のシステム構成図

Fig.1 A construction of our proposed system

に注目する。通常はアイドル状態であるセカンダリにも処理を行わせることで、故障発生時に再処理を行うタブルの数を削減し、Upstream backupのリカバリ時間を短縮する手法を提案する。本章では提案手法を説明し、その動作例を示す。

3.1 動作

提案手法を、通常時の処理の流れ、生存確認の方法、上流ノードが保持する出力キュー内のタブルの削除と再送信に関して順に説明する。また、本手法のシステムの構成図を図2に示す。

処理の流れ

提案手法においては、上流ノードからプライマリとセカンダリへ切り替えてストリームデータを送信する。図2ではこのデータの流を実線の矢印で示す。上流ノードは送信したタブルを保存する出力キューとして、プライマリ用とセカンダリ用の2つのキューを持つ。そして、プライマリに送信したタブルはプライマリ用の出力キューに、セカンダリに送信したタブルはセカンダリ用の出力キューに保存する。プライマリもしくはセカンダリは、上流ノードからのデータを受信すると処理を行い、出力タブルを下流ノードへと送信する。下流ノードは、プライマリかセカンダリのどちらからデータを受信したか判断して、該当するノードへLevel-0 ACKを返す。このとき、Level-0 ACKを返すACK送信間隔はUpstream backupと同様に一定間隔とする。Level-0 ACKを受信したプライマリもしくはセカンダリは、出力が下流ノードで受信され処理が完了したことを上流ノードに知らせるために、Level-1 ACKを上流ノードへ送信する。上流ノードは、プライマリ、セカンダリのどちらからLevel-1 ACKを受信したか判断して、その出力キューから該当するタブルを削除する。

プライマリもしくはセカンダリにおいて故障が発生した場合、後述の生存確認により故障を検知し、生存しているノードは上流ノードに故障発生を知らせるメッセージを送信する。これにより、上流ノードは故障したノードの出力キュー内のタブルを生存しているノードへと送信し、再処理を行うことによって、リカバリを実現する。

生存確認方法

従来の手法においては、セカンダリからプライマリへと一方向に定期的にパケットを送信し、生存確認を行っていた。

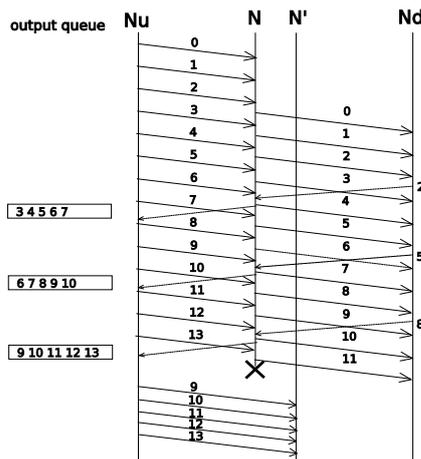


図3 Upstream backup のデータの流れ

Fig.1 Stream data flows in Upstream backup

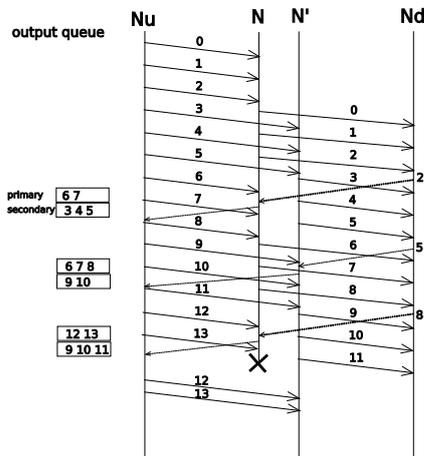


図4 提案手法のデータの流れ

Fig.1 Stream data flows in proposed system

セカンダリがパケットを送信し、それを受け取ったプライマリは応答を返す。一定期間応答が返らずにタイムアウトとなった場合、プライマリに故障が発生したと判断する。

本手法では、プライマリとセカンダリの両方を処理に用いるため、双方向で定期的に生存確認を行う。従来のセカンダリからプライマリへの生存確認に加え、プライマリからセカンダリへの生存確認も行い、どちらのノードで故障が発生した場合でも検知できるようにする。

**キュー内のデータ再送信**

故障が発生し、生存しているノードがそれを検知すると、上流ノードにタブルの再送信のリクエストを送信する。その時、上流ノードは出力キュー内のタブルを再送信しなければならぬ。本手法では、タブルを保存しておくための出力キューをプライマリ用のものと、セカンダリ用のものの2つに分割している。そのため、故障が発生した側のキュー内のタブルを再送信する仕組みが必要となる。プライマリもしくはセカンダリが故障を検知して、再送信のリクエストを送信する。すると上流ノードは、どちらかのノード

で故障が発生したことを識別し、故障したノードの出力キュー内にあるタブルのみ再送信を行う。図2ではこのリカバリ処理時のタブルの流れを破線の矢印で示す。

**3.2 動作例**

図3に従来手法であるUpstream backup, 図4に提案手法それぞれにおけるストリームデータの流れと、上流ノードの出力キュー内のタブルを示す。各タブルの送受信は実線の矢印で表しており、送受信されたタブルのシーケンス番号が付けられている。また、点線の矢印はLevel-0 ACK, Level-1 ACKの各ACKと該当するタブルのシーケンス番号を表している。左側には上流ノードのもつ出力キューを示している。

従来手法では、プライマリに故障が発生したとき、上流ノードのキュー内には(9, 10, 11, 12, 13)の5つのタブルがあり、これらを全てセカンダリに送信しなおして再処理を行う。提案手法ではプライマリにおいて故障が発生したとき、セカンダリ側に送信していたデータ(9, 10, 11)はセカンダリで正常に処理が行われている。そのため再送信する必要がなく、プライマリ側のキュー内にあるデータ(12, 13)のみを再送信する。このように、再送信しなければならないタブルの数が減るので、リカバリ時間の短縮につながると考えられる。

リカバリ時に、故障が発生していないノードへ送信したタブルは再送信しないため、下流ノードが受信した出力結果のタブルの順序が保たれない場合が生じる。演算の種類によってはタブルの順序の変更にはセンシティブなものがあるが、その場合、必要に応じて下流ノードにおいてタブルのソートを行わなければならない。したがって、ソート処理が必要な演算では、下流ノードにおけるコストが増えてしまうことが考えられる。

**4. 実験**

本章では、提案手法の有効性を検証するために行うシミュレーションについて述べ、実験結果について考察する。

**4.1 シミュレーションの概要**

シミュレーション環境として、図1, 2に示すような上流ノード、プライマリ、セカンダリ、下流ノードで構成されるシステムを構築する。これらの計算機の処理能力は全て一様とする。また、各ノード間の通信にはTCP/IPプロトコルを用いる。このシステム内を流れるタブルは、ID部とデータ部で構成される。IDはタブルが生成された順番を示し、データ部は擬似的な数値データとする。また、演算は入力に対して特定の値のみを出力するフィルタ演算を行うと仮定する。フィルタ演算はタブルの順序が変更されても出力結果には影響しないため、下流ノードにおいてソートを行う必要がない。なお、今回は簡略化のため選択率100%のフィルタ演算を用いることとする。提案手法では、上流ノードはプライマリとセカンダリへ交互に切り替えて、ストリームデータを送信する。本実験では、この切り替えを100msの一定間隔で行った。

次にシミュレーションの流れを説明する。システムが定常状態となるように、一定時間ストリームデータを送信し続ける。プライマリもしくはセカンダリの動作を止めることで故障を発生させ、もう一方が故障を知らせるメッセージを上流ノードに送信すると、リカバリ処理を開始する。ネットワークシミュレータns-3[9]を用いてこの環境を構築

表1 シミュレーションに用いたパラメータ  
Table 1 Parameter settings

ACK 送信間隔	25ms
生存確認送信間隔	100ms
タプルサイズ(ID 部, データ部)	(8,50)byte
ネットワークバンド幅	16Mbps
ネットワーク遅延	5ms
演算時間	10 $\mu$ s/tuple
プライマリ/セカンダリ切り替え間隔	100ms

し、システムの振舞いを測定する。

### 4.2 検証項目

従来手法との特性比較に用いる指標として、リカバリ時に再送信されるタプルの数とリカバリ時間を測定する。リカバリ時間は、上流ノードが再送信リクエストのメッセージを受信した時刻から、生存しているノードでタプルの再処理が完了するまでの時間とする。

その際に、パラメータとしてストリームデータの入力レートを変化させる。ACK送信間隔が一定であれば、入力レートの増加によって、上流ノードの出力キュー内に保存されるタプル数も増加する。この出力キュー内のタプル数と、リカバリ時間を測定することで本手法の有効性を検証する。シミュレーションに用いた計算機はCPU Intel (R) Core(TM)2 Duo 3.00GHz, メモリ 4GB, OS Ubuntu 9.04 である。また各パラメータを表1に示す。なお、1回のシミュレーションの時間を1分間とする。システムが定常状態となる時間を30秒間とし、これ以降のランダムな時刻に故障が発生するものと仮定する。このシミュレーションを10回試行した平均値を実験結果とする。

### 4.3 実験結果

実験結果を図5と6に示す。図5は入力レートと故障時に再処理するタプル数の関係を示している。これより、従来手法に比べ、提案手法では再処理しなければならないタプル数が減っていることがわかる。しかし、入力レートが低い間、このタプル数は従来手法と提案手法での差がほとんどない。入力レートが低い時は、新たにタプルが送信されるまでに上流ノードはLevel-1 ACKを受信し、出力キューからタプルが削除されていることが原因であると考えられる。したがって、入力レートが高いほど、リカバリ時に再処理するタプル数を削減できることがわかる。また図6は、入力レートとリカバリ時間の関係を示している。これよりリカバリ時間も短縮されていることがわかる。

入力レートが1000[tuples/s]以下では図5において提案手法と従来手法の差は少ないが、図6のリカバリ時間では2つの手法で大きく差がでている。この原因は、システムの実装時に用いたTCPプロトコルが挙げられる。RFC896[10]によると、TCP/IPプロトコルには、送信側がバッファにデータを一旦取り込んでまとめて送信を行うNagleアルゴリズムが使われている。従来手法では通常処理の際にセカンダリが用いられておらず、リカバリ処理を行う際に初めて利用される。つまり、これまで通信が行われていない相手と

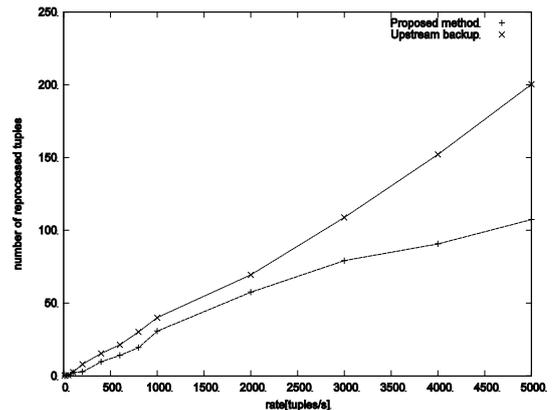


図5 再処理するタプル数

Fig.5 Number of reprocessed tuples

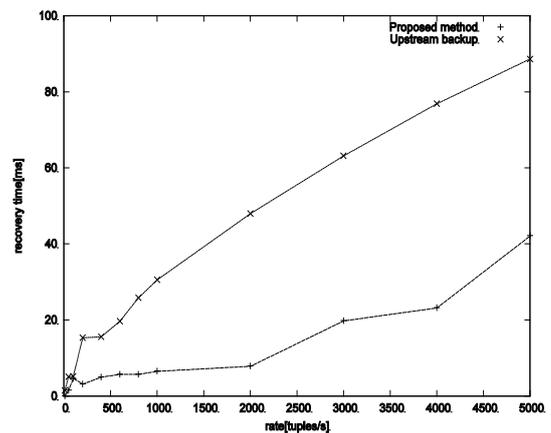


図5 リカバリ時間

Fig.5 Recovery times

通信が開始されるため、Nagleアルゴリズムが働き、上流ノードで送信バッファにデータを一旦取り込む間で遅延が生じる。一方提案手法では、セカンダリとの通信が既に行われているためこのアルゴリズムが働かない。したがってリカバリ時間にはこのような差が生じたと考えられる。

シミュレーションを行った結果、通常処理をプライマリだけでなくセカンダリにも振り分け、タプルを保存するキューを分割することで、再処理するタプル数の削減が可能であると確認できた。また、システムの実装上ではTCP/IPの仕様によってリカバリ時間に影響が出ることがわかった。結果として、提案手法によってリカバリ時間の短縮が可能であると言える。

### 5. おわりに

本研究では、分散ストリーム処理システムの高信頼化手法であるUpstream backupに対してリカバリ時間の短縮のため、通常処理時にアイドル状態となっているセカンダリノードにも処理を振り分ける手法を提案した。そして上流ノード、プライマリ、セカンダリ、下流ノードから構成されるシステムでシミュレーションを行った。これにより、実際にリカバリ時に再送信するタプルの数を削減することが

でき、リカバリ時間の短縮が可能となることがわかった。

今後の課題として、本手法について今回はリカバリ時間のみを議論したが、プライマリとセカンダリに処理を振り分けた場合の、システムのパフォーマンスに関する評価を行う必要があると考えられる。提案手法では、上流ノードにおいてはキューを分割して保存する必要がある、下流ノードにおいては必要に応じてソート処理をしなければならない。これらがパフォーマンスに与える影響を調査する必要がある。演算に関しても、今回は簡略化のためフィルタ演算のみで実験を行った。しかし、より多くの種類の演算で実験を行う必要がある。その際に Join や Aggregate 演算では、ソートを行うことによるペナルティが考えられる。さらに、これらの演算では一定のウィンドウサイズに区切り演算を行う。その特性を考慮して、プライマリとセカンダリへの処理の振り分けを行わなければならない。またシミュレーションでは、プライマリとセカンダリへ送信するストリームデータの切り替え間隔は一定として実験を行った。しかし、処理能力の異なる2つの計算機をプライマリとセカンダリに用いた場合を想定すると、計算機のリソースを考慮した上で処理の振り分けを行う仕組みが必要となる。

実際に分散ストリーム処理システムのインフラ構築を行うことを考えると、リカバリ時間、パフォーマンス、消費電力等の指標を示すことによって各計算機の性能やコストを考慮して、アプリケーションの仕様に適したシステムを設計することが可能になると考えられる。

## [文献]

- [1] Brian Babcock, Shivnath Babu, Mayur Datar, Rajeev Motwani, and Jennifer Widom. Models and issues in data stream systems. In PODS, pp. 1-16 (2002).
- [2] Rajeev Motwani, Jennifer Widom, Arvind Arasu, Brian Babcock, Shivnath Babu, Mayur Datar, Gurmeet Singh Manku, Chris Olston, Justin Rosenstein, and Rohit Varma. Query processing, approximation, and resource management in a data stream management system. In CIDR (2003). Available at <http://www.cidrdb.org/cidr2003>.
- [3] Jeong-Hyon Hwang, Magdalena Balazinska, Alex Rasin, Ugur Cetintemel, Michael Stonebraker, and Stanley B. Zdonik. High-availability algorithms for distributed stream processing. In ICDE, pp. 779-790 (2005).
- [4] Jeong-Hyon Hwang, Sanghoon Cha, Ugur Cetintemel, and Stanley B. Zdonik. Borealis-r: a replication-transparent stream processing system for wide-area monitoring applications. In SIGMOD Conference, pp. 1303-1306 (2008).
- [5] 塩川浩昭, 渡辺陽介, 北川博之, 川島英之. 分散ストリーム処理システムにおける高信頼化手法の提案. In DEIM Forum (2009).
- [6] Daniel J. Abadi, Donald Carney, Ugur Cetintemel, Mitch Cherniack, Christian Convey, Sangdon Lee, Michael Stonebraker, Nesime Tatbul, and Stanley B. Zdonik. Aurora: a new model and architecture for datastream management. VLDB J., Vol. 12, No. 2, pp. 120-139 (2003).
- [7] Sirish Chandrasekaran, Owen Cooper, Amol Deshpande, Michael J. Franklin, Joseph M. Hellerstein, WeiHong, Sailesh Krishnamurthy, Samuel Madden, Vijayshankar Raman, Frederick Reiss, and Mehul A. Shah. TelegraphCQ: Continuous dataow processing for an uncertain world. In CIDR (2003).
- [8] Thomas Repantis and Vana Kalogeraki. Replica placement for high availability in distributed stream processing systems. In DEBS, pp. 181-192 (2008).
- [9] The ns-3 network simulator. <http://www.nsnam.org/>.
- [10] RFC896 Congestion control in IP/TCP internetworks. <http://tools.ietf.org/html/rfc896>.

## 長野 恭子 Kyoko NAGANO

2009 熊本大学工学部数理情報システム工学科卒業。同年から同大学院自然科学研究科情報電気電子工学専攻博士前期課程に在学。ストリーム処理システムの研究に従事。

## 糸川 剛 Tsuyoshi ITOKAWA

熊本大学大学院自然科学研究科情報電気電子工学専攻助教。1997 熊本大学大学院自然科学研究科修士課程修了, 2001 熊本大学大学院博士課程修了。博士(工学)。研究の中心はアルゴリズムとデータ構造の設計と解析。

## 北須賀 輝明 Teruaki KITASUKA

熊本大学大学院自然科学研究科准教授。1993 京都大学工学部情報工学科卒。1995 奈良先端科学技術大学院大学情報科学研究科博士前期課程了。同年シャープ(株)入社。2001 九州大学大学院システム情報科学研究院助手。2006 九州大学博士(工学)。2007年10月より現職。モバイルコンピューティング, 組込みシステム, 並列/分散処理, 計算機アーキテクチャに研究的関心を持つ。

## 有次 正義 Masayoshi ARITSUGI

熊本大学大学院自然科学研究科教授。1991 九州大学工学部情報工学科卒。1996 同大学院了。博士(工学)。同年群馬大学助手。同助教授を経て, 2007 より現職。データベースシステム, 分散並列データ処理等に興味を持つ。