# 複数の一貫性レベルを保証可能な バックエンドベースデータレプリケーション

受付日 2015年6月29日, 採録日 2015年12月7日

概要:システムの可用性、信頼性を実現するためにデータベースの複製(レプリカ)を作成する技術であるレプリケーションが広く利用されている。しかし、アプリケーションごとに要求される最低限の一貫性を保証するレプリケーションプロトコルを個別に実装することは、管理や運用の面でユーザの負担となる。したがって、複数の一貫性レベルを保証可能なレプリケーションプロトコルが求められる。既存研究である Multi-Consistency Data Replication(McRep)は、ミドルウェアベースのレプリケーションであり、かつ複数の一貫性レベルが保証可能である。しかし、クライアント数の増加につれ、ミドルウェア部に配置されたレプリケーション制御を行うサーバ(レプリケータ)が性能のボトルネックとなるという問題がある。本研究では、McRep と同様のレプリケーション制御を実現しつつ、この問題を解決する手法を提案する。具体的には、レプリケータをレプリカのみと通信を行うバックエンド部に配置し、各レプリカもRead-only トランザクションの一貫性制御を行う。評価から、バックエンド部で制御を行うことでレプリケータがボトルネックとなることを回避し、提案手法では Read-heavy なワークロードにおいて、一貫性レベルが One-Copy Serializability の場合、最大スループットが McRep の約 2 倍向上することを確認した。

キーワード:データベースレプリケーション、一貫性レベル、バックエンドベース、ミドルウェア

# Back-end Based Data Replication Supporting Multiple Consistency Models

Atsushi Ohta<sup>1,a)</sup> Masaya Matsuno<sup>1,b)</sup> Ryota Kawashima<sup>1,c)</sup> Hiroshi Matsuo<sup>1,d)</sup>

Received: June 29, 2015, Accepted: December 7, 2015

**Abstract:** Replication is widely used in parallel and distributed systems for reliability and availability. On the other hand, developers have to consider minimum consistency requirement for each application. Therefore, novel replication protocol that ensure multiple consistency models is required. Multi-Consistency Data Replication (McRep) is a middleware-based replication protocol and can support multiple consistency models. However, McRep has a potential problem that a replicator (a server controlling replications) acting as a middleware can be a performance bottleneck. We propose a backend based replication protocol to solve this problem, while ensuring same consistency models. More precisely, we place the replicator on the backend area where the replicator communicates with only replica servers, and extend the replica's role to control the consistency for read-only transactions. We implemented and evaluated both the proposal protocol and the McRep. The results showed that our protocol improve up to approximately twice throughput of transactions at a read-heavy workload in one-copy serializabilty as McRep's.

Keywords: replication, consistency, back-end, middleware

#### 1 名古屋工業大学

Nagoya Institute of Technology, Nagoya, Aichi 466–8555, Japan

- a-ohta@matlab.nitech.ac.jp
- b) matsuno@matlab.nitech.ac.jp
- c) kawa1983@nitech.ac.jp
- d) matsuo@nitech.ac.jp

## 1. はじめに

近年では Google や Amazon, Facebook に代表されるように多様な Web サービスが登場し、広く普及している. これにともなって、サービス利用者も増加し続け、システムも大規模化している.このような状況において、データ

ベース層における性能向上が重要な課題となっている.し かし、データベースサーバ単体の性能向上には限界がある うえ、単一障害点となるという問題がある. そこで、ある データベースの複製を別サーバ上にリアルタイムに作成す る技術であるレプリケーションが広く用いられている. 以 降ではデータベースの複製が存在するサーバをレプリカと 呼ぶ. 複数のレプリカで同一データを管理することで, 負 荷分散が可能となり, 高可用性, 耐故障性を実現すること ができる. レプリケーションを実現する代表的な方式には、 マスタスレーブ方式 [1] とマルチマスタ方式 [2] が存在す る. しかし、マスタスレーブ方式ではマスタレプリカに負 荷が集中する場合があり、マルチマスタ方式では一貫性制 御が複雑になり、性能向上が困難になってしまうという問 題がある. これらに対して、クライアントとレプリカの間 に専用のミドルウェア (以降,レプリケータと呼ぶ)を配置 し、レプリケータがレプリケーションの制御を行うミドル ウェア方式のレプリケーションが研究されている[3],[4]. この方式では、レプリケータがレプリケーション制御を行 うことで特定のレプリカに負荷が集中することを回避し. 一貫性制御も容易にすることができる.しかし、全リクエ ストがレプリケータを経由するために、レプリケータがシ ステム全体の性能ボトルネックとなる場合がある.

また一方で、いずれの方式においても、レプリケーショ ン時にはレプリカ間におけるデータの一貫性をどの程度 保証するべきかという共通の課題が存在する.一般的に は、より強い一貫性を保証するほど性能が犠牲になるた め、アプリケーションが要求する最低限の一貫性を保証す ることが望ましい. そのため、レプリケーション時にレプ リカ間でデータの一貫性をどの程度保証するかによって, 線形化可能性 (Linearizability) や順序一貫性 (Sequential Consistency) など様々な一貫性レベルが提案されている. しかし、特定の一貫性レベルしか保証できない場合は利 用可能なアプリケーションが限られ,汎用性に欠けてし まう. したがって、複数の一貫性レベルを保証可能なレプ リケーションプロトコルが望ましい. Multi-Consistency Data Replication [5] (McRep) は、複数の一貫性レベルを 保証可能なレプリケーション手法だが、ミドルウェア方式 を採用しているため、先述した性能ボトルネックの問題を かかえている.

本研究では、既存手法である McRep を拡張し、ミドルウェア部で行っていたレプリケーション制御を、バックエンド部で実現することでレプリケータがボトルネックとなることを回避し、よりスケールアウトしやすいレプリケーション手法を提案する。本研究で規定するバックエンド部とはクライアントからのリクエストを中継せず、レプリカのみと通信するサーバが配置される領域を指す。提案手法では、レプリケータがバックエンド部で制御を行うことに加え、各レプリカが Read-only トランザクションに関して一貫性制

御を行う.これにより、レプリケータが全リクエストを処理する必要がなくなるため、既存手法と比較してレプリケータの処理量を軽減することができる.評価から、McRepと比較して、提案手法では Read-heavy なワークロードにおいて、一貫性レベルが One-Copy Serializability の場合、レプリケータがボトルネックとなることを回避し、最大スループットが McRep の約 2 倍向上することを確認した.

本論文の構成は以下のとおりである。まず、2章でレプリケーションにおける課題であるレプリカ間の一貫性について説明する。3章では、既存手法である McRep について紹介し、4章で本研究での提案について述べる。そして、5章で評価および考察を行い、6章で関連研究を示し、最後に7章で本研究のまとめと今後の課題を述べる。

## 2. レプリケーションにおける一貫性

レプリケーションを行う際には、レプリカ間の一貫性を どの程度保証するべきかという点も考慮する必要がある。 一般的には、より強い一貫性を保証するほど性能が低下す るため、アプリケーションにとって最低限必要な一貫性を 保証することが望ましい。以下に代表的な一貫性レベルに ついて説明する。

- 線形化可能性 (Linearizability) 線形化可能性 [6], [7] では、実行される一連のトランザクションが実時間で順序付けられる必要があるが、このとき、並行実行されたトランザクション間の順序については実時間順になる必要はなく、レプリカ間で同一の順序となることを保証すればよい.
- 順序一貫性(Sequential Consistency) 順序一貫性 [8] では、同一クライアントが実行した一連のトランザクションは実時間で順序付けられる必要がある。したがって、各クライアントからは一貫性レベルが線形化可能性として観測される。
- 直列化可能性 (One-Copy Serializability) 直列化可能性 [9] では、実行されたすべてのトランザクションが全レプリカにおいて同一順序で逐次実行された結果と同じであることのみを保証すればよく、一連のトランザクションの処理順序に関して制約はない.したがって、各クライアントは古いデータを読み出す可能性があるため、一貫性のない状態を観測する場合がある.

上記以外にも様々な一貫性レベルが提案されているが、一般にはアプリケーションごとに適切な一貫性レベルは異なるため、複数の一貫性レベルを保証可能な汎用性のあるレプリケーションプロトコルが望ましいといえる.

## 3. ミドルウェアベースによる複数一貫性制御 手法とその問題点

本章では、複数の一貫性レベルを保証可能なレプリケー

ション手法である McRep について紹介し、その問題点を 指摘する.

### 3.1 McRep の概要

McRep は、ミドルウェア方式の非同期レプリケーション であり、6つの一貫性レベル (Linearizability, Sequential Consistency, One-Copy Serializability, Session Snapshot Isolation [10], Generalized Snapshot Isolation [1], Causal Consistency [11]) を保証可能とする手法である. McRep の基本的な構造はミドルウェアベースのプロトコルである SRCA [12] に基づく. McRep では、図 1 に示すように、各 クライアントからのリクエストを仲介するレプリケータに おいて、トランザクションの更新結果(トランザクション 実行時に書き込まれたデータ項目)を格納するキューを保 持し、各更新結果には順にバージョン番号を付与する. そ して、このキュー内の特定のエントリを示す4種類のバー ジョン番号を用いて複数の一貫性レベルを制御する. GVN (Global Version Number) はレプリカ全体での最新のバー ジョンを示し、SVN (Session Version Number) は各クラ イアントごとの観測済みのバージョンを示す. また, レプ リカのバージョン (Replica Version Number, RVN) の上 限値および下限値である RVNhi, RVNlo はレプリカにおい てキュー内の更新結果がどこまで反映されているかを管理 するために用いる. このように、McRepではレプリケータ において各クライアントおよび各レプリカに対応するバー ジョン番号を管理する.一方,レプリカではレプリケータ からのリクエストが格納される優先度付きキューを保持す る. そして、各リクエストに付与されているバージョン番 号である QVN (reQuired Version Number) と自身のバー ジョン番号である RVN を用いてトランザクション処理の

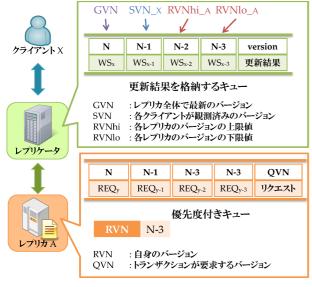


図 1 McRep における構成要素

Fig. 1 Components of the McRep protocol.

制御を行う。QVN はそのトランザクション操作が要求するバージョン番号を示し、自身のバージョン番号である RVN が (RVN  $\geq$  QVN)を満たす場合にキューからエントリを取り出し、対応するトランザクション操作を排他的に実行する。以上のように McRep では、レプリケータにおいて各トランザクションの更新結果をバージョン番号を付与してバッファリングし、各レプリカがそのバージョン番号に基づいて実行制御を行うことで複数の一貫性レベルを保証可能とする。

## 3.2 一貫性レベルの判定

レプリケータは現在設定されている一貫性レベルごとに 指定される条件を満たすレプリカを選択し、そのうちの1 つのレプリカへとトランザクション操作を転送することで 一貫性レベルを満たすレプリカでの実行を保証する.一貫 性レベルごとに選択されるレプリカの条件を**表 1** に示す.

また、レプリカは、トランザクション実行後、トランザク ションがアクセスしたデータ項目の情報とそのトランザク ション実行時のレプリカのバージョンをレプリケータへ返 す. レプリケータでは、そのデータ項目とキュー内の更新 結果のうち、トランザクション実行時のバージョン以上の ものと競合検出を行う、競合する場合は、トランザクショ ン実行時のレプリカのバージョンより大きい値の RVNhi を 持つレプリカを選択し、再度トランザクションを送信する. このトランザクションの再送信回数が一定の値を超えたら, トランザクションの競合が起こりコミットに失敗したこと をクライアントへ伝える. 競合検出の際, Linearizability, Sequential Consistency, そして One-Copy Serializability では,応答のデータ項目すべてについて競合を検査する が、Session Snapshot Isolation と Generalized Snapshot Isolation では、応答のデータ項目のうち更新操作によるも ののみについて競合を検査する. Causal Consistency の場 合は、競合検出そのものを行わない.

#### 3.3 動作フロー

McRep における動作フローを図 2 に示す。レプリケータはクライアントからリクエストを受け取ると、GVN、SVN、RVNhi の 3 つのバージョン番号を用いて、設定され

表 1 一貫性レベルごとの選択されるレプリカの条件

**Table 1** Conditions of a replica to be selected from each level of consistency.

一貫性レベル	選択されるレプリカの条件
Linearizability	RVNhi = GVN
Sequential consistency	$RVNhi \ge SVN$
One-Copy Serializability	$RVNhi \ge 0$
Session Snapshot Isolation	$RVNhi \ge SVN$
Generalized Snapshot Isolation	$RVNhi \ge 0$
Causal Consistency	$RVNhi \ge SVN$

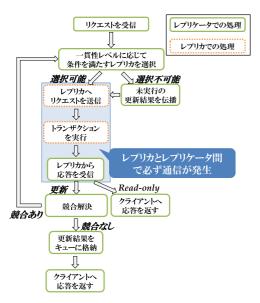


図 2 McRep における動作フロー

Fig. 2 A processing flow of the McRep.

ている一貫性レベルを満たすレプリカを選択する. 条件を 満たすレプリカが存在する場合は、そのレプリカヘクライ アントのリクエストを送信する. 条件を満たすレプリカが 存在しない場合は、任意のレプリカを選択し、レプリケー タ内のキューに格納されている更新結果のうち,選択した レプリカのデータベースに書き込まれていないもの. つま り、RVNhi より大きいバージョンを持つ更新結果を伝搬す る. 更新結果を受け取ったレプリカは、その更新結果をコ ミットし、自身のバージョン番号である RVN をインクリ メントし、(RVN = GVN) を満たすようになる. その後, レプリケータはそのレプリカへクライアントのリクエスト を送信する. そして、レプリカは、そのトランザクション 処理を実行し、レプリケータへ応答を返す.このとき、-貫性レベルが Causal Consistency であれば実行したトラン ザクションの更新結果を即時にデータベースへ書き込み, コミットを行うが、それ以外であればこの時点では更新結 果をデータベースへは書き込まず、コミットは行わない、 そして、レプリケータはレプリカから応答を受け取ると、 トランザクションの種類が Read-only トランザクションで ある場合は単純にクライアントへ応答を返す.一方,更新 トランザクションである場合は、一貫性レベルが Causal Consistency であれば受け取った更新結果を即時にキュー へ格納した後, クライアントへ応答を返す. それ以外であ れば, 競合解決を行った後で, 更新結果をキューへ格納し, クライアントへ応答を返す. また, この一連の処理におい て、レプリケータ内で管理している4種類のバージョン番 号は**表 2** に示すように更新される.

#### 3.4 問題点

McRep は、トランザクション処理を各レプリカにおい

表 2 各バージョン番号の更新動作

Table 2 Update timing of each version number.

タイミング	更新動作	
Read-only トランザクション	CVN(CVN DVN)	
に対する応答時	$SVN = \max(SVN, RVN)$	
更新トランザクション	SVN = ++GVN	
に対する応答時		
更新結果の伝搬に対する応答時	RVNlo = RVN	
更新結果の伝搬時	RVNhi = GVN	

て並行実行可能であり、かつレプリカ間で同期的に実行する必要がない非同期レプリケーションである。そのため、レプリカ数に応じて性能がスケールアウトすることが期待される。そのうえ、6つの一貫性レベルをサポートしており、汎用性のあるレプリケーションプロトコルであるといえる。しかし、McRep はミドルウェア方式のレプリケーションであるため、クライアントからの全リクエストがレプリケータを経由する必要がある。その結果、クライアント数の増加につれてレプリケータがシステム全体のボトルネックとなり、性能向上が困難になるという問題がある。

## 4. 提案手法

本章では、ミドルウェアベースではなくバックエンドベースで McRep と同様のレプリケーション制御を実現する手法を提案する.

#### 4.1 概要

既存手法である McRep は、クライアントの全リクエス トがレプリケータを経由するため、レプリケータが性能の ボトルネックとなるという問題点が存在した. これを解決 するため、提案手法では図3に示すように、レプリケー タがミドルウェア部ではなくバックエンド部で動作するよ うに拡張する. McRep では、レプリケータがクライアン トからの全リクエストに関して一貫性制御を集中的に行っ ており、レプリカは単にレプリケータの要求どおりに処理 を行う方式となっていた. これに対し、提案手法ではレプ リケータだけでなくレプリカにおいても SVN と RVNhi の バージョン管理を行い,一貫性制御を行うように拡張する. これにより、レプリケータが集中的に行っていた一貫性制 御を分担して行うことが可能となり、レプリケータの処理 量が軽減される. そのうえ, レプリケータがバックエンド 部で動作することで,設定されている一貫性レベルを満た すレプリカが存在するときはレプリケータが Read-only リ クエストを処理する必要がないため、レプリケータの処理 量を軽減することができる.

#### 4.2 一貫性制御時の動作の違い

本節では、McRep で行っていた一貫性制御時の動作と

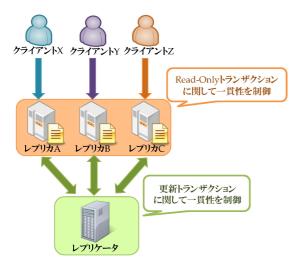


図 3 バックエンドベースのレプリケーション

Fig. 3 An architecture of the proposed method.

提案手法における一貫性制御時の動作を比較し、その違い について説明する.

• 既存手法 (McRep) での一貫性制御 既存手法である McRep において一貫性を制御する際 には、レプリケータは能動的に動作する. 具体的には、 レプリケータはクライアントからトランザクション要 求を受け取ると、まず一貫性レベルに応じて条件を満た すレプリカを選択する. 条件を満たすレプリカが存在 しない場合は、レプリケータ主導で未実行の更新結果を レプリカへ伝搬することで、一貫性レベルを保証する.

## • 提案手法での一貫性制御

一方、提案手法において一貫性を制御する際には、レプリカはクライアントからトランザクション要求を受け取ると、まず自身が一貫性レベルを満たしているかを判断する。一貫性レベルを満たしていない場合は、未実行の更新結果をレプリケータに要求することで一貫性レベルを保証する。すなわち、レプリカ主導で一貫性制御を行い、レプリケータはレプリカからの要求に応じて受動的に動作する。クライアントからトランザクション要求を受けたレプリカは、一貫性レベルを満たしている場合はレプリケータを介さずに処理を行うため、レプリケータの処理量を軽減することができる。

以上のように、提案手法ではレプリカとレプリケータで 分担して一貫性制御を行う.以降、この一貫性制御時の動 作を実現するための、レプリカおよびレプリケータ内にお ける具体的な制御について述べる.

#### 4.3 レプリカにおける制御

提案手法においてレプリカが行うのは,クライアントからのリクエスト処理とレプリケータからのリクエスト処理 および応答処理である.

## 4.3.1 クライアントからのリクエスト処理

レプリカはクライアントからリクエストを受け取ると、まず自身が一貫性レベルを満たすかどうか判断する.このために、McRepがレプリケータで管理していたバージョン番号 (GVN, SVN, RVNhi) のうち、SVN と RVNhiの管理をレプリカが担当する.各バージョン番号は次のように管理する.

#### • GVN

GVN はレプリカ全体で最新のバージョン番号を示す 必要がある。これをレプリカで管理するためには同期 をとる必要があり、通信コストが大きくなってしまう。 そのため、GVN に関してはレプリカで管理は行わず、 McRep と同様にレプリケータが管理を行う。

### • SVN

SVN は各クライアントが観測済みのバージョン番号 であるため、クライアントが自身の SVN を記憶し、 それをリクエストに付与することで、レプリカ間で同 期をとる必要がなく,各レプリカごとに独立して管理 可能である. SVN は、表 2 に示すように、各トラン ザクションの応答時に更新され, それぞれ更新時に は RVN および GVN が必要となる. RVN に関しては 元々各レプリカごとに管理しているため、特に問題は 発生しない.一方で、GVN に関してはレプリカ内で 管理していないため、提案手法では各更新トランザク ションにおいて、レプリケータが更新結果をキューへ 格納し、レプリカへ応答を返す際に GVN の値を付与 するように拡張する. つまり、レプリカは自身が管理 する SVN の値を、更新トランザクションの応答時に 受信した GVN の値に更新することで管理する. そし て、レプリカはトランザクション応答に SVN を付与 することで、クライアントはこれを自身の SVN とし て記憶することができる.

## • RVNhi

RVNhi は、各レプリカごとのバージョンの上限値を示すものであり、レプリカごとに独立して管理可能である。RVNhi は、表 2 に示すように、更新結果の伝搬時に更新されるが、GVN が必要となるため、このままでは管理できない。そのため、提案手法ではレプリカへの更新結果の伝搬時においても、レプリケータがGVN の値を付与するように拡張する。

以上のようにレプリカで SVN と RVNhi を管理することで、自身が一貫性レベルを満たしているかを判断することが可能となる。Linearizability の場合はレプリカ単独での判断はできないため、レプリカはトランザクションを受信次第レプリケータへ転送し、レプリケータにおいて GVN と等しい RVNhi を持つレプリカを選択しトランザクションを転送する。GVN 以上のレプリカが存在しない場合は、任意のレプリカにキュー内の更新結果を伝搬しレプリカ

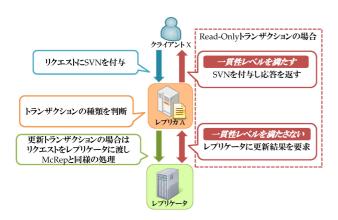


図 4 レプリカにおけるクライアントからのリクエスト処理 **Fig. 4** A process of the request from a client.

のバージョンを更新したうえで、そのレプリカにトランザ クションを転送する. また、レプリカで一貫性制御を行う 際、トランザクションの種類が更新トランザクションの場 合は, たとえ一貫性レベルが判断できたとしても, 必ずレ プリケータによる競合解決およびトランザクションの更 新結果のキューへの格納が必要になる. そのため, レプリ ケータを経由するリクエスト数の削減は期待できない. 一 方で、Read-only トランザクションの場合は、一貫性レベ ルを満たしていれば、レプリケータを介さずに処理する ことができる. まとめると、図4のように、クライアン トは SVN を付与してレプリカにリクエストを送信し、レ プリカでは受信したトランザクションの種類を判断する. Read-only トランザクションの場合は、自身のレプリカが 一貫性レベルを満たしていれば、トランザクションを実行 し即座にクライアントへ応答を返す. そうでなければ、レ プリケータに未実行の更新結果を要求する. 更新トランザ クションの場合は、リクエストをそのままレプリケータへ 渡し、McRep と同様の処理をする.

## 4.3.2 レプリケータからのメッセージ処理

レプリカはレプリケータから更新トランザクションリクエスト, 更新トランザクションリクエストの応答, そして 伝搬された更新結果を受信する. 各メッセージは以下のように処理する.

 更新トランザクションリクエストの処理 レプリカはレプリケータから更新トランザクションの リクエストを受信すると、McRep と同様に処理し、レ プリケータへ応答を返す. すなわち、Causal Consistency の場合はトランザクション実行後、コミットを 行い、RVN をインクリメントした後、トランザクション実行時に書き込まれたデータ項目である write セットと RVN を応答に付与する. Causal Consistency 以 外の場合はコミットを行わず、トランザクション実行 時に読み書きされたデータ項目である read/write セットと RVN を応答に付与する.

- 更新トランザクションリクエストの応答の処理 この場合も、受け取った応答には GVN が付与されて いるため、自身が管理する SVN を受け取った GVN に 更新する. その後、リクエスト元クライアントへ応答 を返す.
- 伝搬された更新結果の処理 この場合も McRep と同様に、レプリケータから受け 取った更新結果をコミットし、RVN をインクリメン トした後、その値を付与してレプリケータへ応答を返 す. ただし、前項で述べたように、この応答には更新 結果に加えて GVN の値も付与されているため、更新 結果をコミットする前に、自身が管理する RVNhi を GVN の値に更新する。

#### 4.4 レプリケータにおける制御

提案手法においてレプリケータは、レプリカから更新トランザクションのリクエストおよび更新結果の伝搬要求を 受け取る。各リクエストは以下のように処理される。

● 更新トランザクションの場合
レプリケータはレプリカから更新トランザクションの
リクエストを受信すると、McRepと同様に動作する。
つまり、一貫性レベルを満たすレプリカを選択し、そ
のレプリカへリクエストを送信する。レプリカ選択時
には、4.3.1 項で述べたように、リクエストに付与され
ている SVN を用いる。そして、レプリカから応答を
受け取り、競合解決を行い、トランザクションの更新
結果をキューへ格納し、要求元レプリカへ応答を返す。
また、応答を返す際にも 4.3.1 項で述べたように、レプリカ内で SVN を管理するため GVN の値を付与する。

● 更新結果の伝搬要求の場合 この場合は、RVNhi および RVNlo を用いて、要求元 レプリカにおいて実行されていない更新結果をすべて 伝搬する. ただし、このときも、4.3.1 項で述べたよう に、レプリカ内で RVNhi を管理するため、更新結果 に加えて GVN の値も付与する. そして各バージョン 番号は McRep と同様に更新される. つまり、表 2 に 示すように、更新結果を伝搬した後、要求元レプリカ の RVNhi を GVN と同じエントリを指すように更新 し、レプリカからの応答を受信後、RVNlo をレプリカ からの応答に含まれる RVN と同じバージョン番号を 指すように更新する.

またこれらの処理に加え、レプリケータは McRep と同様に、キュー内の更新結果を反映していないレプリカへ定期的に伝搬し、全レプリカへ伝搬済みとなった更新結果をキューから削除する.

#### 4.5 動作フロー

以上で述べた提案手法の具体的な動作フローを図 5 に示

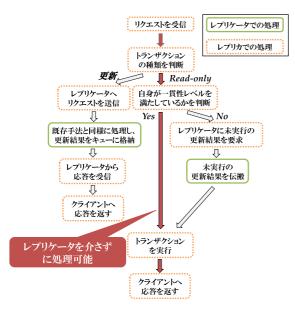


図 5 提案手法における動作フロー

 ${\bf Fig.~5} \quad {\bf The~processing~flow~of~the~proposal~method}.$ 

す. 提案手法では、クライアントからリクエストを受信し たレプリカはまず, そのトランザクションの種類を判別し, 更新トランザクションである場合は、レプリカでは一貫性 制御は行わず、そのままレプリケータへリクエストを渡 す. 一方, Read-only トランザクションである場合は, レ プリカで一貫性制御を行うため、まず要求されている一貫 性レベルを満たしているかどうかを判断する.このとき, Linearizability の場合はつねに条件を満たさないと判断し, それ以外の一貫性レベルであれば SVN と RVNhi から判断 する. そして、一貫性レベルを満たしている場合は、受け 取ったトランザクション処理を実行し、即時にクライアン トへ応答を返すことができる.一方で、一貫性レベルを満 たしていない場合は、レプリケータに対して自身の RVN が  $(RVN \ge GVN)$  を満たすように更新結果の伝搬を要求 する. レプリケータから受信した更新結果をコミットした 後, 受け取ったトランザクションを実行し, クライアント へ応答を返す。図2の動作フローと比較すると、更新トラ ンザクションに関しては同等の処理を行うが、Read-only トランザクションに関しては、レプリカが一貫性レベルを 満たしている限り、レプリケータを介することなく処理可 能である点が大きく異なる. つまり、提案手法では McRep と比較して、Read-only トランザクション実行時における レプリケータとの通信回数が削減され、レプリケータの処 理量を軽減することができる.

#### 5. 評価実験

本章では、ミドルウェアベースのレプリケーションである McRep とバックエンドベースのレプリケーションである提案手法を比較し、提案手法ではレプリケータがボトルネックとなることを回避し、性能がスケールアウ

表 3 評価環境

Table 3 The specification of PosrgeSQL, pgpool-II and client nodes.

	${\bf Postgre SQL}$	pgpool-II
version	9.3.5	3.3.3
OS	Linux 3.5.0-23-amd64 Ubuntu 12.04 Server	
CPU	Intel(R) Core(TM) i5 3470 @ 3.2 GHz	
Memory	$16\mathrm{GB}$	
Network	1000BASE-T	

トすることを示す。実装において、レプリカとして Post-greSQL [13] を用い、レプリケータとして pgpool-II [14] を用いた。PostgreSQL/pgpool-II はマルチプロセス構成であり、クライアントごとに別々のプロセスで処理される。そのため、McRep および提案手法におけるレプリケータ内のキューは、全プロセスから観測できるように共有メモリ上にリングバッファとして実装した。したがって、バッファサイズに制限があるため、バッファが一杯になった際には更新トランザクションが待たされることとなり、バッファサイズによって更新トランザクションの性能が大きく左右されることとなる。そのため、以降の評価においては待機が発生しない十分なバッファサイズで行った。

#### 5.1 評価環境

評価環境を表 3 に示す. ベンチマークには PostgreSQL 付属の pgbench を用いた. この際, 100,000 エントリを持つテーブルを用意し, 更新トランザクションではこのテーブルから1つのエントリを選択する UPDATE 文を, Read-only トランザクションでは同一のテーブルから1つ選択し読み出す SELECT 文を使用した. 以降の評価は, McRep, 提案手法についてそれぞれ行った. 評価内容として, 更新トランザクションと Read-only トランザクションの比率が1対9,5対5および9対1の場合において,クライアント数およびレプリカ数を変化させた場合におけるスループットを評価した. また, 以降の評価結果はいずれも pgbenchで60秒間トランザクションを実行し,それを10回繰り返した際のスループットの平均を測定した.

#### 5.2 クライアント数による影響

まず、レプリカ数を固定してクライアント数を変化させた場合のスループットを評価した。レプリカ数を6で固定し、pgbenchのクライアント数を50から300まで変化させた際の評価結果を図6上部に示す。一貫性レベルがSequential Consistency および Linearizability の場合は、SVN またはQVNが更新されるたびに、レプリカのRVNを更新する必要がある。このためトランザクションを受信してから実行するまでの間、レプリケータから更新結果が送信され、それらをデータベースへ書き込むまでの待機時間が発生する。この待機時間のためにMcRep、提案手法

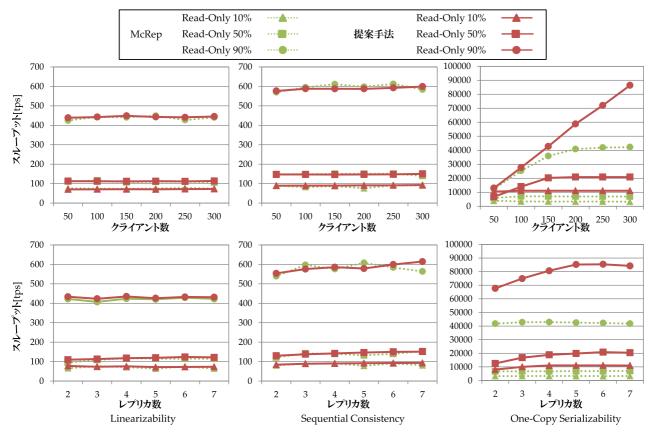


図 6 クライアント数およびレプリカ数に応じたスループットの推移

Fig. 6 Throughputs of different number of clients and replicas.

ともにクライアント数が増加してもスループットが向上しない。一方,一貫性レベルが One-Copy Serializability の場合は,この待機時間は発生しないため,両手法ともにクライアント数の増加にともないスループットが向上している。しかし,McRep では 200 クライアントで性能が頭打ちとなっている一方で,提案手法では 300 クライアントにおいても性能が向上している。さらに,最大スループットも提案手法では McRep の約 2 倍に向上している。これは,McRep がミドルウェア方式のレプリケーションであるために,pgpool-II がボトルネックとなっているのに対し,提案手法ではバックエンドベースで制御を行うことでボトルネックとなることを回避できたためである。

一貫性レベルが One-Copy Serializability かつ Readheavy なワークロード (Read-only トランザクション 90%) でのレプリケータおよび各レプリカにおける CPU 使用率とネットワーク I/O バンド幅の平均値を図 7 に示す.他の条件はスループット測定時と同様である.この条件下において提案手法では更新操作のみがレプリケータを経由するため,レプリケータにおける CPU 使用率およびネットワーク I/O バンド幅は,McRep と比較し,それぞれ最大で約 1/4,1/20 まで削減されている.一方,レプリカにおいては,クライアント数が増加すると CPU 使用率,ネットワークバンド幅ともに,提案手法のほうが McRep よりも増加している.これは,レプリケータがボトルネックと

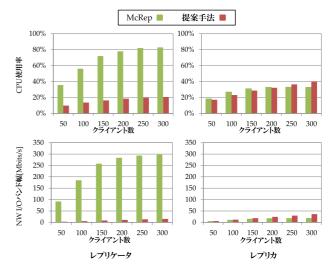


図 7 Read-heavy なワークロードにおける CPU 使用率およびネットワークバンド幅 (One-Copy Serializability)

Fig. 7 CPU usages and network bandwidths in a read-heavy workload (One-Copy Serializability).

なることを回避したことにより、レプリカにおけるトランザクションの処理数が増加したためである.

## 5.3 レプリカ数による影響

次に、クライアント数を固定し、レプリカ数を変化させた場合のスループットを評価した、評価では、十分な負荷

を与えるために、pgbench でクライアントを 300 とし、レプリカ数を 2 台から 7 台まで変化させた場合のスループットの推移を測定した。評価結果を図 6 下部に示す。一貫性レベルが Sequential Consistency および Linearizability の場合は、レプリカ数の増加にともなう顕著なスループットの向上は見られない。これは 5.2 節と同様の理由による。一方、一貫性レベルが One-Copy Serializability の場合は、提案手法ではレプリカ数の増加にともないスループットの向上が見られるが、McRep では見られない。これは、バックエンドベースでレプリケーション制御を行うことで、レプリケータがボトルネックとなることを回避できているためだと考えられる。

以上から、クラウド上に展開された多くのアプリケーションにとって望ましい一貫性レベルである One-Copy Serializability [15] において、提案手法ではレプリケータがボトルネックとなることを回避し、Read-heavy なワークロードにおいてレプリカ数に応じて性能がスケールアウトすることが確認できる.

#### 5.4 提案手法における遅延オーバヘッドの評価

提案手法はクライアントが直接レプリカと通信するため、McRepと比較してレプリカを選択する自由度がない。McRepではどのクライアントからのトランザクション要求もレプリケータを中継するため、一貫性レベルを満たすレプリカが1つでも存在すれば、即座にトランザクション要求を送信できる。一方、提案手法ではクライアントが通信するレプリカのバージョンが一貫性レベルを満たさない場合は、レプリケータを経由した後、一貫性レベルを満たす他のレプリカにトランザクション要求が送信され、その後、またレプリケータを経由し、元のレプリカに返される。この場合のトランザクション要求の通信経路は図8に示すように提案手法のほうが長くなる。

上記の場合において、McRep、提案手法での通信遅延の

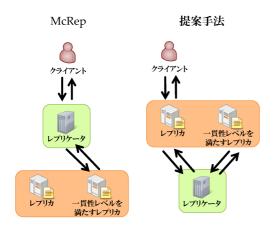


図 8 McRep と提案手法の間で通信経路に差が現れる状況

Fig. 8 A situation in which a difference appears in a communication path between McRep and proposed protocol.

オーバヘッドを評価するために、図8に示す通信経路を通るRead-only/更新トラザクション要求を100個発行し、その遅延の平均を計測した.評価環境は5.1節と同様である.評価結果は、応答遅延の平均は、Read-onlyトランザクションでは、提案手法で3.3ミリ秒、McRepで2.3ミリ秒、更新トランザクションでは、提案手法で3.7ミリ秒、McRepで2.8ミリ秒となった. どちらの場合でも遅延の差はたかだか1ミリ秒程度であるため、その影響は小さい.さらに、レプリケータが更新結果を伝搬する間隔を調整することで、このような状況を減らすことが可能である.

## 6. 関連研究

開発コストの削減やアプリケーションの性能が保証する 一貫性レベルにより受ける影響を最小化するため、複数の 一貫性を制御可能なレプリケーションプロトコルが研究さ れている. PRACTI[16] はデータの一貫性レベルを柔軟に 設定可能な部分レプリケーションシステムで, データに付 与されたバージョン番号や更新時間を指標として一貫性を 保証する.しかし、この手法は、McRep や提案手法のよ うにレプリケーション制御を行うサーバが存在しないた め,一貫性制御のためにロックによる同期が必要となる. Garcia-Recuero ら [17] は利用者が、データ項目ごとにレプ リケーションを作成するタイミングを指定することで利用 者の要求に応じたデータの一貫性を保証する手法を提案し ている. ただし、この手法は、データの更新回数もしくは 時間のみをレプリケーションを作成するタイミングの指標 として用いるため、McRep や提案手法のように更新操作 の時間的順序関係を考慮していない.

またミドルウェアを用いた手法として、AQuA [18] はレプリカ間で更新操作を全順序かもしくは FIFO 順序で順序付けするかを選択できる。Ganymed [4] は、ANSI SQL で定められた分離レベル [19] のうち 2 種類を保証できる。また、BaseCON [20] は Linearizability、Sequential Consistency および One-Copy Serializability の一貫性が保証可能な手法だが、McRep や提案手法とは違い、レプリカ間でバージョン番号の差異がある状態を許可しない。さらに、提案手法および McRep は、ミドルウェアを用いた上記の 4 つの手法よりも多くの一貫性レベルを保証することができる。

#### 7. おわりに

本研究では、既存のミドルウェア方式のレプリケーション手法である McRep を拡張し、ミドルウェアベースではなく、バックエンドベースで同様のレプリケーション制御を実現する手法を提案した。ミドルウェア方式のレプリケーションにおいては、クライアントからの全リクエストがレプリケータを経由する必要があり、レプリケータがシステム全体のボトルネックとなるという問題が存在した。これに対し、提案手法ではレプリケータがバックエンドで

制御を行うことで、全リクエストがレプリケータを経由する必要がなくなるうえ、レプリケータが集中的に行っていた一貫性制御を各レプリカにおいても行うように拡張することで、レプリケータにおける処理量を軽減することができる。評価から、提案手法は既存手法と比較して、レプリケータがボトルネックとなることを回避し、Read-heavyなワークロードにおいては性能がスケールアウトすることを確認した。

また、McRep、提案手法ともに、レプリケータが単一障害点となるという問題があるため、今後の課題として、特定のサーバが単一障害点とならないレプリケーションモデルによる複数一貫性制御手法の実現があげられる。他にも、5.4 節で述べたような状況に対応するため、より現実的なベンチマークでの評価による、提案手法および既存手法における最適な更新結果の伝搬間隔の調査も必要である。

謝辞 本研究の一部は、科研費基盤研究 (C) 24500113, (C) 15K00168 による.

## 参考文献

- Elnikety, S., Zwaenepoel, W. and Pedone, F.: Database Replication Using Generalized Snapshot Isolation, Proc. IEEE 24th Symp. on Reliable Distributed Systems, SRDS '05, pp.73–84 (2005).
- [2] Kemme, B. and Alonso, G.: Don't Be Lazy, Be Consistent: Postgres-R, A New Way to Implement Database Replication, Proc. 26th Int'l Conf. on Very Large Data Bases, VLDB'00, pp.134–143 (2000).
- [3] Krishnamurthy, S., Sanders, W.H. and Cukier, M.: An Adaptive Quality of Service Aware Middleware for Replicated Services, *IEEE Trans. Parallel Distrib. Syst.*, Vol.14, No.11, pp.1112–1125 (2003).
- [4] Plattner, C. and Alonso, G.: Ganymed: Scalable Replication for Transactional Web Applications, Proc. 5th ACM/IFIP/USENIX Int'l Conf. on Middleware, Middleware '04, pp.155-174 (2004).
- [5] Al-Ekram, R. and Holt, R.: Multi-consistency Data Replication, Proc. IEEE 16th Int'l Conf. on Parallel and Distributed Systems (ICPADS 2010), pp.568–577 (2010).
- [6] Herlihy, M.P. and Wing, J.M.: Linearizability: A Correctness Condition for Concurrent Objects, ACM Trans. Program. Lang. Syst., Vol.12, No.3, pp.463–492 (1990).
- [7] Riegel, T., Fetzer, C., Sturzrehm, H. and Felber, P.: From Causal to Z-linearizable Transactional Memory, Proc. 26th Annual ACM Symposium on Principles of Distributed Computing, PODC '07, pp.340–341 (2007).
- [8] Lamport, L.: How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs, *IEEE Trans. Comput.*, Vol.28, No.9, pp.690–691 (1979).
- [9] Bernstein, P. and Goodman, N.: A proof technique for concurrency control and recovery algorithms for replicated databases, *Distributed Computing*, Vol.2, No.1, pp.32–44 (1987).
- [10] Daudjee, K. and Salem, K.: Lazy Database Replication with Snapshot Isolation, Proc. 32nd Int'l Conf. on Very Large Data Bases, VLDB '06, pp.715–726 (2006).
- [11] Raynal, M., Thia-Kime, G. and Ahamad, M.: From serializable to causal transactions for collaborative applica-

- tions, Proc. 23rd Conf. EUROMICRO 97, New Frontiers of Information Technology, pp.314–321 (1997).
- [12] Lin, Y., Kemme, B., Patiño Martínez, M. and Jiménez-Peris, R.: Middleware Based Data Replication Providing Snapshot Isolation, Proc. 2005 ACM SIGMOD Int'l Conf. on Management of Data, SIG-MOD '05, pp.419–430 (2005).
- [13] Momjian, B.: PostgreSQL: Introduction and concepts, Vol.192, Addison-Wesley, New York (2001).
- [14] pgpool-II, available from (http://www.pgpool.net/).
- [15] Fetai, I. and Schuldt, H.: SO-1SR: Towards a Self-optimizing One-copy Serializability Protocol for Data Management in the Cloud, Proc. 5th Int'l Workshop on Cloud Data Management, CloudDB '13, pp.11–18 (2013).
- [16] Dahlin, M., Gao, L., Nayate, A., Venkataramana, A., Yalagandula, P. and Zheng, J.: PRACTI replication, NSDI (May 2006), Vol.12, p.80 (2006).
- [17] Garcia-Recuero, A., Esteves, S. and Veiga, L.: Quality-of-data for consistency levels in geo-replicated cloud data stores, Proc. IEEE 5th Int'l Conf. on Cloud Computing Technology and Science (CloudCom), Vol.1, pp.164–170 (2013).
- [18] Krishnamurthy, S., Sanders, W.H. and Cukier, M.: An adaptive quality of service aware middleware for replicated services, *IEEE Trans. Parallel and Distributed* Systems, Vol.14, No.11, pp.1112–1125 (2003).
- [19] Berenson, H., Bernstein, P., Gray, J., Melton, J., O'Neil, E. and O'Neil, P.: A Critique of ANSI SQL Isolation Levels, SIGMOD Rec., Vol.24, No.2, pp.1–10 (1995).
- [20] Zuikevičiūtė, V. and Pedone, F.: Correctness criteria for database replication: Theoretical and practical aspects, On the Move to Meaningful Internet Systems: OTM 2008, pp.639-656, Springer (2008).



#### 太田 篤 (学生会員)

1991 年生. 2014 年名古屋工業大学工 学部情報工学科卒業. 現在,同大学大 学院工学研究科創成シミュレーション 工学専攻博士前期課程在籍.



## 松野 雅也

1990年生. 2014年名古屋工業大学大学院工学研究科創成シミュレーション工学専攻博士前期課程修了. 現在, 日興システムソリューションズ.



## 川島 龍太 (正会員)

2007年岩手県立大学大学院ソフトウェア情報学研究科博士前期課程修了. 2010年総合研究大学院大学複合科学研究科修了. 博士 (情報学). 同年株式会社 ACCESS 入社. 2013年名古屋工業大学大学院工学研究科助教,現在

に至る. 主に SDN, システムソフトウェアに関する研究に 従事. IEEE 会員.



## 松尾 啓志 (正会員)

1983 年名古屋工業大学工学部情報工学科卒業. 1985 年同大学大学院修士課程修了. 同年松下電器産業(株)入社. 1989 年名古屋工業大学大学院博士課程修了. 同年名古屋工業大学工学部電気情報工学科助手. 講師, 助教授

を経て,2003年同大学大学院教授,現在に至る.分散システム,画像認識,分散協調処理に関する研究に従事.工学博士.人工知能学会,IEEE 各会員.