

SAN 結合 PC クラスタ上での動的資源割り当てを用いた
並列データマイニング処理

†合田 和生 ‡田村 孝之 §小口 正人 †喜連川 優

† 東京大学生産技術研究所

‡ 三菱電機株式会社

§ 中央大学研究開発機構

概要

次世代のデータベースプラットフォームとして PC クラスタが注目されているが、並列分散システムの中ではノードごとの負荷の偏りが並列効果を阻害することが古くから知られ、研究が行われてきた。特に PC クラスタではノード数が膨大で異質のノードが混在することが想定されるため、実運用の為に非均質環境に於ける効果的な負荷分散機構が求められている。一方、従来利用されて来た PC クラスタではストレージが CPU ノードに従属する Shared Nothing 方式のディスクアクセスが主流であったが、次世代のストレージ技術として注目されている Storage Area Network(SAN) の登場により、Shared Disk 方式のアクセス方法が検討されている。そこで、本論文では並列相関ルール抽出処理を例に、SAN を適用した PC クラスタにおける動的負荷分散機構について設計および実装による評価実験を行う。さらに、実行時に系のバウンド状態から CPU およびディスク資源を動的に拡張する機能を実現し、1 ノード/1 ディスクから開始して 16 ノード/16 ディスクへと実行時に CPU ・ディスク資源を拡張する実験を行う。

Parallel Data Mining with Dynamic Resource Allocation
on SAN-connected PC Cluster

†GODA Kazuo, ‡TAMURA Takayuki, §OGUCHI Masato and †KITSUREGAWA Masaru

†Institute of Industrial Science, The University of Tokyo

‡Mitsubishi Electric Corporation

§Research and Development Initiative, Chuo University

Abstract

The PC cluster system is becoming remarkable as a next-generation database platform. But on such a parallel distributed system, the load skew of nodes aggravates the parallel effectivity. Especially the PC cluster system has lots of various nodes, so efficient load balancing technology in heterogeneous environment is required for practical use. In addition, although shared-nothing-style storage access has been major in the conventional PC cluster system, today shared-disk-style access is being researched and developed since Storage Area Network(SAN) came on stage. In this paper, parallel association rule mining is picked up and we designs and implements dynamic load balancing system on our SAN-connected PC cluster. Load skew across the cluster system is decreased and resources such as CPUs and storage are dynamically allocated so that the execution time can be improved.

1 序論

近年 PC クラスタはそのスケーラビリティとコストパフォーマンスの面から、次世代の大規模並列計算機システムの中心的な役割を果たし、多くのデータベースアプリケーションのプラットフォームとなるべく注目されている。一方、並列分散システムの下では、演算処理にともなうノード間の負荷の偏りが並列効果を阻害することが古くから知られ研究されて来た。データベースにおけるスキューに関しては、Walton[3]らによって分類がなされて、近年問い合わせ実行時に動的に負荷分散処理を行う機構の研究 [12][8] が行われて来た。

ところで従来利用されてきた PC クラスタでは、全てのストレージドライブが各ノードに従属する Shared Nothing 方式が主であったが、次世代のストレージ技術として注目されている SAN(Storage Area Network) の登場により、共有ディスク (Shared Disk) 方式のディスクアクセス手法が検討されている。また PC クラスタシステムが普及するに伴い、資産の効率的利用という側面から高速なノードと低速なノードが混在した環境での負荷分散機構が求められている。

そこで、本論文ではデータベースアプリケーションの一つである、データマイニングアプリケーションにおける相関ルール抽出処理を例に、SAN を適用した PC クラスタにおける動的負荷分散機構について設計を行い、実装による実験を行う。その中で、非均質環境においてノード間の負荷バランスを調節することによりボトルネックを除去し並列効果を高めると共に、実行時に動的に CPU やディスク資源を割り当てることにより実行時間の短縮を実現する。

本論文ではまず、第 2 章において SAN 技術と実験に用いた SAN 結合 PC クラスタシステムについて述べ、第 3 章においてデータマイニングにおける相関ルール抽出処理を解説する。第 4 章では SAN を利用した負荷分散機構の設計を行いその評価実験から有効性を確認する。その後、第 5 章では負荷分散機構を利用して動的 CPU 投入と動的デクラスタ [13] を併用することで、系のバウンディングファクタにより CPU およびディスク双方の資源を動的に割り当てる実験を行う。最後に第 6 章でまとめを行う。

2 SAN と SAN 結合 PC クラスタ

2.1 Storage Area Network(SAN)

SAN(Storage Area Network) は 1 台のホストコンピュータと従属的な関係を保って来たストレージデバイスを、専用のネットワークによって N:N に結合するものである。従来ストレージの共有は IP ベースの LAN 経由で NFS(Network File System) 等の NAS(Network Attached Storage) 方式を用いて行われており、LAN の輻輳の原因となっていた。

特にデータベースでは巨大なシーケンシャルファイルをたびたび扱うため、NAS 方式のストレージ共有は、TCP/IP のプロトコルソフトウェア処理による CPU 負荷が無視できない程となる。一方、表 1 の比較が示す通り、LAN(Gigabit Ethernet) では 1500B 毎に CPU が介入する必要があるが、現行 SAN で支配的な FC(Fibre Channel) では CPU の I/O 命令一つで 128MB のデータを扱えるなど、シーケンシャルデータの通信に優位な設計がなされている為、ストレージアクセスの CPU コストを低く抑えることができる。

現在、SAN では仮想化 (Virtualization[2]) と呼ばれる技術が注目されている。仮想化は接続された複数の多様なストレージデバイスをアプリケーションが効率的に扱うための技術である。例えば、VERITAS[9] はそのファイルシステムで複数のネットワークでまたがった物理ディスクを一つの論理ボリュームに見せる仮想化を行っている。また、Vicom の Service Virtualization Router[7] や DataCore の SAN Symphony[5] は FC ネットワーク上のルータやスイッチにおける In-Band 型のソリューションである。一方、Compaq の VersaStor[4] の様な HBA(Host Bus Adapter) における仮想化技術も見られる。これらの技術は現在限定的な環境下でのみ有効であり、今後 iSCSI の登場によりストレージのネットワーク化が進み、多様な環境で様々なサービスを提供する仮想化が登場すると見られる。

表 1: Design comparison of SAN and LAN[10][11]

| | SAN | LAN |
|---------------|---|------------------------|
| Media | Fibre Channel | Gigabit Ethernet |
| Throughput | 1.063Gbps (2-10Gbps) | 1.25Gbps (-10Gbps) |
| Transfer unit | frame/sequence/exchange (128MB/sequence) | frame (1500B/frame) |
| Error control | Hardware CRC | - |
| Routing | static | dynamic/static |

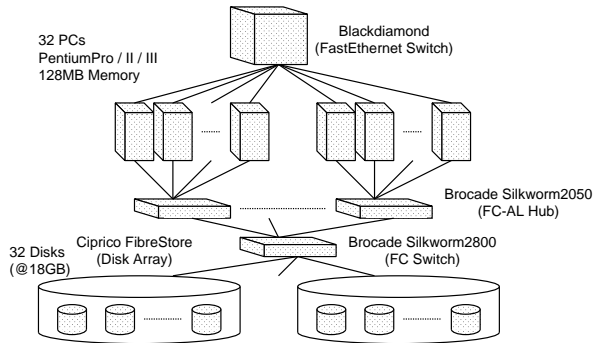


図 1: Overview of SAN-connected PC cluster

表 2: Configuration of PC nodes

| PC ノード共通性能 | |
|--------------------|---------------------|
| Main Memory | 128Mbytes |
| OS | Solaris 8 for x86 |
| FastEthernet NIC | Intel EtherPro/100+ |
| Fibre Channel HBA | Emulex LP8000 |
| PC ノード別性能 | |
| CPU | Chipset |
| Pentium III 800MHz | Intel 440BX AGPset |
| Pentium II 450MHz | Intel 440BX AGPset |
| PentiumPro 200MHz | Intel 440FX |

2.2 SAN 結合 PC クラスタ

実験システムとして 32 ノードの PentiumPro/II/III を搭載した PC および FC ディスクアレイを FastEthernet および Fibre Channel によって接続した SAN 結合 PC クラスタを構築してある。システムの概要を図 1 に、各 PC ノードの所元は表 2 に示す。

3 並列データマイニング

3.1 相関ルール

本資料ではデータベースアプリケーションの一つとして、相関ルール抽出処理を対象とする。相関ルールは以下の様に定義される。アイテム集合を $I = \{i_1, i_2, \dots, i_m\}$ 、トランザクションデータベースを $D = \{t_1, t_2, \dots, t_n\} (t_i \subseteq I)$ とした際に、各要素 t_i をアイテム集合と呼び、 k 個の組合せの物を長さ k のアイテム集合と呼ぶ。相関ルールは $X \Rightarrow Y$ で表現され、 $X, Y \subseteq I, X \cap Y = \emptyset$ を意味する。相関ルールは支持度 $sub(X \Rightarrow Y)$ と確信度 $conf(X \Rightarrow Y)$ の二つのパラメータを持ち、前者は D 全体に対し、 X

と Y をともに持つ確率 $sub(X \cap Y)$ と、後者は更に $sub(X \cap Y)/sub(X)$ と定義される。

相関ルール抽出処理はトランザクションデータベース D に対して支持度 $sub(X \cap Y)$ の最小値および確信度 $sub(X \cap Y)/sub(X)$ の最小値が与えられた際に、これらを満足するルールを見出すことである。この処理は以下の二つのステップによって行われる。

1. 与えられた最小支持度を満たすアイテム集合 (ラージアイテム集合) を全て抽出する。
2. 得られたラージアイテム集合から最小確信度を満たす相関ルールを得る。

この第 2 ステップは限られた個数のルールをフィルタする処理であるため、比較的軽負荷の処理であるのに対し、第 1 ステップは巨大なトランザクションデータベースを繰り返し検索し支持度を調査するため、非常に多くの時間を必要とする演算である。このため、相関ルールの抽出アルゴリズムは第 1 ステップの効率化に焦点を当てている。

3.2 相関ルール抽出処理の並列処理

この相関ルール抽出の代表的なアルゴリズムとしては、IBM アルマデン研究所の Agrawal[1] による Apriori が良く知られている。Apriori では k 個のアイテムの組合せを k -itemset、長さ k のラージアイテム集合を L_k 、長さ k の候補アイテム集合を C_k とする。そこで、 $k(\geq 2)$ に対して以下の処理を行う。

1. L_{k-1} から C_k を作成する。
2. トランザクションデータベース D を検索し、各候補アイテム集合の支持度を計測する。
3. C_k の中から最小支持度を満足するものを取り出し、 L_k とする。

上記で長さ L_k を求める手続きをパス k と呼び、 L_k が空になるまで、繰り返される。

本論文では上記 Apriori をもとに提案された並列アルゴリズム HPA(Hash Partitioned Apriori)[6] を利用する。HPA ではクラスタ内の各ノードに SEND プロセスおよび RECV プロセスの二つのプロセスを配置し、 L_k が空になるまで以下の手順でパス k を繰り返す。

1. SEND プロセスは前パス $k-1$ で得られたラージアイテム集合 L_{k-1} を基に長さ k の候補アイテム集合 C_k を作成し、ネットワークを介してハッシュ分散により RECV プロセスに送信し、ハッシュ表に挿入する。
2. SEND プロセスはディスク上のトランザクションデータベースから長さ k のアイテムの組合せを逐次作成し、ハッシュ分散により RECV プロセスに送信する。RECV プロセスではハッシュ表を用いて候補アイテムを数え上げを行う。
3. トランザクションデータベース全てが検査された後、RECV プロセスは最小支持度を満たすラージアイテム集合 L_k を求め、全 SEND プロセスにブロードキャストする。

上記で計算機に対しては 2. に於ける SEND の長さ k のアイテムの組合せ作成処理と RECV のハッシュ表検索が大きな負荷となりえ、特に最も多くの候補アイテム集合が発生するパス 2 が最も CPU に負荷を与え、マイニング実行時間の大半を占める。対して、パス 3 以降の後半のパスは徐々に CPU への負荷を失い、結果、処理はディスクの I/O 性能に依存する。この性質から、HPA では第一にパス 2 の実行時間を改善を行い、続いて後半のパスのディスク I/O 性能を向上させることが全体の高速化に結びつく。

4 SAN を利用した動的負荷分散機構

4.1 SAN 適用負荷分散システムの概要

HPA アルゴリズムを先述の SAN 結合 PC クラスタシステムで実施するため、従来 Shared Nothing に於いて行われて来た負荷分散機構 [8] をもとに設計を行う。図 2 に本システムの論理構成を示す。

図中マイニングノード (node#[0-2]) の HPA app の部分で実際にマイニング処理を行う SEND/RECV プロセスが動作する。SEND プロセスは SAN 上のトランザクションデータベースにアクセスする際、下層の LVM(Logical Volume Manager) を用いてアクセスを行う。LVM はアプリケーションから呼ばれる API(Application Programming Interface) の形で実

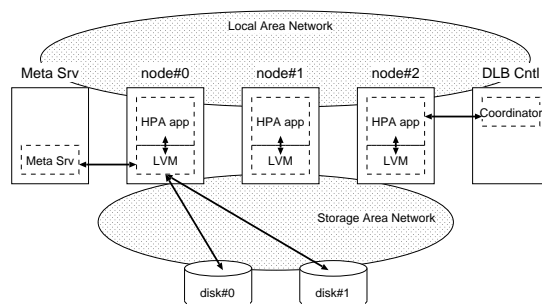


図 2: Overview of load balancing system for HPA on the SAN-connected PC cluster

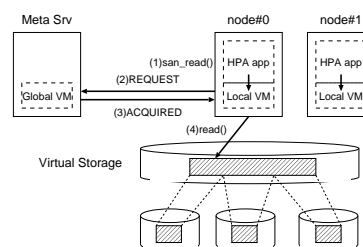


図 3: Disk access procedure with SAN Meta Server

装されており、別ノードに用意する SAN メタサーバ (Meta Srv) と協調する形で、通常のディスクアクセスシステムコール (read() 等) に代わるサービスを提供する。一方、負荷分散コントローラ (DLB Cntl) は Shared Nothing 環境で用いられる物とほぼ同様で、HPA のプロセスと協調して負荷情報を収集し、負荷均衡化の為にマイグレーションプランニングやトリガ実行等の行う。

4.2 SAN メタサーバによる仮想化

SAN 環境に於いてはストレージデバイスが複数のホストコンピュータによって共有される。このため、ホストコンピュータ間でディスクアクセスのメタ情報の一貫性を保持する必要がある。そこで、本研究では先述の SAN メタサーバを用いてメタ情報の一元管理を行い、一貫性の保持を行うとともに、後述のデクラスタリングのようなサービスを行う。

図 3 にその具体的な動作として read() 処理の例を示して説明する。まず、HPA アプリケーションからはじめて san_read() が呼び出されると、LVM はメタサーバに問い合わせを行う。メタサーバは自身の管理しているメタ情報に基づいてプランを立てて、SAN デバイス上の領域をある粒度でマイニングノード用にロックを行い、その領域のデバイス情報やオフセット情報な

どのメタ情報をマイニングノードに通知する。それを受け、マイニングノードのLVMはSANデバイスへのread()を行い、必要なデータを取得し、上位のAPI呼び出し元に受け渡す。2回目以降のsan_read()呼び出しに対しては、自身がロックしている領域が残っている場合はそのままデバイスへのread()を行い、領域が枯渇した場合は改めてメタサーバへの問い合わせを行う。

この機能によりSAN上の複数ディスクのファイルは、複数のホストコンピュータから仮想的に一つのファイルとして共有されることになる。このため、ディスク間のデータ量の偏り、例えば関係データベースに於けるTPS(Table Placement Skew)[3]や実行中のCPU負荷によって発生したデータ量の差異等は、SANの仮想化機構によって吸収することができる。

4.3 Migrationによる負荷分散

並列分散環境におけるアプリケーションの負荷調整はメモリ上の演算データや演算スレッドのMigrationによって行われるが、Shared Nothing環境においてはノード間の実行時間の差異を解消するため、ディスク上のデータ量を考慮する必要がある。しかし、演算データのMigrationによるディスク間データ量の調整はスキューが小さい場合は可能であるが、大きくデータ量がびりびりな場合には必ずしも十分でなかった。[8]ではさらにディスク上のデータをLAN経由で読み出すことでデータスキューを解消する方式(Transaction Migration)を導入していた。しかし、Transaction Migrationを利用してもなお、未利用のCPUノードに処理を拡張するケース等に於いては、ネットワークトラフィックが偏るなどに理由により十分対応しきれていなかった。

本研究ではSANの適用により、負荷分散コントローラはデータ量の差異を考慮する必要がない。そのためCPUバウンド領域に於いて、それぞれの時刻に於ける各ノードのCPU負荷が、相互にボトルネックにならないように調節することで済む。HPAに於いてはRECVプロセスの優先度が高く、特定のノードでRECVプロセスのみがCPU資源を使いきってしまう場合、RECVプロセスは他のノードから送られてくるデータを十分処理しきれず、結果全ノードの性能を低下させていると考えられるからである。

このため、負荷分散コントローラの制御方式は以下

の様に定式化することができる。

ノード*i*に於けるノード全体、RECV, SENDプロセスのCPU使用率をそれぞれ $L_i, L_{RECV_i}, L_{SEND_i}$ とする。 L_{RECV_i} はRECVプロセスが受信するデータ流量(V_{Ri})に比例し、 L_{SEND_i} はSENDプロセスが送信するデータ流量(V_{Si})に比例するため、以下のように表すことができる。

$$\begin{aligned} L_i &= L_{RECV_i} + L_{SEND_i} \\ L_{RECV_i} &= \alpha \gamma_i V_{Ri} \\ &\sim \alpha \gamma_i C_i \sum V_{Si} \\ L_{SEND_i} &= \beta \gamma_i V_{Si} \end{aligned}$$

ここに、 α, β はRECV, SENDの処理の負荷係数、 γ_i はCPU係数、 C_i はRECVプロセスのハッシュテーブル重み係数とする。添字*i*はノード番号を表す。また、

$$L_i \leq 1$$

の条件を満たす必要がある。

この時、系内でRECVプロセスがボトルネックにならないためには全ノードで $L_{RECV_i} < 1 - \epsilon^1$ が満たされるべく制御を行う必要がある。このため、負荷分散コントローラは以下の手続きを定期的実施する。

1. 負荷分散コントローラはマイニングノードから統計情報を収集する。
2. 得られた統計情報(L_{RECV_i})から、各ノードの $\alpha \gamma_i$ を算出する。 $L_{RECV_i} < 1 - \epsilon$ を満たさないノードが存在する場合、マイニングノード間でのハッシュテーブルの再配置(Candidate Migration)が必要であると判断し、 C_i を操作変数として、全ノードの L_{RECV_i} がほぼ均衡化される方向にマイグレーション計画を建てる。
3. 負荷分散コントローラが全マイニングノードにマイグレーション計画を通知し、それを受け取ったマイニングノード間でハッシュラインのマイグレーションを行う。

¹ ϵ は本研究の実験では5%程度取っている。

表 3: CPU placement

| Case | Node #0 | Node #1 | Node #2 | Node #3 |
|------|----------------------|----------------------|----------------------|----------------------|
| c0 | PentiumIII 800MHz | PentiumIII 800MHz | PentiumIII 800MHz | PentiumIII 800MHz |
| c1 | PentiumIII 800MHz | PentiumII 450MHz | PentiumII 450MHz | PentiumII 450MHz |
| c2 | PentiumIII 800MHz | PentiumIII 800MHz | PentiumIII 800MHz | PentiumII 450MHz |
| c3 | PentiumIII 800MHz | PentiumIII 800MHz | PentiumIII 800MHz | PentiumPro 200MHz |

表 4: Data placement

| Case | Disk1 | Disk2 | Disk3 | Disk4 |
|------|---------------------|---------------------|---------------------|------------------------|
| d0 | 1,000,000 (84MB) | 1,000,000 (84MB) | 1,000,000 (84MB) | 1,000,000 (84MB) |
| d1 | 200,000 (16.8MB) | 200,000 (16.8MB) | 200,000 (16.8MB) | 3,400,000 (285.6MB) |

Number of transactions / Size of data volume

4.4 CPU バウンド領域での評価実験

4.4.1 4 ノード非均質環境

上記の負荷分散機構を評価するために、表 3 の PC ノードの組合せ、表 4 のトランザクションデータベース配置を用意した。この時、トランザクションデータベース内のアイテム数は 5000 とし、1 トランザクションあたりの平均アイテム数は 20 とした。また最小支持度は 0.7% とした。ディスクアクセスのバッファサイズは 64KB、SAN メタサーバのロックを行う粒度は 512KB とした。

各 CPU スキューおよびデータスキューの組合せに対し、SAN 環境下で Candidate Migration を無効にした場合、有効にした場合それぞれのケースのパス 2 の実行時間を測定した。表 5 に示す。この時、CPU/ディスク共スキューのない、c0d0 に於ける Shared Nothing で無制御のケースを 100 とし、必ずしも正確な比較にはならないが CPU クロック数を以って正規化を行った。比較のために [8] で述べられている Shared Nothing 環境に於ける負荷制御手法を検証実験した結果を併せて掲載する。表 5 で SAN Enable が SAN 適用の負荷分散、SAN Disable が [8] で述べられている Shared Nothing 下の負荷分散を表す。[8] 方式ではディスク上のデータ量を評価関数に含めて Candidate Migration を行い、解決できない場合は LAN 経由でディスク上のデータの均衡化を図る (Transaction Migration)。このため、SAN、Shanred Nothing 間で Candidate Migration の定義は異なる。† は Candidate Migration、‡ が Transaction

表 5: Normalized execution time of Pass 2

| Case | SAN | Contorol | d0 | d1 |
|------|---------|----------------------------|----------------------|--------------------|
| c0 | Disable | - Cand. Cand.+Trans. | 100 100† 100†‡ | 170 136 136‡ |
| | Enable | - Cand. | 98 98† | 102 102† |
| c1 | Disable | - Cand. Cand.+Trans. | 100 94 94‡ | 173 126 126 |
| | Enable | - Cand. | 93 93† | 95 95† |
| c2 | Disable | - Cand. Cand.+Trans. | 129 104 104 | 223 158 155 |
| | Enable | - Cand. | 108 97 | 112 100 |
| c3 | Disable | - Cand. Cand.+Trans. | 220 121 113 | 395 270 259 |
| | Enable | - Cand. | 175 98 | 176 99 |

Migration のトリガが引かれなかったことを表す。

データスキューに関しては、d0 と d1 の比較から、SAN の適用により実行時間の悪化を回避していることが分かる。つまり、SAN の適用によってデータスキューは容易に解決できることが分かる。一方、CPU のスキューに関しては、上記中では c2,3 のケースが、一台の低性能ノードが全体のボトルネックとなる苛酷なケースであるが、SAN 環境下では Candidate Migration と併せて適切に負荷制御が行われていることが分かる。

また、c1 のような緩やかな CPU スキューに関しては、Candidate Migration の必要がなくとも負荷制御が可能であることが分かる。c1d0 で SAN 下で Candidate Migration を行わないケースの実行トレースを図 4 に示す。この図では下から Node [0-3] の各リソースを表し、太線は SEND のディスク read スループット、点線は SEND の CPU 利用率 L_{SEND_i} 、実線は SEND と RECV の CPU 利用率 $L_{SEND_i} + L_{RECV_i}$ を表す。時刻約 10-300 秒がパス 2 である。ハッシュテーブルは等量分散されているので、相対的に速い CPU を持つ Node 0 の L_{RECV_i} は少なくなるが、余剰 CPU リソースを使用して SEND が多くのデータ処理を行っていることが分かる。Node [1-3] の低速ノードでは CPU のほとんどを RECV 処理が使用しているが、独占はしていないので Node 0 に対してボトルネックにならず、全ノードの CPU が有効に利用されている。c1 のようなケースに於いてはアプリケーションが専用の負荷分散機構を有しない場合でも SAN のみによって自動制

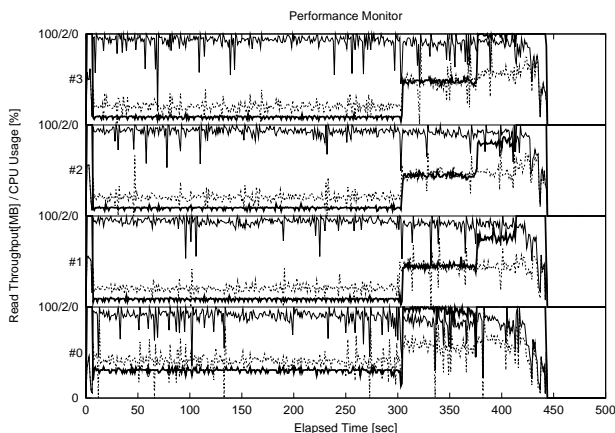
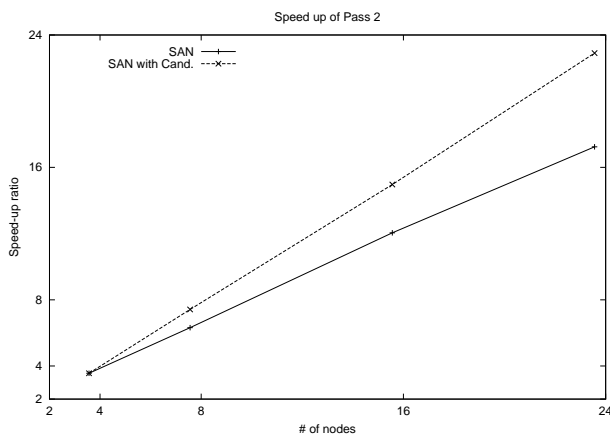


図 4: Execution trace (Case:c1d0, Without Candidate Migration, On SAN)



PentiumIII(800MHz)*N 台+PentiumII(450MHz)*1 台の組合せ

図 5: Speed-up curve of Pass 2

御が可能である。

4.4.2 スケーラビリティ

次にスケーラビリティを調査するために、先の d0 のディスク配置の元で、「PentiumIII(800MHz) * N 台 + PentiumII(450MHz) * 1 台」の組合せを用意し、パス 2 の実行時間を計測した。スピードアップ曲線を図 5 に示す。この図に於いては横軸のノード数は PentiumIII(800MHz) を 1 とし、クロック周波数で正規化を行ったものを用いている。

この図から前節の負荷分散手法は非常によいスケーラビリティを実現し、多ノード環境に於いても有効であることが分かる。

5 動的資源割り当て

5.1 実行時 CPU 投入

PC クラスタは多数のノードを有しているが、アプリケーションの要請によって必要なリソースを適宜決定する必要がある。例えば HPA アプリケーションに於いてはパス 2 が非常に多くの CPU 資源を消費し、実行時間全体に対して支配的である。このようなアプリケーションに対して適切なノード数を割り当てる機能が必要になってくる。

従来 Shared Nothing 環境ではデータがノードに依存していたため、ノード数の変更に関してはディスク上のデータの再配置や [8] に於ける Transaction Migration のように非常にコストのかかる制御が必要であった。一方、SAN 環境ではストレージと CPU ノードは独立であるため、柔軟な CPU の資源投入が可能である。ここでは、先述の SAN 適用の負荷分散機構を利用して、実行時に追加的に CPU ノードを投入し、CPU バウンド時の実行時間を短縮させる制御を行う。

5.1.1 制御方式

実行時 CPU 投入機構は負荷分散コントローラの機能として設計する。まず、アプリケーション開始時には利用が想定される全てのノードにアプリケーション駆動の為にプロセスを生成しておく。そのうち当初は少数ノードのみが HPA アプリケーションの処理を行う。その後は以下の手続きを繰り返す。

1. 負荷分散コントローラはマイニングノードから統計情報を収集する。
2. 得られた統計情報から、系全体が CPU バウンドしていると判断され、投入可能な CPU ノードが存在する場合、新規投入ノードに投入処理を開始するよう通知する。
3. 通知を受けたノードは直ちに、SAN メタサーバを介してディスクをマウントし、SEND プロセスが処理を開始する。
4. 負荷分散コントローラは直ちに Candidate Migration を利用して、新規投入ノードにハッシュラインの配布を試みる。この時、新規投入ノードでは

表 6: CPU placement

| Node #[0-1] | Node #2 | Node #[3-5] | Node #[6-8] |
|----------------------|----------------------|---------------------|----------------------|
| PentiumIII 800MHz | PentiumPro 200MHz | PentiumII 450MHz | PentiumIII 800MHz |

表 7: Data placement

| Disk1 |
|-----------------------|
| 16,000,000 (1.3GB) |

Number of transactions / Size of data volume

$C_i = 0$ であるため、統計情報 (L_{RECV_i}) から $\gamma_i\alpha$ が算出できず、マイグレーション時の C_i 配布量が決定できない。そのため、ごく少数のハッシュラインを配布し $\gamma_i\alpha$ 値が得られるのを待つ。

5. $\gamma_i\alpha$ 値が得られたのち、改めて Candidate Migration を利用して負荷の調整を行う。

5.1.2 評価実験

実行時 CPU 投入機構の性能評価を行うために、表 6 の種々の PC が混在する PC の組み合わせ、表 7 のデータ配置を用意した。初期ノード数は 1(Node #0) としてスタートした。各時点で投入されるノード数は運用ノード数と等量ずつとした。

図 6 に実行トレースを示す。パス 1 が約 55 秒程度行われた後、パス 2 が開始し、システムは CPU バウンドを判断、直ちに 1, 2, 4 とノードを追加していることが分かる。新規ノード投入後は直ちに SEND が処理を開始するとともに、Candidate Migration によって RECV にハッシュラインが配られ、負荷が均衡化している。追加ノードを投入しない 1 ノードの場合はパス 2 の実行時間は 3145 秒かかっていたが、動的投入により 580 秒にまで改善し、5.5 倍のスピードアップの結果となった。CPU 拡張による単純クロック比は 6.9 倍であり、逐次投入している期間の存在を考慮すると、PentiumPro 200MHz から PentiumIII 800MHz が混在する非均質環境に於いて、動的 CPU 投入処理が有効に機能していることが分かる。

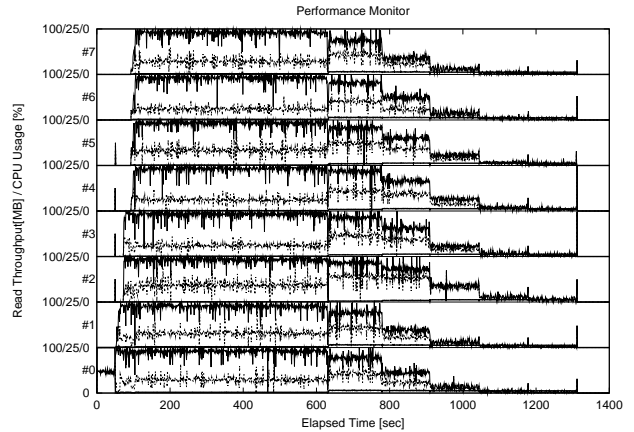


図 6: Execution trace (Dynamic resource expansion up to 8CPUs, Minimum support:0.7%)

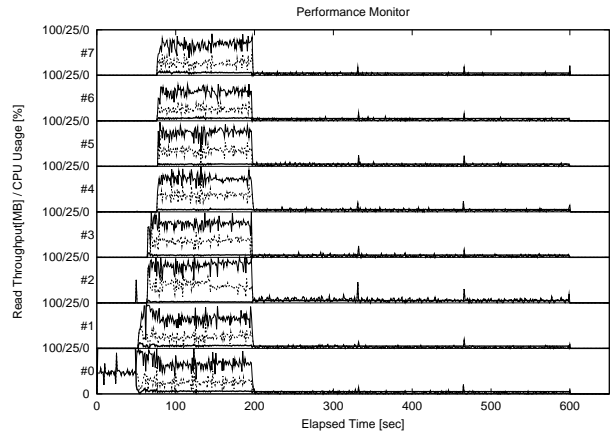


図 7: Execution trace (Dynamic resource expansion up to 8CPUs, Minimum support:1.6%)

5.2 ディスク I/O バウンドと動的デクラスタリング

以上は系が CPU バウンドしている領域に限った議論であった。例えば、前節の実行時 CPU 投入を利用すると CPU 数の増大に伴い、系がディスク I/O バウンド化する状況が存在する。前節の実験に於いても、最小支持度を 1.6%と変化させると結果は図 7 の様になる。ノード数が 8 に達した時点でパス 2 がディスク I/O バウンド化し、そのため CPU が有効に利用されていない。また、パス 2 の実行時間の改善により、パス 3 以降のディスク I/O バウンドのパスが実行時間全体に与える影響も無視できなくなる。

本節ではこのようなケースのディスク I/O バウンド問題の解決のため、動的デクラスタリング機構 [13] を SAN メタサーバの機能として加える。これはデータが

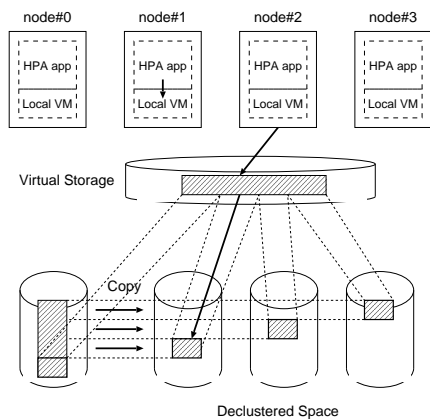


図 8: Overview of dynamic declustering on SAN

存在しているストレージデバイスからデータを分割して未使用のストレージ空間にコピーを行い、並列にアクセスを行うことで、ディスク I/O の帯域を拡張するもので、HPA の様な繰り返しストレージデバイスをシーケンシャルアクセスするアプリケーションに非常に有効である。図 8 に動的デクラスタリングの概念図を示す。

そこで、動的デクラスタリング機構を SAN メタサーバおよび LVM の機能として追加を行った。その制御方式について説明する。

1. パス 1 に於いては、SAN メタサーバから LVM に対して ACQUIRED メッセージを通知する際に、デクラスタを作成する旨 (DODECLUSTER) とコピーの作成先メタ情報を併せて送信する。
2. DODECLUSTER メッセージを受けた SEND プロセス (LVM) はコピー元デバイスからトランザクションデータベースを読み込み、処理を行うとともに未使用デバイスにデクラスタコピーを作成する。
3. パス 2 開始以後、SEND プロセス (LVM) はノードのバウンディングファクタ (CPU バウンドかディスク I/O バウンドか) を判断し、SAN メタサーバに通知する。通知を受け SAN メタサーバは、系全体がディスク I/O バウンドと判断されると、自身の管理するメタ情報を更新してコピー元から Declustered Space へと切替えを行う。
4. この切替えにより、次回の REQUEST 以降、マイニングノードの Declustered Space へのアクセスが活性化する。

5.3 バウンディングスイッチによる動的資源投入

5.3.1 バウンディングスイッチと資源

実行時 CPU 投入機構および動的デクラスタリングを用いることで、系がどの資源の性能にバウンドした状態であるかを適切に判断して、資源の拡張を行うことが可能になる。つまり、CPU バウンド状態にある時には実行時 CPU 投入による CPU 資源の拡張を、ディスク I/O バウンド時には動的デクラスタリングによる帯域拡張を行うことが可能である。

5.3.2 動的資源投入の評価実験

実行時 CPU 投入機構および動的デクラスタリング双方を利用した動的資源投入の実験を行う。実験には同じく表 6、表 7 に示すものを用いたが、CPU ノード数は 800MHz のものをさらに 8 台加え 16 台まで拡大し、未使用ディスク 15 台を加え 16 台のディスクを用意した。また、最小支持度は 1.6%とした。

実験結果の実行トレースを図 9 に示す。CPU のみを拡張する図 7 のケースでは CPU を 8 台まで投入した時点で系がディスク I/O バウンド化していたが、図 9 では動的デクラスタリングを適用することで、直ちに SAN メタサーバにより Declustered Space へのアクセスの活性化が行われる。このため、I/O 帯域が拡大し再び CPU バウンド化している。そこで、再び負荷分散コントローラが CPU の動的投入を指示するため、最終的に 16 ノードまで拡大するとともに、パス 3 以降のディスク I/O バウンド時の性能改善から図 7 と比較して実行時間が 610 秒から 169 秒と大幅に短縮されていることが分かる。

図 10 には動的拡張前の 1CPU/1DISK のケースと、CPU およびディスク双方 16 台まで拡張した図 9 のケースの各パスの実行時間の変化を示す。

このように、独立した二つの機構を組み合わせることで CPU およびディスク I/O 双方の性能限界を検知し、必要な資源を投入することができる。

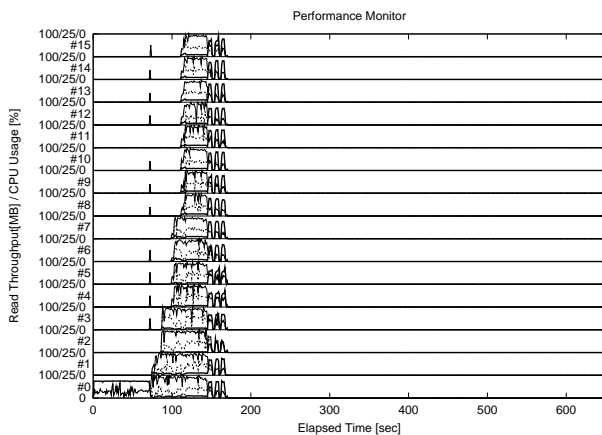


図 9: Execution trace (Dynamic resource expansion up to 16CPUs and 16DISKS, Minimum support:1.6%)

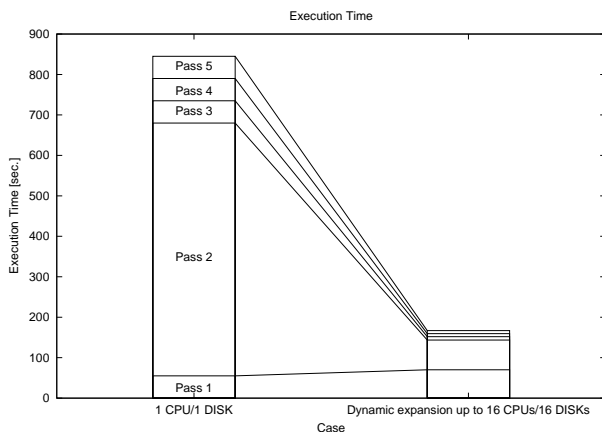


図 10: Comparison of execution time

6 結論

SAN を適用した PC クラスタに於ける相関ルール抽出処理に関して、動的負荷分散機構の設計を行い、実装実験によりその有効性を確認した。データスキューに関しては SAN の仮想化機構によりほぼその弊害を除外することができた。また、苛酷な CPU スキューに関しても Shared Nothing に比べ実行時間が改善し、多ノード環境でも有効であることが分かった。加えて、緩やかな CPU スキューに関しては Candidate Migration が不要で、SAN のみで解決できることも分かった。

その上で負荷分散機構に、動的 CPU 投入と動的デクラスタリングを加えて実装し、非均質な PC の組み合わせ環境の元で実験を行った。動的 CPU 投入は負荷分散コントローラ、動的デクラスタリングは SAN メタサーバによってそれぞれ独立に制御されるが、双方とも有効に機能することが実験によって確認された。

今後は SAN のみで解決できるスキューの範囲を明確にした上で、その効果を検証する予定である。SAN のみでスキューが解決できることはアプリケーション毎に開発する必要のある負荷分散機構を共通の下位レイヤで制御可能にするものであり、数多くのデータインテンシブアプリケーションに対し、アプリケーションの書き換え無しに適応可能と考えられる。また、本論文で提案した Candidate Migration はアプリケーションからの知識が必要となるが、その修正もそれほど大きくなく、アプリケーションへの変更を最小化するツールに関して今後研究を進めていきたい。

参考文献

- [1] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In *Proceedings of the Twentieth International Conference on Very Large Data Bases*, 1994.
- [2] Charles T. Clark. *The Virtualization of Storage*. White Papers, TidalWire <http://www.tidalwire.com/>, 2001.
- [3] R. M. J. Christopher B. Walton, Alferd G. Dale. A taxonomy and performance model of data skew effects in parallel joins. In *Proceedings of the Seventeenth International Conference on Very Large Data Bases*, 1991.
- [4] Compaq. Compaq Innovation to Transform Enterprise Storage Deployment, Utilization and Management.
- [5] SAN Symphony. DataCore Software, <http://www.datacore.com/>.
- [6] T. Shintani and M. Kitsuregawa. Hash based parallel algorithm for mining association rules. In *Proceedings of Parallel and Distributed Information Systems*, 1996.
- [7] SV Router. Vicom Systemns, http://www.vicom.com/pdf/bro_svrouters.pdf.
- [8] M. Tamura and M. Kitsuregawa. Dynamic load balancing for parallel association rule mining on heterogeneous pc cluster systems. In *Proceedings of the Twenty-fifth International Conference on Very Large Data Bases*, 1999.
- [9] VERITAS Software. <http://www.veritas.com/>.
- [10] コンピュータエージ社. データストレージレポート 2001. 月刊コンピュータピア, 2001.
- [11] ファイバチャンネル協議会. ファイバチャンネル技術解説書. 論創社, 2001.
- [12] 安井, 田村, 小口, 喜連川. 並列 DBMS に於ける動的負荷分散機構の実装. Number 61 in 99, pages 393-398. 情報処理学会, 1999.
- [13] 小口, 喜連川. SAN 統合 PC クラスタ上の並列データマイニングのための動的データ・デクラスタリング. In 情報処理学会データベースシステム研究報告, 2001.