

1 R-04 PVMによる中規模クラスタシステムの負荷分散と効率について

黒住祥祐 平石裕実 竹内 勉 大本 英徹

京都産業大学工学部

1. まえがき

多数のPC (パソコン) をLANで接続した分散処理はこれまでも広く行なわれている。しかし、クラスタシステムとして、実用化されているものは、数台から数十台のものが多い。

PCの高性能化とLANの高速化により、量産効果の高いこれらの製品を使ったクラスタマシンに注目する。本文では、約500台のPCをPVM (Parallel Virtual Machine) によりクラスタシステムとして構成し、その負荷分散と効率について報告する。

2. KSU Linux クラスタ

本学では1999年4月から508台のPCを7教室に設備し、全学の情報教育に活用している。CPU400MHz、メモリ128MB、ディスク8GB、LAN100Base-TxのPCにLinuxとNTのOSを搭載している。通常は授業使用と自由使用により、1台のPCとして使用する。この運用と保守のために、LANにより集中的に電源投入切断と状態監視のほか、OSの更新も可能である。クラスタサーバ、NFSサーバ、WEBサーバなどは常時稼働している。

この管理機能は学内のすべてのPCからも利用資格があれば利用できる。授業時間は教育用に優先であり、ほぼ1:2の割合でLinuxの方が少ないが、夜間は全台をLinuxとして立ち上げることができる。これらをPVMによりクラスタシステムとし、最大508台のLinuxクラスタを可能としている。

3. クラスタシステムのベンチマーク

本クラスタシステムは専用ではなく、ユーザや利用OSが勝手に変わり、不意に電源が切られるような過酷な条件でクラスタシステムを構成する。短時間で全PCの状況を把握し、PVMの動作効率などの程度であるかを知るためにベンチマークプログラムを作成した。ベンチマークの機能は

- (1) 全台のOS利用状況一覧
- (2) PVMの開始と終了の処理時間予測
- (3) クラスタシステム処理能力の予測

である。

Load Share and Efficiency for Medium Scale Cluster System using PVM

Yoshisuke KUROZUMI, Hiromi HIRAISHI, Tsutomu TAKEUCHI and Eitetsu OOMOTO
Kyoto Sangyo University Faculty of Engineering

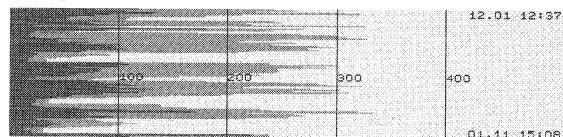
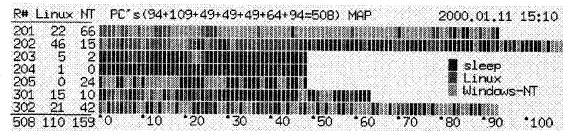


図1 OS利用状況

508台のOS利用状況の一覧を図1に示す。各PCは未使用、Linux、NTの3状態のいずれかである。Linuxとして使われていれば、そのホスト名によりPVMを立ち上げる。

PVMの開始と終了はPVMコマンドで行い、普通は1台あたり約0.1秒で復帰する。しかし、この間にPCの電源断などの状態変更があると、PVMは数分間のタイムアウトまで待ち、全体で数10分もかかることがある。この状況をチェックすることができる。

クラスタシステムの処理能力はCPU、メモリ、ディスク、ネットワークの総合性能である。従来からある並列処理用のベンチマークは連立方程式計算や流体解析であり、明らかに本システムのような条件では高性能は望めない。本システムが得意なのは、完全並列に近い問題とか、台数や負荷の変動に簡単に対応できる問題に限られる。そこで、つぎのベンチマークを考えた。

- (1) 素数計算
- (2) 全メモリ検索
- (3) 同時転送

これらの処理を各種のパラメータで実行した結果をつぎに紹介する。

4. 素数計算

素数計算は理論と応用に魅力ある問題であり、多くの試みがなされている。ここでは、数nまでの素数を求め、その個数を並列処理により計算する。素数を求める方法にエラトステネス法があり、1台のPCによる処理時間は $n=10^9$ までの素数で約200秒である。処理時間は $n \cdot \log n$ に比例することが知られている。これ以上になると32ビットのPCでは桁数とメモリ制限で計算困難となる。

そこで、m台のPCで並列処理するため数nをm分割し、それぞれのPCで素数の個数を求め合計する。この方法で処理した結果を図2に示す。

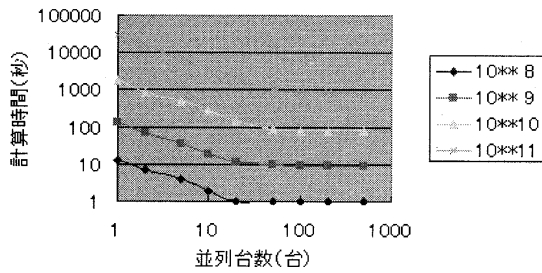


図2 並列台数と計算時間

横軸はPCの台数(1-500台)で縦軸は並列処理による計算時間(秒)である。数nの範囲は $10^{**}8$ から $10^{**}11$ までである。PCの台数が20-50台までは順調に台数に比例して性能が向上するが、そのあとは台数を増してもほとんど性能向上はない。

クラスタシステムには転送時間を含めPVMなどのクラスタOSのオーバーヘッドが常に存在する。この問題ではオーバーヘッドはつぎようになる。

分割時間	送信時間	素数計算時間	受信時間
2%	4%	90%	4%

この例の合計約10%のオーバーヘッドでは上のような現象は起きない。

エラトステネス法では大量のメモリを使用する。数nが $10^{**}9$ であれば125MB、 $10^{**}12$ であれば12GB必要であり、並列処理により、データを分散させる必要がある。

たとえば、 $10^{**}9$ までの計算で、等分割10台で並列処理し、900000000-1000000000の間の素数を計算すると12.5MBのメモリを使い、処理時間は21秒となる。分割数を大きくし、最後の分割範囲のみに注目し、プロットしたグラフを図3に示す。

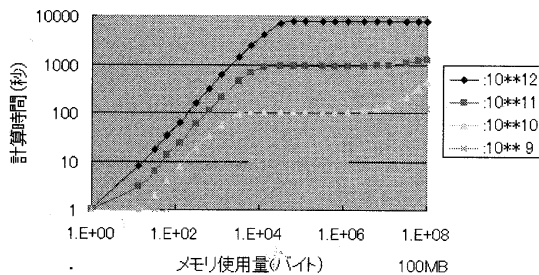


図3 メモリ使用量と計算時間

横軸はメモリ使用量(1K-100MB)で縦軸は1台のPCによる計算時間である。数値の範囲は $10^{**}9$ から $10^{**}12$ までである。メモリ使用量が500-50000B(50KB)まではメモリ使用量に比例して計算時間が増加する。しかし、そこから3M-100MBまでは計算時間は増加しない。その後、メモリ使用量が増すと、またまた増加に転じる。

最初の増加部分はキャッシュでの処理であり、つぎのフラット部分はキャッシュとメモリでの処理であり、最

後の増加部分はディスクスワップ処理となる。大量のメモリをとびとびにアクセスする問題ではこの例のような現象となり、2次元配列の参照でも起こりうる。

このようにPC台数に対して処理時間が非線形になる問題では、台数が増えれば速くなるとは限らず、最高効率を与える最小台数が存在することになる。これを求める負荷分散アルゴリズムを開発すれば大規模クラスタを多ユーザが効果的に分散使用できるであろう。

5. 全メモリ検索

クラスタシステムの効用にメモリの拡大がある。1台のPCは128MBのメモリであるが、500台あれば、60GBを超える容量となる。たとえば、50GBの文章を全PCのメモリ上に作成し、単純な全文検索したところ4秒でヒット件数の結果を得た。1台110MBの検索は2秒であり、500台のクラスタオーバーヘッドが2秒であった。

この処理は並列化効率のよい例で、50GBのメモリを持つ1台のPCでは1000秒であり、250倍の効率となる。単純全文検索ではなく、処理時間の多い複雑な画像検索ではさらに効率アップとなる。

6. 同時転送

クラスタシステムでは、PC間のデータ転送時間が問題である。100Base-TxのLANで接続しているため、少量データを瞬間的に送ることは不得意で、1パケットあたり約5m秒の転送時間である。LANは本質的にシリアルとなるためパケット数に比例して処理時間は増加する。

また、並列処理の効果をあげるためには、同時に処理終了にすることである。このとき、最後にサーバに結果を集めるために直列処理となり、数バイトのデータでも500台で約2秒必要である。大量データを転送して終了するときは、サーバに500台のPCから一齐にデータを送ると衝突のため極端に遅くなることがあり、転送時間分ずらす工夫が必要となる。

7. あとがき

中規模クラスタシステムの利用状況を述べた。PCやLANはますます高性能化するであろう。しかもコストパフォーマンスは改善され、普及するに違いない。

クラスタシステムとしては、ラックマウントPCを組み込んだ専用クラスタと本システムのように、各台が単独にも使われる共用クラスタが普及すると思われる。とくに共用クラスタでは、突然の電源投入切断や負荷の激変があり、これに耐えしかもこれを活用できるクラスタOSを開発したいし、また、その出現を期待する。

参考文献

- [1] Al Geist, et al "PVM: Parallel Virtual Machine" MIT Press, 1994
- [2] Rajkumar Buyya "High Performance Cluster Computing" Prentice Hall, 1999