

2J-07 高信頼かつ高性能な非同期メッセージキューの実現手法

田中 俊介 小山 貴夫 日下 貴義 吉谷 文徳 松田 栄之

株式会社 NTT データ 情報科学研究所

e-mail: {shun, kym, kusaka, yositani, matu}@rd.nttdata.co.jp

1. はじめに

OLTP(OnLine Transaction Processing)システムでは ACID 性質(原子性・一貫性・隔離性・耐久性)を満たす必要がある。TP モニタもしくは非同期メッセージキュー(以下キューと略す)は ACID 性質の一部を提供するミドルウェアである。TP モニタは同期型であるため、サーバ停止時にはクライアントでのメッセージ送信要求は異常終了する。キューは非同期型であるため、サーバ停止時にもメッセージ送信要求を実行可能である。メッセージはサーバが再起動した時点で転送される[1]。

キューでは転送するメッセージをログとしてディスクに記録するため、ディスク I/O がボトルネックとなり、スループットが上がらないケースが多い。

本報告では、高信頼かつ高性能なキューの実現方法を提案する。また、本方式と従来方式の性能比較結果について報告する。

2. 従来キュー

2.1. 従来キューの処理方式

2台のサーバ間の通信を、従来キューⁱ⁾を利用して行う場合を図1に示す。

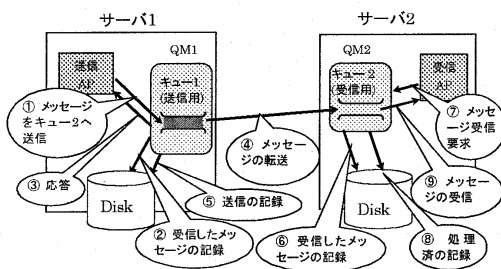


図1 従来キューでの処理

始めに、送信 AP が同一サーバ内の QM(キューマネージャの略)にメッセージの送信要求を行う。送信

要求を受けた QM1 は、受信したメッセージをログに記録し、送信 AP に応答を返す。その後、QM2 にメッセージを転送する。メッセージを受信した QM2 はログ記録を行う。受信 AP が送信要求を行うと、QM2 がメッセージを渡し、ログに処理終了を記録する。

2.2. 従来キューの問題点

最新鋭のサーバ機では、CPU の速度に比べるとディスク I/O が遅いため、OLTP システムではディスク I/O がボトルネックになる場合が多い。

従来キューでは、1つのメッセージを送信側 QM と受信側 QM の両方でログ記録を行う。ディスク I/O が多く発生するため、スループットが上がらないケースが多い。

3. 新しいキュー実現手法の提案

3.1. 目的

キューにおいて、トランザクションの耐久性を維持しながら処理性能(メッセージ転送のスループット)を高めることを目的とする。QM のメッセージ転送処理に関してログ記録の回数を減らすことで処理性能向上を目指す。

3.2. 処理方式

本報告のキューでは、1つのメッセージは、送信側 QM もしくは受信側 QM のどちらか一方のみでログの記録を行う。ログの記録は2つの QM で交互に行う。送信 AP からのメッセージ転送要求に対して、1つ目のメッセージは送信側 QM のみでログ記録を行い、2つ目のメッセージは受信側の QM のみでログ記録を行う。3つ目のメッセージは再び送信側の QM のみでログ記録を行い、以後、送信側 QM、受信側 QM で交互にログ記録を行う(図2参照)。

2種類のメッセージ転送処理の各々について図2に示す。

3.3. 信頼性(障害時の処理)

信頼性に関しては以下のようなアプローチを取る。

受信側サーバに障害が発生した場合の処理を図3に示す。QM2 でログ記録を行う処理(図2の偶数番号メッセージの処理)では、QM1 は QM2 でのログ記録が完

A Design of a Reliable and Scalable Asynchronous Message Queue

Shunsuke Tanaka, Takao Koyama, Takayoshi Kusaka, Fuminori Yoshitani, Shigeyuki Matsuda
Laboratory for Information Technology,
NTT DATA CORPORATION

i)従来キューの処理方式は、代表的なキューの実装である IBM 社の MQSeries[®][2]の処理方式に基づいて記述した。MQSeries[®]は IBM 社の商標。

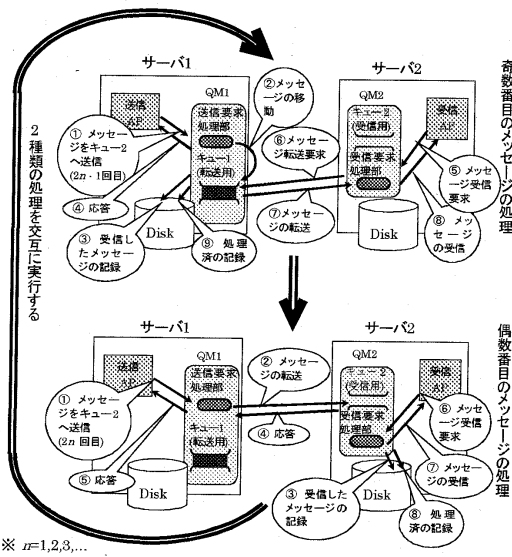


図2 本報告のキューの処理方式

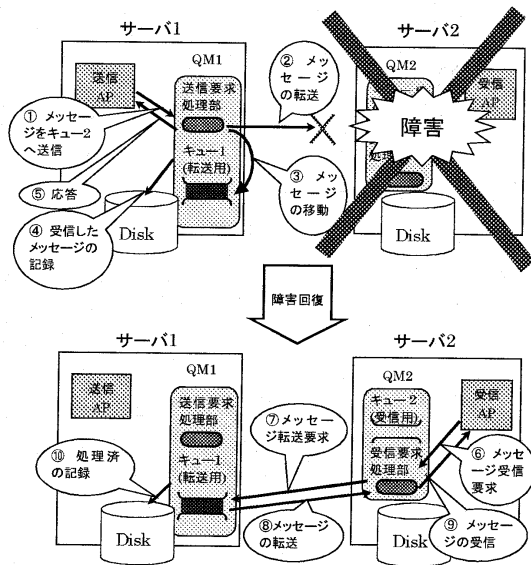


図3 障害時の処理の流れ

了したことを確認してから送信 AP に応答を返す。サーバ 2 に障害が発生した場合、QM1 から QM2 へのメッセージ転送は失敗し、QM1 はタイムアウトによって失敗を検知する。そのとき、QM1 は自サーバのディスクにログ記録を行い、ログ記録が完了してから送信 AP に応答を返す。以上のように、受信側サーバ障害時には全てのメッセージを送信側 QM でログに記録するため、本方式ではトランザクションの耐久性を実現可能である。

送信側サーバに障害が発生した場合に関しては、従来のキューと同様に、特別な対策は必要でない。これは送信 AP がメッセージ転送要求を行えないことから、QM で処理を継続する必要があるためである。

4. 評価

従来のキューと同等の処理方式を持つ QM と本方式の処理方式を持つ QM のプロトタイプを実装した。2 種類の QM のそれぞれに対して、2 台のサーバの間でメッセージ転送を行い、送信 AP・受信 AP の数、メッセージサイズを変化させた場合のスループットを測定した。

測定した結果を図 4 に示す。図 4 から、本方式の方が従来方式よりも高いスループットが得られることが確認できた。特に、AP 多重度 2 以上の場合にその差が顕著に現れるという結果が得られた。

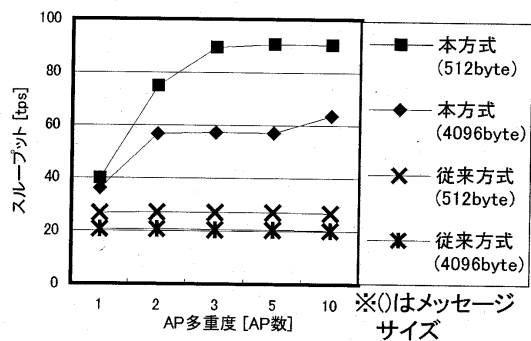


図4 性能比較の実験結果

5. おわりに

従来の高信頼非同期メッセージキューでは、全てのメッセージを送信側マシンと受信側マシンの両方でログに記録するため、ディスク I/O が多く発生し、スループットが上がらないという問題がある。本報告では、各メッセージを送信側もしくは受信側のどちらか片方のみでログ記録を行う処理方式を提案した。また、提案した方式に基づいてプロトタイプを実装し、従来の処理方式との性能比較から、本方式の優位性を確認した。また、本方式が性能の確保と共に十分な信頼性が実現できることを示した。

【参考文献】

- [1] Philip B. and Eric N., トランザクション処理システム入門, 日経 BP 社, Mar. 1998
- [2] MQ Series® Planning Guide, IBM, Apr. 1999

ii) 実験は以下のマシンを 2 台使用して実施した。
 Sun Ultra1, CPU:UltraSparc167MHz×1, Memory:256Mbyte,
 Disk:FastSCSI 内蔵 Disk 2.1G×1, Ethernet:100Base-T,
 OS:Solaris 2.6, コンパイラ:gcc 2.8.1