

計算機ネットワークに自律的に適応する分散システムの検討

谷江健 角川裕次 藤田聡 山下雅史
広島大学工学部

分散システムの最適な並列性や実行形態は、分散システムが実行される計算機ネットワークの状態に大きく依存する。本研究は、計算機ネットワークの状態に依存して、その並列度や実行形態を自律的に変える分散システムの可能性について議論する。

A STUDY OF DISTRIBUTED SYSTEM AUTONOMOUSLY ADAPTED TO COMPUTER NETWORKS

Takeshi Tanie, Hirotsugu Kakugawa, Satoshi Fujita, Masafumi Yamashita
Faculty of Engineering, Hiroshima University

Optimal parallelism and execution form of distributed systems depend on the state of computer networks. This paper suggests potentialities of the distributed systems which changes it's parallelism and execution form autonomously.

1 はじめに

複数の計算機を、ネットワークでつないで相互に通信できるような計算機ネットワークが普及してきた。このためユーザーはネットワークを通して手軽に他の計算機資源を利用することが出来るようになってきている。そこで、複数の計算機を用いて問題を分散して解くことで計算の応答時間を短縮することができる。我々は、処理手順は同一でデータのみを分割するアルゴリズムである分割統治法に基づき分散化を行うことにする。

その際に注意しなければならない点が二つ挙げられる。一つは、通信遅延によるオーバーヘッドのため問題を分散したからといって計算の応答時間が短くなるとは限らないことである。二つ目は、計算機資源を占有できないことである。これは、マルチタスクマルチユーザーシステムでの計算機プロセッサは複数のユーザと共有して使われるため計算機の負荷が重くなり本来の計算能力より劣った計算機を使用することになる。したがって、負荷の状態によりアルゴリズムの動作を変えることで応答時間を短縮できるかも知れない。この二点より単純に並列度を上げれば計算の応答時間が短縮されるわけではないことがわかる。シーケンシャル実行が最も応答時間が短い場合もあるのである。それらは計算機ネットワークの状態に依存するので分散の仕方をユーザーが決定するのは困難であるし手がかかる。本稿では、アプリケーション自身が計算機ネットワークに適応して自分自身の動作を変えることで自律的に最適な分散システムを構築する可能性を検討する。

自律的に最適な並列性を獲得するアルゴリズムの評価のために 10Mbps の Ethernet でつながれた同機能同性能な計算機群を用いてマージソートを行った。2章では、並列実行時間がインスタンスサイズと計算機負荷により決定できることを示す。3章では、自律的に並列度を変えるためのアルゴリズムを示す。4章で、そのアルゴリズムを実験的に評価する。5章で、本研究のまとめを述べる。

2 並列実行時間

分割統治法によるアルゴリズムの実行時間が計算機の負荷とインスタンスサイズに依存することを示す。マージソートを例題として、その実行時間を実

験式と理論式から見積もる。 n をインスタンスサイズ、 l を計算機負荷、 p を計算機台数として、ソートの実行時間を $T_{sort}(n, l, p)$ で表す。また、ソートアルゴリズム中のマージ操作時間を $T_{merge}(n, l)$ 、通信時間を $T_{comm}(n, l)$ 、そして部分問題の遠隔起動時間を $T_{rezeq}(l)$ とする。マージソートの計算構造は2分木で表せる。2分木の各階層のマージ操作の計算複雑度は $\frac{1}{2}n$ であり、2分木の高さは $\lg n$ だからシーケンシャルマージソートの計算複雑度は $\frac{1}{2}n \lceil \lg n \rceil$ である。

2.1 モデル化

並列度 4 のマージソートを図 1 のような計算構造で行う。ノード中の番号はホスト番号を表しており、親と子が同じ番号ならばローカルな再帰呼出しを表す。そうでなければリモートな再帰呼出しである。並列度 2 では、図 1 の深さ 1 までの 2 分木となる。チャート図を図 2 に示す。 $t_0 \sim t_1$ は $T_{rezeq}(l)$ 、 t_1

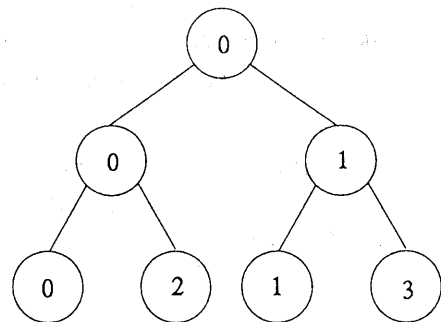


図 1: 計算構造と計算機との対応

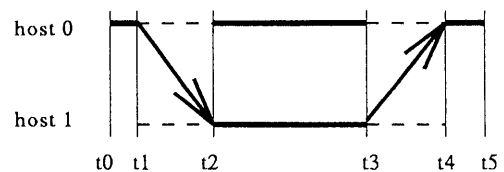


図 2: 並列度 2 のチャート図

$\sim t_2$ は $T_{comm}(\frac{n}{2}, l)$ 、 $t_2 \sim t_3$ は $T_{sort}(\frac{n}{2}, l, 1)$ 、 $t_3 \sim t_4$ は $T_{comm}(\frac{n}{2}, l)$ 、 $t_4 \sim t_5$ は $T_{merge}(\frac{n}{2}, l)$ である。チャー

ト図 (図 2) から $T_{\text{sort}}(n, l, 2)$ を導くと次式となる。

$$T_{\text{sort}}(n, l, 2) = T_{\text{sort}}\left(\frac{n}{2}, l, 1\right) + 2T_{\text{comm}}\left(\frac{n}{2}, l\right) + T_{\text{recec}}(l) + T_{\text{merge}}\left(\frac{n}{2}, l\right)$$

イーサネット上での通信は 1 本のバスを計算機群が共有して使うため同時に複数の通信はできない。したがって、並列度 4 のときのチャート図は図 3 の様になる。 $t_0 \sim t_1$ は $T_{\text{recec}}(l)$ 、 $t_1 \sim t_2$ は $T_{\text{comm}}\left(\frac{n}{2}, l\right)$ 、 $t_2 \sim t_3$ は $T_{\text{recec}}(l)$ 、 $t_3 \sim t_4$ と $t_4 \sim t_5$ は $T_{\text{comm}}\left(\frac{n}{4}, l\right)$ 、 $t_5 \sim t_6$ は $T_{\text{sort}}\left(\frac{n}{4}, l, 1\right)$ 、 $t_6 \sim t_7$ は $T_{\text{comm}}\left(\frac{n}{4}, l\right)$ 、 $t_7 \sim t_8$ は $T_{\text{merge}}\left(\frac{n}{4}, l\right)$ 、 $t_8 \sim t_9$ は $T_{\text{comm}}\left(\frac{n}{2}, l\right)$ 、 $t_9 \sim t_{10}$ は $T_{\text{merge}}\left(\frac{n}{2}, l\right)$ である。チャート図 (図 3) から $T_{\text{sort}}(n, l, 4)$ を導くと次式となる。

$$T_{\text{sort}}(n, l, 4) = T_{\text{sort}}\left(\frac{n}{4}, l, 1\right) + 2T_{\text{comm}}\left(\frac{n}{2}, l\right) + 3T_{\text{comm}}\left(\frac{n}{4}, l\right) + 2T_{\text{recec}}(l) + T_{\text{merge}}\left(\frac{n}{4}, l\right) + T_{\text{merge}}\left(\frac{n}{2}, l\right)$$

2.2 実験式

実験によりシーケンシャル実行時間、マージ操作時間、通信時間、リモート起動時間を求めて次のような実験式を作成した。

- シーケンシャル実行時間

$$T_{\text{sort}}(n, l, 1) = 1.35(l + 1)n \lg n \text{ [}\mu\text{sec]}$$

- マージ操作時間

$$T_{\text{merge}}(n, l) = T_{\text{sort}}(n, l, 1) / \lg n \text{ [}\mu\text{sec]}$$

- 通信時間

計算機負荷に比例する定数と通信媒体のみに関係する定数が送信すべきインスタンスサイズに比例すると次のような式となる。定数は実験により決定した。図 4 がその実験結果である。

$$T_{\text{comm}}(n, l) = (2.355l + 5.846)n \text{ [}\mu\text{sec]}$$

- リモート起動時間

$$T_{\text{recec}}(l) = 0.6(l + 1) \text{ [sec]}$$

2.3 結果

図 5 の実線は各々実際のマージソートを実行してその応答時間を計測したものである。そして破線は先のモデル式をプロットしたものである。同様に図 6 をみても破線は実線とほぼ一致しておりモデル化が正しいことを示している。しかしながら、インスタンスサイズが大きければ通信は最大パケット長の繰り返し送信なので問題ないが、インスタンスサイズが小さいと可変パケット長となるのでばらつきが大きくなりモデルどおりには行かない。

3 アルゴリズム動作

3.1 基本的な考え方

アルゴリズムにサイズ n のインスタンスを与えるとする。アルゴリズムはまず初めにこのインスタンスを半分に分けて並列度 2 で実行を試みる。具体的には $\frac{n}{2}$ を関数の再帰呼出し (ローカルな再帰呼出し) によって割り当て、残りの $\frac{n}{2}$ をリモートホスト上に割り当てて実行 (リモートな再帰呼出し) する。そのとき分散したインスタンスサイズ $\frac{n}{2}$ をスレッシュホールド (θ とする) として記憶しておく。ここで、シーケンシャルな実行時間 (T_{seq} とする) が分かっているならば、並列度 2 のパラレル実行時間 (T_{par} とする) は自分自身の実行時間を計測することで得られるから、この二つの実行時間を比較することでサイズ n のインスタンスの並列化の効果を判定できる。もし並列化の効果がなければ、次回の実行からはこのインスタンスサイズではシーケンシャル実行をすれば良い。逆にパラレル実行時間の方が短ければ、並列化の効果があるので次回の実行では更に並列度を上げて実行してみる。

1. もし ($T_{\text{seq}} \leq T_{\text{par}}$) ならば θ を大きくして集中化傾向を強める。
2. もし ($T_{\text{seq}} > T_{\text{par}}$) ならば θ を小さくして分散化傾向を強める。

3.2 スレッシュホールドの設定と更新

実行時に与えられるインスタンスサイズを n とする。初めは並列度 2 で実行する必要があるため θ の初期値は $\frac{n}{2}$ とする。また、直前の実行時に並列化の効果があったかどうか覚えておくために $lower$ と $upper$ の 2 変数を導入する。この 2 変数の初期値は負値とする。それ以降の実行では、以下のように θ を設定する。アプリケーション起動後、 θ を設定するために $lower$ と $upper$ を保持しているファイルを読み込んで次を実行する。

case 1: ($lower \geq 0$) かつ ($upper < 0$) ならば

$$\theta \leftarrow 2 \times \theta$$

case 2: ($lower < 0$) かつ ($upper \geq 0$) ならば

$$\theta \leftarrow \frac{1}{2} \times \theta$$

case 3: ($lower \geq 0$) かつ ($upper \geq 0$) ならば

$$\theta \leftarrow \frac{1}{2}(lower + upper)$$

プログラムの再帰呼出し部分を θ を閾値として用いて再帰呼出しをローカルかリモートにするかを決定する。そしてアプリケーション実行後に次のようにして二分探索の要領で最適な θ の範囲を決定する。もしシーケンシャル実行時間がパラレル実行時間より短ければ θ はもっと大きくてもいいのでこの θ の値を $lower$ に保存する。逆に、パラレル実行時間の方が短ければ、 θ は分散化傾向を強めるため小さくすべきなので $upper$ に保存する。但し、シーケンシャル動作しかなかったときは何も更新せずに終了する。次回実行時に読み出して θ を設定するために更新した $lower$ と $upper$ はファイルに保存しておく。

3.3 シーケンシャル実行時間の推定

シーケンシャル実行時間を推定するためにプログラミングユーザにインスタンスサイズ n におけるそのアルゴリズムの計算量を n の関数の形で記述してもらう。与えられた関数を $f(n)$ 、インスタンスサイズを n として並列度2の実行が行われたときに同じインスタンスサイズのシーケンシャル実行時間を推定する。この時、子供のノードは $n/2$ のインスタンスサイズのシーケンシャル実行を行っていると思わせる。アルゴリズムがこの子供の実行時間を自分自身で計測すれば、 $f(n/2)$ に対応する実行時間 $T(n/2)$ が得られる。インスタンスサイズ n のときの実行時間 $T(n)$ は、

$$T(n) = \frac{f(n)}{f(n/2)} \times T(n/2)$$

として計算により求められる。

4 実験

4.1 実験環境

10MbpsのEthernetでつながれた同機能同性能なワークステーション4台を用いて実験を行った。機種は、Sun Microsystems社のAS4080/30EGXワークステーション(CPU:Super SPARC, クロック周波数:36MHz, cache:36KB, main memory:36MB)である。記述言語はC言語を用いた。

4.2 実験方法

整数のソーティングをマージソートにより行う。インスタンスサイズ n を10万、20万、30万、40万、50万まで変化させる。そして計算機の負荷を4台とも同程度となるように調節して、負荷値 l を0、1、2、3、4、5まで変化させた。それぞれシーケンシャル実行時間とパラレル実行時間(並列度2と4)および本研究で提案するオートノマス実行時間を測定する。実験は10回繰り返して、その平均値をグラフにプロットした。

4.3 実験結果

実行時間が計算機の負荷とインスタンスサイズより決まると仮定すると、 l が2のときは図7より、インスタンスサイズが0~16万まではシーケンシャル動作が最適であり、15万~40万では並列度2のパラレル動作が最適であり、40万より大きな所では並列度4のパラレル動作が最適であることが分かる。図9と図10はオートノマス実行が自律的に最適な並列度を獲得している様子を表している。図9では1回目の実行(並列度2)で並列化の効果があつたので、2回目の実行では並列度を更に上げて並列度4で実行している。その結果、3回目の実行で最適な並列度は2であることが分かり、それ以降は並列度2で実行している。図8は、図9の実行時のスレッシュホールドの動きを表しています。図7で、シーケンシャルと並列度2のグラフがクロスするインスタンスサイズはおおよそ16万なので、スレッシュホールドは8万程度となるはずである。図8では、スレッシュホールドは2回目の実行で5万から10万の範囲となつて、5回目の実行ではおおよそ7万に収束している。

5 おわりに

本研究によって計算機ネットワークの状態により最適な並列度が存在することがわかり、アルゴリズムが自律的に変化するシステムの有効性を確認した。今後は計算機負荷の動的変動に対する適応について研究する予定である。

参考文献

- [1] H.E.Bal, J.G.Steiner, A.S.Tanenbaum, "Programming Languages for Distributed Computing Systems," ACM Computing Surveys, Vol.21, no.3, Sept, 1989.
- [2] K.Schwan, R.Ramnath, S.Vasudevan, and D.Ogle, "A Language and System for the Construction and Tuning of Parallel Programs," IEEE Trans. Software Eng., Vol.14, no.4, pp.455-471, Apr, 1988.
- [3] 前川 守, 所 真理雄, 清水 謙多郎, "分散オペレーティングシステム UNIX の次にくるもの," 共立出版株式会社, 1991.

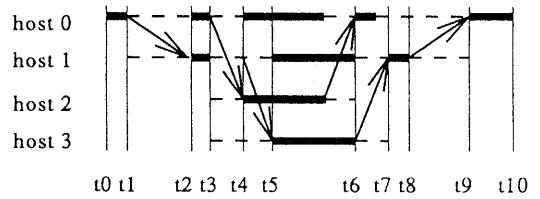


図 3: 並列度 4 のチャート図

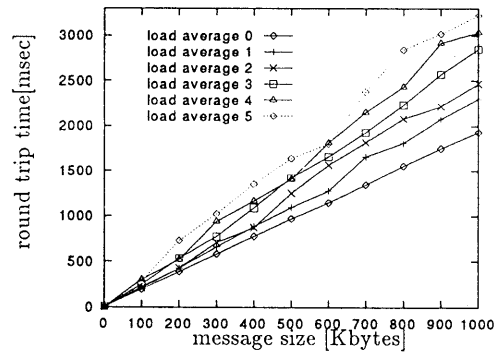


図 4: 通信時間

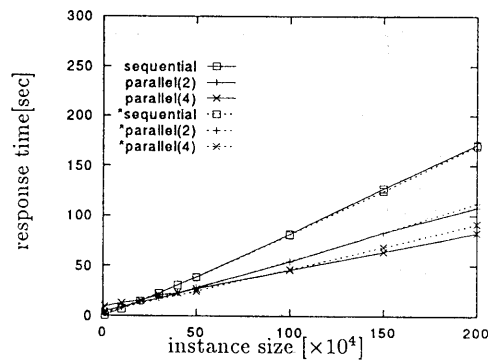


図 5: $l = 2$

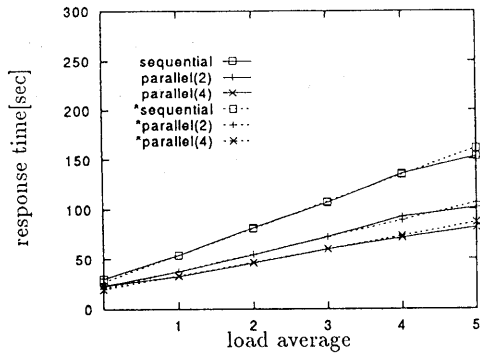


図 6: $n = 1000000$

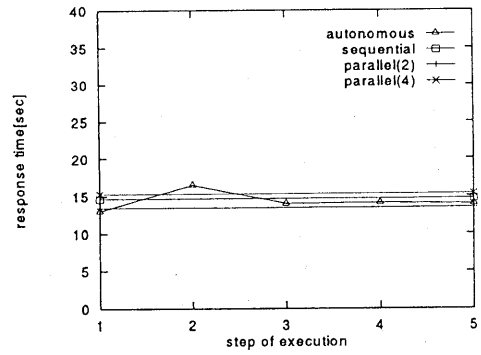


図 9: オートノマス実行 ($n = 200000, l = 2$)

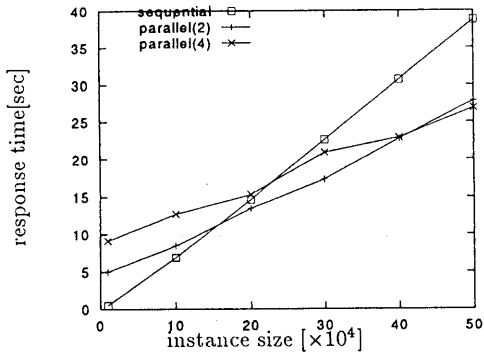


図 7: $l = 2$

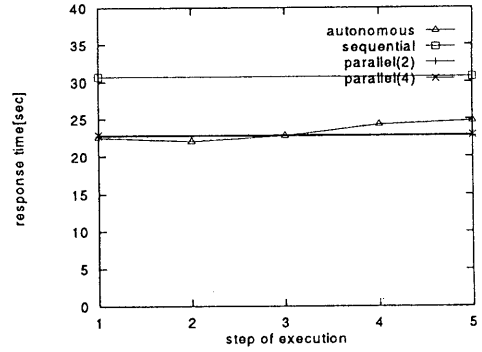


図 10: オートノマス実行 ($n = 400000, l = 2$)

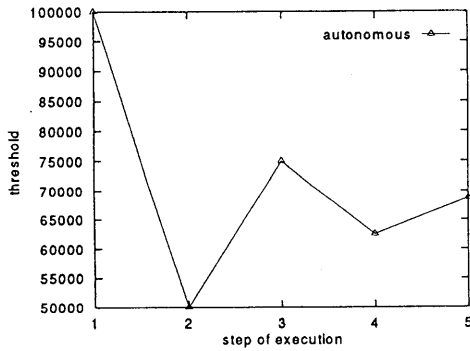


図 8: スレッシュヨルド変化 ($n = 200000, l = 2$)

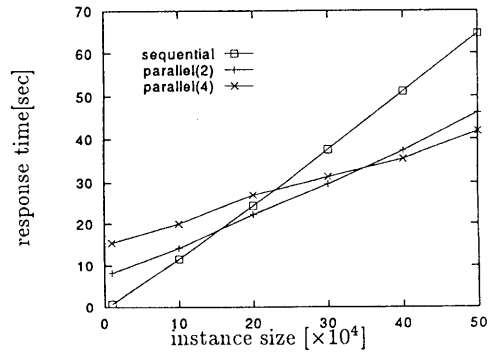


図 11: $l = 4$