

ホストマシンの NUMA トポロジーに基づく vCPU スケジューリング

林 遼[†] 味曾野 雅史[†] 品川 高廣[†]

1. はじめに

マルチコアの普及により、NUMA (Non-Uniform Memory Access) アーキテクチャのマシンはデータセンター等でも広く使われるようになった。NUMA は、ローカルノード内のメモリアクセスは高速、リモートノードへのメモリアクセスは低速という非対称的構造になっており、メモリアクセスの局所性を利用して、頻繁に参照するデータをローカルノードに置くことで、全体としてメモリアクセスの高速化を実現している。従って、メモリアクセスが頻繁なアプリケーションを効率よく実行するためには、プロセスのスケジューラが NUMA ノードを考慮する必要がある。

しかし、現在の仮想化環境では NUMA アーキテクチャを十分に活用することは難しい。まず、通常の仮想化環境ではホストマシンの NUMA トポロジーは VM (仮想マシン) からは隠蔽されており、ゲスト OS が NUMA ノードに基づくスケジューリングをおこなうことは難しい。また、VMM (仮想マシンモニタ) が NUMA ノードに基づくスケジューリングをおこなう手法¹⁾も提案されているが、ゲスト OS の内部状態を用いずに効率の良いスケジューリングをおこなうことは難しい。pCPU (物理 CPU) を vCPU (仮想 CPU) にピンニングすることで、VM が NUMA ノードをまたがないように設定する手法もあるが、VM に割り当てられる pCPU の数が固定化されてしまう。

本研究では、ゲスト OS が NUMA に基づく効率の良いスケジューリングをおこなうために、ホストマシンの NUMA トポロジーを再現した VM を作成する。一方、VM の負荷変動時に利用する pCPU 数を柔軟に増減するために、ゲスト OS 内のエージェントで vCPU を専有・解放して、ゲスト OS がプロセスをスケジューリング可能な vCPU の数を実質的に増減する。解放した vCPU に対応する pCPU を別の VM に割り当てることで、VM 間での負荷分散を実現する。

本稿では初期評価として、ホストの NUMA ノードを意識した vCPU 割り当てをするかどうかで、メモリアクセス中心のアプリケーションの実行時間に差が出るのかどうかを評価した。

2. 提案手法

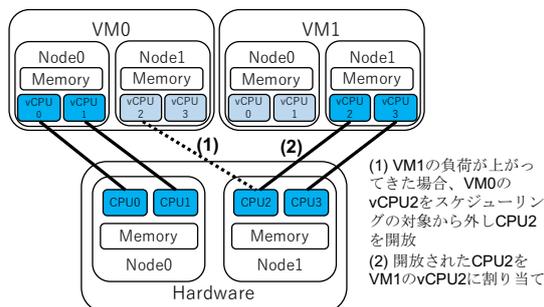


図 1 NUMA トポロジーを意識した動的な物理コアの受け渡し

本研究では、NUMA アーキテクチャをもつホストマシン上に複数の VM が存在する環境を想定する。まずそれぞれの VM にホストマシンの NUMA トポロジーを再現して、vCPU と pCPU の 1 対 1 のピンニングをおこなう。vCPU の総和が pCPU よりも多い場合は pCPU と vCPU は一対多の対応関係となってしまうが、VM の負荷に応じてゲスト OS で vCPU を動的にスケジューリング対象から外して、対応する pCPU を VM 間で受け渡しすることで、実質的に pCPU に対応する vCPU の数は常に高々 1 個であるようにする。受け渡しの際は、切り離すコアと受け取るコアの NUMA ノードを意識して、切り離す際はリモートノードのコアを、受け取る際はローカルノードのコアを優先して受け渡す。これによって、メモリアクセスを全体で最適化して、メモリアクセス中心のアプリケーションの実行時間短縮を実現する。

図 1 の例では、VM0 と VM1 のそれぞれにホストマシンの NUMA トポロジーが再現されており、初期状態では VM0 には pCPU0,1,2, VM1 には pCPU3 が実質的に割り当てられている (図 1(1))。このとき、VM0 のゲスト OS には NUMA トポロジーが見えているため、NUMA ノードを跨いでも効率の良いスケジューリングを行うことができる。ここで、VM1 の負荷が上昇したため vCPU の割当数を増やすことを考える。まず、VM0 でリモートノードである pCPU2 を選択し、対応する vCPU2 をスケジューリングの対

[†] 東京大学
The University of Tokyo

象から外すことでpCPU2を実質的に解放する。次に、VM1でvCPU2をスケジューリング対象に加えることで、pCPU2を専有して利用可能にする(図1(2))。一連の処理を動的に行い、不必要なメモリアクセス遅延要因を生まないvCPUの負荷分散を実現する。

3. 予備評価実験

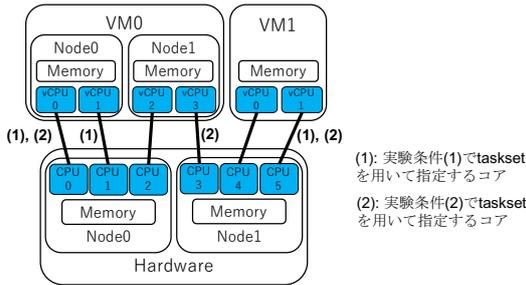


図2 実験条件 (1), (2)

VM上でホストマシンのNUMAトポロジーを再現することで、ゲストOS上のアプリケーションの実行時間に差が出るのかを確かめるために、まずは静的ピンニングによる予備評価実験をおこなった。この実験では、VMに1つのNUMAノード内のpCPUを割り当てた場合と、複数のNUMAノードのpCPUを割り当てた場合で実行時間の差を測定した。実験に用いたホストマシンのCPUはIntel Xeon X5690 × 2、各CPUは6コア12スレッドで計24個のpCPU、2ソケットのそれぞれがNUMAノードである。

実験では、ホストマシン上でCentOS7を動かして、KVMを用いて2つのVMを作成した。VM0にはNode0の12個+Node1の4個のpCPUをピンニングしたvCPU16個、VM1にはNode1の8個のpCPUをピンニングしたvCPU8個を割り当てた。それぞれのVMのメモリは4GiB、ストレージは20GBである。

実験では、VM1でメモリ負荷をかけながら、VM0でベンチマークを実行した。メモリ負荷をかけるのに用いたのはstress-ngコマンドで、ベンチマークはPARSEC³⁾のcannealである。VM1はvCPU6個をtasksetで指定してstress-ngを実行し常にNode1のメモリに負荷をかける状態である。この状態で、VM0において以下のそれぞれの条件を満たすvCPU6個をtasksetで指定して実験をおこなった(図2)。

- (1) Node0にピンニングされたvCPUのみ
- (2) Node0にピンニングされたvCPUとNode1にピンニングされたvCPUそれぞれ3個ずつ

図3に実験結果を示す。NUMA-awareは(1)の条件、NUMA-unawareは(2)の条件での結果を示し、同じベンチマークを5回実行して実行時間をそれぞれ計測した。実験結果から、NUMAを意識したvCPU

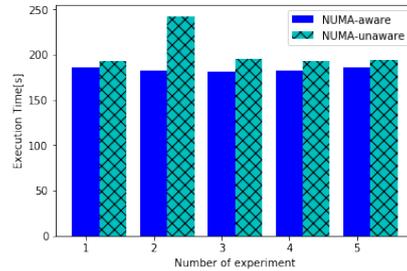


図3 canneal 実行結果

割り当てをした場合の方が実行時間が安定しており、平均8.9%、最大24.6%削減できていることがわかる。

4. 関連研究

vProbe¹⁾では、VMMで一定の周期ごとに各vCPUで実行されているアプリケーションのメモリアクセス量を性能カウンタで取得して、各NUMAノードに再分配している。提案手法では、VMにNUMAノードを再現することで、ゲストOSによるNUMAノードに基づくスケジューリングを可能にしている。

vScale²⁾では、マイクロ秒単位でvCPUの数を増減できる機構により、OSとVMMの二重スケジューリングの問題を避けつつVM間のvCPUの受け渡しを柔軟におこなえるようにしている。しかし、vScaleはNUMAアーキテクチャを考慮していない。

5. まとめと今後の予定

NUMAアーキテクチャが仮想化環境では十分サポートされていない問題に対して、VM上でホストマシンのNUMAトポロジーを再現したvCPUの割り当てを動的に行う手法を提案し、その予備実験として静的にNUMAを再現した割り当てをした場合の性能評価をおこなった。実験によって、NUMAを再現したvCPUを割り当てをした場合はそうでない場合に比べて、ベンチマークの実行速度が平均8.9%程度高速化するという結果が得られた。

今回は静的に割り当てた場合の評価を行なったが、今後は動的にvCPUを受け渡す手法を実装して、その性能を評価する予定である。

参考文献

- 1) S. Wu, et al. "vProbe: Scheduling Virtual Machines on NUMA Systems", In Proc. of 2016 IEEE International Conference on Cluster Computing (CLUSTER), 2016
- 2) L. Cheng, et al. "vScale: Automatic and Efficient Processor Scaling for SMP Virtual Machines", In Proc. 11th European Conference on Computer Systems (EuroSys '16), 2016.
- 3) <http://parsec.cs.princeton.edu/>