

マルチプロセッサシステムにおける スケーラビリティボトルネックの分析手法

NEC 共通基盤ソフトウェア研究所

堀川 隆

はじめに

近年の IT システムは、マルチコア CPU の普及に伴って CPU 資源が強化されていますので、CPU 以外の要素が性能上のボトルネック（以下、ボトルネック）になることがあります。特に、排他制御のためのロック（論理資源）が新たなボトルネックになると、物理資源（CPU, Disk, Network, Memory）を強化するという単純な対策では性能は向上しない、すなわちシステムにスケーラビリティがないことから、その対応は困難になりがちです。

このチュートリアルでは、このような状況への体系的な対処として、イベント・トレースを利用した方法を紹介します。事例として、この方法をデータベース処理のベンチマーク・プログラムに適用し、1)マルチコア・マシンの処理時間に大きく影響していたロックを特定できたこと、2)プログラム構造を見直すことでスループット性能を向上できたこと、を紹介します。

マルチプロセッサ特有のボトルネック

プログラムにおいて複数のプロセッサが同時に実行すると処理結果が不正となるプログラム部分（クリティカル・セクション）の保護は、ロックなどによる排他制御によって実現されています。クリティカルセクションを実行しようとして複数のプロセッサが競合すると、ロックを獲得できなかったスレッドはその解放を待ち合わせるため、必要な処理が進まなくなります。これが、ロックがボトルネックになるメカニズムです。

排他制御のバリエーションには、1)保護対象データの読み込み操作については同時実行を可能とするための **Read-Write Lock** (`rw_lock`)、2)並行実行する処理（スレッド）の数を性能上の観点から制限するためのもの (`conc_lock`)、があります。

複数のプロセスやスレッドが多数のデータを共有して行う処理が基本となっているデータベース管理システム（DBMS）の場合、それらのデータへのアクセスを制御する排他制御の影響が大きくなり、効率的な並列実行が阻害されがちです。IT システムの重要な土台である DBMS は、排他制御によるボトルネックが発現する典型例といえます。

イベント・トレースによるボトルネックの分析

DBMS プログラムへのプローブ設置

イベント・トレースは、測定対象システムの動作を時系列として記録したデータで、測定点（プローブ）を設置したプログラムを実行させて採取します。ここでは、物理資源の使用状況を把握する目的でカーネル内に設置したプローブに加え、各種ロックについての待ち状況を把握する目的で、測定対象DBMSのmysql¹内部にもプローブを設置しました。プローブを設置した関数数を図1に示します。

```
mutex mutex_enter_func()
rw_lock rw_lock_s_lock_func(),
        rw_lock_x_lock_func()
conc_lock innodb_srv_conc_enter_innodb()
```

図1 プローブを設置した関数

ベンチマークの実行方法と結果

ここでは、16プロセッサ（4core CPU ×4）のマシンによるDBT-1ベンチマーク・プログラムの実行を測定しました。mysqlやDBT-1の構築やベンチマーク実行方法、および、mysqlの実行パラメータの設定は、OSS iPediaにある資料に沿って実施しました。mysqlの実行パラメータの内、並行実行処理の上限値（innodb_thread_concurrency）は、ロック待ち状況に大きく影響するので、6、

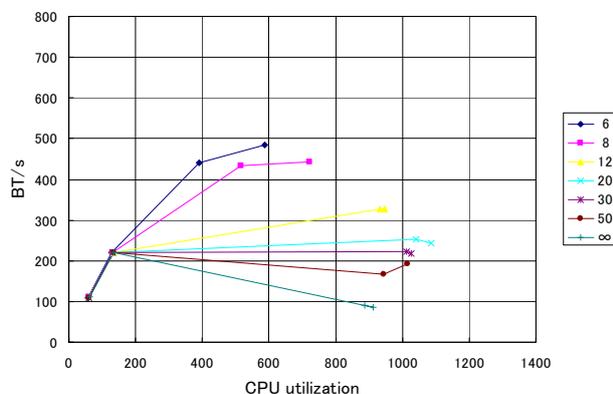


図2 CPU 使用率とスループット

8, 12, 20, 30, 50, 無制限と変化させて測定を実施しました。また、性能上の観点から、予め立ち上げておくスレッド数（thread_cache_size）を200としました。

ベンチマーク実行結果として、innodb_thread_concurrency 毎にプロットしたCPU使用率（x軸）とスループット（y軸）の関係を図2に示します。CPU使用率は、1CPU分を100%としていますので、16CPUでは最大1600%になります。スループット値は、DBT-1の結果集計プログラムが出力するBT/s（BogoTransaction/秒）です。図からは、並行実行処理の上限値を6や8に制限した場合を除き、CPU使用（消費）がスループット向上に役立っていないことが分かります。

ロック待ち時間の分析

スペースの関係上、並行実行処理数が無制限のケース、および最も高いスループットを示したケース（並列実行処理の上限値は6）の2種類の結果を図3に示します。

¹ mysql のバージョンは 5.0.75、使用したストレージエンジンは innodb

並行実行処理数の上限値が有限の場合（図3左）、conc_lockの待ち時間が大半となりますが、本質的なボトルネックは、conc_lockによって保護された処理に隠されていると考えられます。並行実行処理数が無制限のケース（図3右）は、その本質的なボトルネックを明確にするためのものです。これより、rw_lockであるbtr_search_latchの状態を管理する変数への排他アクセスを実現するためのmutexに関する競合が最大のボトルネックとなっていることが分かりました。

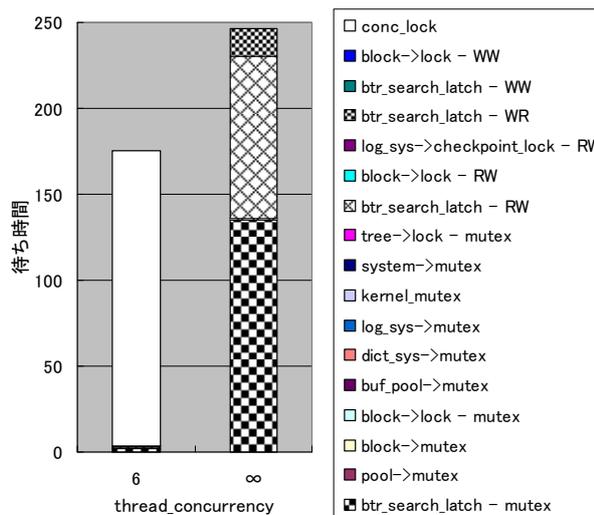


図3 ロック待ち時間の内訳

凡例には、待ち時間がゼロでなかったlockを示しています。大半のlockは待ち時間がゼロに近いため、棒グラフでは存在していないように見えます。

チューニング実施および効果の検証

チューニング内容 DBMS

問題となっているmutexで保護されている変数は、幸いにも、1)アクセス方法が単純、2)アクセスしている箇所は少ない、3)データのサイズが小さい、という特徴がありました。このため、rw_lockに関するmutexを不要とする改造は容易と考え、Compare And Swap (CAS) 命令を利用したlock-freeな方式を実装しました。

具体的には、1)mutexによる排他制御を行ってアクセスするrw_lock構造体中のメンバをメンバをCAS操作可能な64ビットに集約、2)rw_lock変数(lock)をアクセスする部分を図4に示す構造に変更、という書き換えを行ないました。なお、1)に際して、メンバの1つであるwriter_threadは64bitデータ(ポインタ)のpthread_self()から32bitデータのgettid()を使用するように変更しています。

```

(a) 元のコード
mutex_enter(&(lock->mutex));
lockのrw_status部分を更新;
mutex_exit(&(lock->mutex));

(b) CASによりlock-freeとしたコード
while(1) {
    new = old = lockのrw_status部分;
    newを更新
    if(CAS(&lock, old, new)が成功) break;
}

```

図4 lock-freeなrw_statusの更新

(b)は、理論的にはll/sc(Load-Link/Store-Conditional)として定式化されています。局所変数のoldにlockのrw_status部分を代入しているのがll操作、CAS命令がsc操作に対応しています。

効果

改造したmysqlでDBT-1を実行させたときのCPU使用率とスループットの関係を図5に示します。図2との比較から、負荷が高い領域でのスループットが改善されていることが分かり

ます。これにより、イベント・トレースを利用した測定・分析は、ボトルネックを正しく特定できていたといえます。

おわりに

排他制御（ロック）がボトルネックとなっている状況への体系立てた対処として、イベント・トレースを利用する方法を紹介しました。個々のボトルネックはシステム基盤やアプリケーションによって異なる様相を呈しますが、その対応に際して基本となる考え方は共通しています。本稿がマルチコア・システムの性能をチューニングする際の参考になれば幸いです。

参考文献

- [1] Database Test Suite, <http://sourceforge.net/projects/osldbdt/>
- [2] MySQLに対応した評価ツールDBT-1を利用したハードリソース変更によるパフォーマンスへの影響の考察, <http://ossipedia.ipa.go.jp/capacity/EV0612260303/>
- [3] OSS 技術開発・評価コンソーシアム, 「OSS 性能・信頼性評価 / 障害解析ツール開発」, DB層～OSDL DBT-1/3 による DBMS 評価編～
<http://www.ipa.go.jp/software/open/forum/development/download/051115/db-dbt.pdf>
- [4] Compare-and-swap, <http://en.wikipedia.org/wiki/Compare-and-swap>
- [5] Load-Link/Store-Conditional,
<http://en.wikipedia.org/wiki/Load-Link/Store-Conditional>
- [6] 堀川 隆、DBMS におけるスケーラビリティボトルネックの分析、情報処理学会研究報告、2009年8月、Vol. 2009-EVA-29, No. 2.

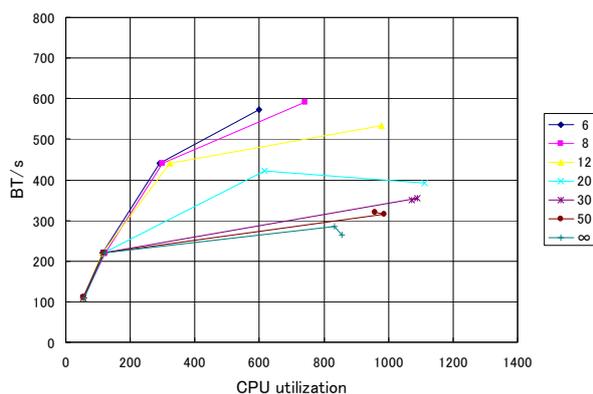


図5 CPU 使用率とスループット
(チューニング後)