

アクセス計算量：新しい並列計算量の枠組みの提案

横山 大作[†] 近山 隆[†]

近年の計算機ハードウェアの進化により、既存の計算量理論ではアルゴリズムの解析に十分とはいえない局面が多くなってきている。計算機のメモリ階層は深化し、RAMモデルと現実とは著しく乖離している。また、クラスタや大規模分散計算など、各計算機間の通信遅延やその違いを無視できない並列環境が一般的になっているが、現状の並列計算コストモデルは、通信遅延の差異を考慮しないものや特定のネットワークポロジに特化したものしか存在しない。我々は、単一計算機内のメモリ階層から計算機間のネットワーク遅延の差異までを統一的に記述できる計算量モデル「アクセス計算量モデル」を提案した。このモデルは、計算のコストの本質は演算ではなく通信にこそ存在するという立場をとる。このモデルは十分簡潔なものであり、並列アルゴリズムの計算量を解析的に求めることができることを、いくつかのアルゴリズムで実証することができた。また、bitonic sort アルゴリズムの実計算機上での実行と比較することで、このモデルの妥当性を示すことができた。

Access Complexity: A New Framework for Complexity of Parallel Computation

DAISAKU YOKOYAMA[†] and TAKASHI CHIKAYAMA[†]

Recent advances in computer hardware have been making existing computational complexity theory inappropriate for many cases. The random access memory (RAM) model was made unrealistic by large speed gap between the processing units and main memory systems. Distributed computing environments have obsoleted traditional models for parallel computation due to non-negligible diversity in communication delay. In this paper, we propose a new framework for computational complexity, named *access complexity*, in which the cost lies in data transfer rather than computation itself. The model tries to capture all levels of system hierarchy, from cache systems to globally distributed environments. It models these diverse access costs in a simple and uniform way. We apply the model to analyze some parallel algorithms, to show that the model can analyze well-known algorithms easily. We also show the appropriateness of the model, through comparing the predicted and the measured performance of bitonic sort algorithm.

1. はじめに

あるアルゴリズムの計算コスト解析には、すべてのメモリが単位時間でアクセスできるとするRAMモデル¹⁾が多数の場合で使われ続けてきた。しかし、近年の計算機アーキテクチャは、少量の速いメモリと大量の遅いメモリが存在するという階層構造を持ち、しかもそのメモリ階層は年々深くなり続けている。このため、モデルでの計算コスト見積りと、現実の計算機での実行結果との乖離が大きなものとなってしまう。RAMモデルはもはや適用できないのが現実である。

たとえば、図1はSPARC III 1.2GHz, 8CPUの計算機において配列の要素ヘランダムにアクセスする

際の1要素あたりのアクセス時間を測定したものである。横軸が要素数のlogをとったものであり、系列の詳細は後述する。図から分かるとおり、要素数が増えたとアクセス時間は1,000倍近くまで増大する。この振舞いは、とてもRAMモデルでは表現できない。

これに対し、メモリへのアクセスコストが一律ではないとするモデルがいくつか提案されてきた。Hierarchical Memory Model (HMM)²⁾やUniform Memory Hierarchy Model (UMH)³⁾がそれである。HMMは、メモリのアクセスコストがアドレス x に対して $f(x)$ で与えられるというものであり、メモリの中でアドレスの小さい部分はCPUに近く、速くアクセスすることができるが、大量のメモリを使おうと思うと次第にCPUから遠い、遅いメモリを使うようになる、ということ表現できるようにしたモデルである。UMHは、大きさやレイテンシが再帰的に決められた一連の

[†] 東京大学新領域創成科学研究科

Graduate School of Frontier Science, the University of Tokyo

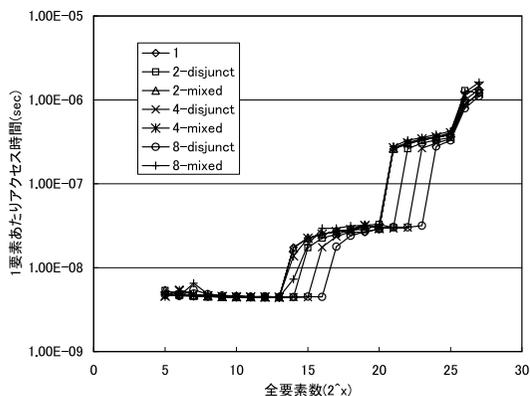


図1 ランダムアクセステスト (SPARC)
Fig. 1 Random access test on SPARC.

メモリモジュール群というものを考え、それらがバスで相互に接続されたものがメモリである、とするモデルである。つまり、CPU の近くには小さくて速いメモリモジュールがあり、次第に大きくて遅いモジュールがつながっていく、という構造である。これらのモデルはアルゴリズムの解析に成果をあげてきたが、並列アルゴリズムは対象としていなかった。

並列計算のモデルでは、 $\text{Log}P$ ⁴⁾ がよく知られている。これは、計算機間が4つのパラメータを用いて表現された均一な通信路で相互接続されている、とするモデルである。しかし、現在普及しつつある大規模な分散計算環境では、各計算機間の通信路が均一ではなく、遅延が大きいプロセッサとそうでないプロセッサ、すなわちプロセッサの「遠い近い」が計算時間に無視できない影響を与えてしまうため、このモデルも現実との乖離が大きくなりつつある。LogPの発展モデルとして、さらにパラメータを付け加え、詳細に通信路特性を記述しようとする試みもいくつか見られるが、これらはすべて通信路のモデル化にとどまっている。

また、Coarse Grained Multicomputers (CGM)⁵⁾ や、Bulk-synchronous Parallel (BSP)⁶⁾ といった並列計算モデルもある。これらは、ある程度の量のメモリをローカルに持つプロセッサが、何らかのトポロジを持つネットワークで相互に接続されている、と計算環境のモデル化を行い、そのうえで「ローカルな計算フェーズに必要なコスト」と「通信フェーズに必要なコスト」を別々に考え、加算することによって全体の計算量を把握しようというモデルである。

また、BSPを基本に、通信路の振舞いや計算機性能が場所によって異なるような計算機環境をモデル化する Heterogeneous BSP⁷⁾ も提案されている。計算フェーズ、通信フェーズが、並列計算全体で同期して

いる BSP モデルであるため、通信が最も遅いところに揃うことで通信路の構成に起因する振舞いをとらえようとしている。ただし、このモデルでの通信時間はプロセッサ入口のバンド幅と送受信するデータ量によってのみ決定されるというシンプルなものであり、通信路の構造、通信の局所性などは考慮していない。

これらのモデルのように、単体計算機内の計算と計算機間の通信とを別々に考慮する場合、どうしてもモデルの組合せによる複雑さが生じる。また、単体のCPU内のメモリ階層による影響を表現するモデルと、CPU間の通信路を精密に表現するような別々のモデルを組み合わせるという複雑な手順をとって、ある特定の計算環境を精度良く表現することが可能になったとしても、単体のCPUの構成が違う、あるいはネットワークポロジが違う環境を考えようとなると、モデルの再設定、および計算コストの評価のやり直しが必要となってしまう。計算機構成のトレンドは年々変化しており、トレンドが変化するたびにモデルを変更しなければならないのでは、アルゴリズムの設計に対する指針としては使いにくい。

さらに、非常に多数の計算機が結合した環境で、すべての通信路を独立に扱うことはもはや困難になる場合などもある。このように、並列計算環境を「複数の計算機がネットワークでつながったもの」としてとらえることは、モデルの複雑化を招き、またこれからおおいに発展すると思われる大規模分散環境に対応できないアプローチであると考えられる。

我々は「計算のコストとは、演算にかかる時間にあるのではなく、演算に必要なデータを取得し、演算結果を格納するという通信過程こそ存在する」という基本理念に基づいて、計算量を統一的にとらえ直そうと考えた。メモリ階層はデータへのアクセス遅延時間の違いが生み出すものであり、ネットワークポロジの違いを考慮すべきなもの、他の計算機にあるデータへのアクセス遅延時間の違いが無視できないからである。

この基本理念をもとに、「計算機の内側と外側の区別、単体計算機内の計算とネットワーク上の通信の区別をなくし、それらを統合した1つの仮想計算機として表現したモデル」によって、新しい計算量モデルを構築することを考えた。このモデルではすべてのメモリ間に距離 x が定義できるものとし、それぞれのメモリ間の通信コストが「距離による単純な関数 $f(x)$ 」で表されるものと仮定する。この $f(x)$ によって、任意の単体計算機内部のメモリ階層、任意のトポロジによるネットワークの振舞いを、統一的にかつ簡潔に表

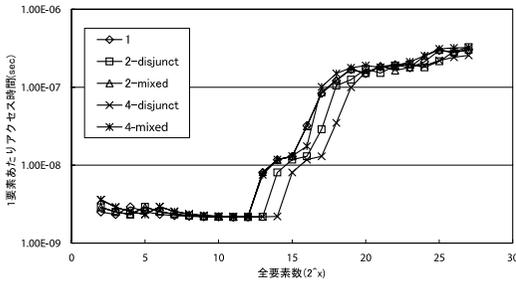


図2 ランダムアクセステスト (Opteron)
Fig.2 Random access test on Opteron.

現しうのではないか、というのが我々の主張である。これにより、アルゴリズム設計の際、計算コスト評価の労力が著しく緩和される。

このモデルは、特定の計算機環境における計算時間をきわめて精密に予測するというよりは、多くの並列計算機環境で共通した振舞いをよく予測できるような、一般性の高い指標を目指している。図2は、図1と同じ実験を Opteron 1.4 GHz, 4 CPU の計算機で実行した結果である。メモリアクセス時間の増大はキャッシュなどの構造によって異なるので、構造の異なる2つの環境では振舞いもかなり違っているが、全般的な増大の仕方に共通点もある。さらに、この実験はローカルメモリの大きさで制限されているが、ネットワークでつながった別の計算機のメモリを使ってランダムアクセス実験を行えば、このグラフの先もやはりアクセス時間は増大していくと予想される。その際も、計算機のネットワーク性能やネットワークポロジの違いによって、細かい振舞いは環境によって変化する。しかし、環境に依存しない増大の仕方を、ローカルメモリからネットワーク越しのメモリの範囲まですべて、 $f(x)$ という単純な関数で表現するよう割り切ってしまう、というのがこのモデルの基本的な方針である。

Rauber らの Locality measure⁸⁾ は、メモリアクセス履歴を用いた簡単な指標で、並列アルゴリズムの実行環境に依存しない実行性能を表現しようとするものであるが、我々のモデルも同じ方向を指向している。ただし、Locality measure がメモリアクセスの履歴から間接的な定性的要素を取り出しているのに対し、我々のモデルはより直接的に、計算機の構造を意識した指標を作ろうとしている。

このような並列計算コストのモデルを、我々は「アクセス計算量」と名付けた。本稿では、このアクセス計算量の概念と、仮想機械のアーキテクチャの設計、仮想機械語の設計を示す。また、いくつかの並列アルゴリズムへの計算量解析への適用を通してモデルの妥

当性と適用可能性について示す。

2. アーキテクチャと計算の概観

2.1 モデル概要

アクセス計算量とは、以下のようなモデルに従う計算量である。

- メモリはある密度で連続的に（現在のところ1次元で）分布したものであるとする。
- 計算する資源はメモリ上の至るところに存在している。メモリ上には「計算実体」と呼ばれる、計算を行う「もの」が存在し、メモリ上の好きな地点で計算を行うことができる。
- プログラムの実行時間はメモリアクセスの時間が支配的であり、演算そのものにかかる時間は無視できるとする。ただし、演算のためには必ずメモリ上からデータを持ってくる必要があるため、見掛け上無限に速い演算ができるわけではない。
- メモリへのアクセスには「格納された場所」から「計算する場所」までの「距離」に従ったコストがかかる。 x だけ離れた場所のメモリにアクセスするときは、指令からある時間 $f(x)$ だけたつた後にアクセスが行われる。 x の点でのアクセスにはある一定の時間 $l, l > 0$ がかかる。さらに、読んだデータを送り返す時間または ack を返す時間が $f(x)$ だけかかるため、メモリアクセス命令全体では $2f(x) + l$ だけの時間がかかる。
ack を返さない書き込み命令はない。
- $f(x)$ は単調増加関数であり、 $f(0) = 0$ が成り立つ。また、 $0 < x < y$ としたとき $f(y) < f(x) + f(y-x)$ が成り立つ、すなわち上に凸な関数である。これはつまり、どこかで通信を中継することで速いアクセスが可能になるということはない、ということを示している。
- メモリアクセスそのものの衝突は考えない。メモリアクセスの際に使われる通信路が混雑することで、アクセス競合の振舞いをとらえようとする。これは、データを近傍に並べるのではなく、適宜分散させて計算した方が速い場合があることを反映させようという狙いである。
- 計算実体は「プログラムカウンタ」と「実行場所」のみを持っており、PCを進めつつ、メモリ中で実行場所を移動させて計算を行う。計算実体自身にレジスタのようなメモリはない。
- 命令を取ってくる時のコストは考えない。
- 計算の実行場所が移るときもメモリアクセスと同じ時間コスト $f(x) + l$ がかかる。メモリアクセス

を計算実体が追い越すようなことはできない。

- 通信路にはある有限の容量があり、容量を超えた通信についてはペナルティが追加された時間コストがかかるものとする。これにより、計算主体自身は任意の並列性で存在できたとしても、計算全体は通信路の混雑により速度向上に制限がかかる。

今、1CPU しかない普通の計算機をアクセス計算量モデルで表現すると、1次元のメモリ上に計算実体が1つだけ存在する状態となる。計算実体の直近のメモリは、アクセスする際の距離が小さい、読み書きのレイテンシが最も短い場所であり、現実の計算機ではレジスタにあたる。その隣には、少し距離が大きくなった L1 キャッシュに相当するメモリがあり、以後だんだん遠くに L2 キャッシュ、メインメモリが置かれていると考えられる。現実の計算機ではレジスタ、L1、L2 キャッシュ、メインメモリとレイテンシが増大するに従って容量は大きくなる。このレイテンシと容量に合わせて $f(x)$ を定義すれば、現実の計算機のアクセスコストを模倣することができる。一般的には、レイテンシと容量の変化の度合いは容量の方が累乗的に増大する傾向にあり、HMM の研究によると $f(x) = \log(x)$ 程度が適切だといわれている。もちろん、近年はレイテンシの増大が激しくなっているので、 $f(x)$ をどのように設定したらよいかはまだ研究の余地がある。

現実の計算機では、メインメモリにアクセスするとアクセス対象の近傍のメモリがキャッシュメモリに移され、以後のアクセスが高速になるが、アクセス計算量モデルではキャッシュはアルゴリズム設計者が明示的に使用しなければならない。つまり、このモデルの仮想機械が持つブロック転送の機能を用いて遠いメモリから近くのメモリへとデータを移し、近くのメモリを使って計算を行う、という動作を明示的に行わせる必要がある。ただ、このモデルでは計算実体が自由に場所を移せるため、データを取ってくる代わりに、計算実体自身をデータの近くへと移すことでキャッシュと同じ効果が得られる。

次に、SMP マシンをこのモデルで表現すると、メモリ上にプロセッサ数だけの計算実体が置かれた状態になる。図 1 と図 2 の実験を例に、SMP マシンの動作がこのモデルでどのように表されるかを説明する。

図 1、図 2 の実験は、連続したメモリ領域の要素が 1 本のリストとしてつながっており、リストをたどって全要素を読むことで全メモリ領域をランダムに 1 回ずつ読む、ということを繰り返すものである。系列 1 は 1 スレッドで、2- から始まる系列は 2 スレッドで、のように複数スレッドで実験を行い、複数スレッドでも全

要素数 N は 1 スレッドと変わらず、各スレッドが (スレッド数 P として) N/P 個の要素を読む。-disjunct の系列は、 N の連続メモリ領域を N/P 個ずつの連続領域に分け、各スレッドはそれぞれの部分領域の中だけでランダムアクセスをした場合で、-mixed の系列は各スレッドが N の中からランダムに N/P 個のアクセスをした場合となる。

1 スレッドの場合、わざわざレイテンシの大きいメモリを使う必要はないので、計算実体は N の要素のどこかにいる。アクセス対象メモリまでの平均距離は $N/2$ である。-disjunct の場合、計算実体は N/P の要素のどこかにいるのが自然だろう。このとき、計算実体からアクセス対象メモリまでの平均距離は $N/2P$ である。それ以外は 1 スレッド実行のときと違いがない。そのため、図を見ると分かるとおり、-disjunct のアクセス時間は要素数 $1/P$ の 1 スレッド実行時のグラフときれいに同じ形をしている。一方、-mixed の場合、計算実体の位置は -disjunct と同じでも、アクセス対象までの平均距離は 1 スレッド実行と同じ $N/2$ となる。そのため、-mixed のグラフは 1 スレッド実行のグラフときれいに重なる。このように、このモデルは並列実行の様子を比較的直接的な形で把握できることを目指している。

ところで、メモリが 1 次元であることによる現実との乖離ももちろん存在する。SMP のモデル化で、3 個並んだ計算実体が互いに相手の近くのメモリを読むとき端の 2 つの間の方がレイテンシが大きくなることになるが、これは現実計算機では考えにくい。クラスタのような分散環境でも同様である。また、分散環境の場合、ネットワークポロジによっては、計算機間に複数の通信経路がある場合や、通信方向によってレイテンシに差がある場合などもありうる。このような現実を、1 次元のメモリモデルでは扱うことができない。しかし、たとえばメモリと通信路を 2 次元以上に拡張し、2 地点間の距離と通信経路を適切に定めることができれば、このような現実も扱うことが可能になる。ただし、多次元の拡張を行うと距離の定義が難しく、解析が困難になると予測される。多次元にする事で「迂回する」通信路を作ることができるのも、通信路の混雑の解析を困難にする。そもそも、一般的なネットワークは何次元あれば表現できるのかもよく分からない。以上のような理由で、現状のアクセス計算量モデルでは、メモリと通信路は 1 次元であると仮定する。多次元への拡張は今後の課題としたい。

2.2 複層の通信路モデル

上述のモデルでは、並列実行時の計算の振舞いは通

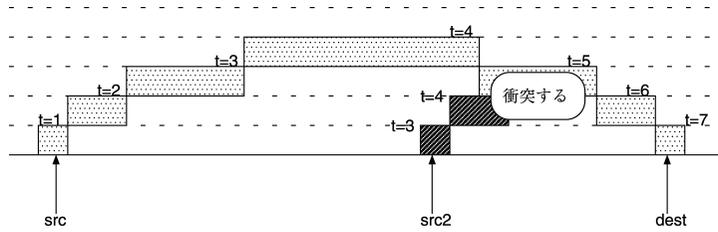


図 3 通信路と衝突
Fig. 3 Communication channel and collision.

信路によって大きく規定される．計算と同様，通信路の混雑の局所性も計算量を見積もる際に大きな影響を与える，という立場から，我々のモデルではすべての通信をパケットととらえ，パケットの動きを追跡することで通信路の振舞いを把握することにした．以下に，この「複層の通信路モデル」の概要を示す．

- メモリが1次元であるモデルの場合，通信路は正方向行き/負方向行きの2つが独立して存在する．正方向と負方向に送られるパケット相互は衝突しない．
- 通信路は層状に積み重なっている．通信の衝突は各層内でのみ発生する．各層の通信速度は，底の方ほど遅い．通信は一番底の通信路からスタートし，一定時間ごとに隣の層へと移動していく．このとき，世界全体で共有されたパルスが存在すると仮定している．一定時間ごとに世界全体でパルスが発生し，そのパルスに従って通信路がいっせいに状態を変える．
- 各通信路の速度は $f(x)$ によって規定される．たとえば $f(x) = O(\log(x))$ の場合，高さ h の通信層の速度を $v(h)$ ， k を1より大きい適切な定数として， $v(h) = kv(h-1)$ とすればよい(図3参照)．
- 通信は，現在使っている通信層の速度 $v(h)$ ，データ長 n に対して， $(n+1)v(h)$ に相当する広さの領域を占有する．1はパケットが持つ定数のコストを示しており， n が大きいパケットを用いるほど通信路の利用効率が上がることを表現している．
- パケットは任意のデータ長を一度に送ることができる．複数ワードを送信する際は，送信するデータが置かれた場所を次々にパケットが訪問し，次第に長いパケットに成長していく．データを受信する側でも同様に次第にパケットが短くなりながらメモリ上に書き込みが行われる(図4参照)．
- 次の単位時間に占有するべき領域が空いていないとき，パケットは空きを待つ．空きが足りない場合はできる限り進もうとする．2つのパケットが衝突しているときは，底の通信路へ降りていく(パケットが減る方向の)パケットを優先させる(図5参照)．この通信路モデルは，いわばキャッシュメモリからメ

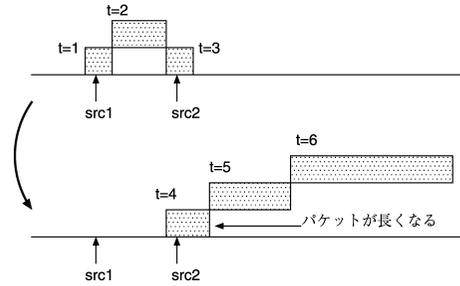


図 4 2ワードを1パケットにまとめて送信する場合
Fig. 4 Sending two words together in one packet.

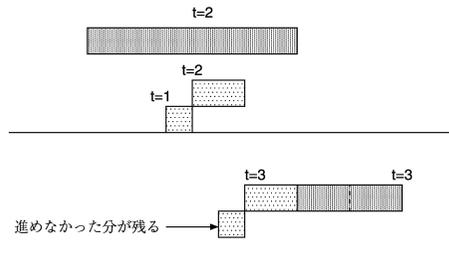


図 5 通信路がふさがっている場合
Fig. 5 Blocked communication.

インメモリへと連なるメモリバス，さらに Local Area Network, Wide Area Network へとという階層的な物理的通信路のモデル化を狙ったものである．一番底の層が CPU 内部にあるレジスタやキャッシュとのバス，次の層がメインメモリのためのバスで，その上に順に LAN, WAN と連なっている，という意識である．たとえば，LAN 経由で隣の計算機と通信をする際，データはメモリからメモリバス，LAN を通り，相手の計算機のメモリバスへと移行する．LAN の層でパケットの衝突が起きたとき，複数の計算機の持つ広い範囲のメモリ領域の上の通信層が影響を受けていることになる．モデル上で，上に位置する層の方が速度が速いが1つのパケットがより広範囲に影響を及ぼすようになる，という設定はこのような状況を表現しようとしたものである．ただし，LAN の通信層が輻輳を起こしているときでも，単体計算機内の(キャッシュ)メ

モリアクセスは影響を受けない．この現実が，モデル上では，上の層の通信衝突が下の層には影響を与えないという設定で表現されている．このような層構造は，CPU 内部のメモリバスでも，大規模なネットワーク環境でも，共通して見られる構造であると考えている．

2.3 通信負荷累積モデル

我々は，このような計算モデルに基づく仮想機械を設計（後述，3 章）し，このモデルでの計算の様子を把握できるシミュレータを構築した．また，この仮想機械上でのプログラミング手法を提案し，いくつかのアルゴリズムをシミュレータ上で実行することで計算モデルの妥当性や有効性を検討してきた．その結果，複層の通信路モデルは現実のシステムを確かによく表現するが，解析的に計算量を求めるにはやや複雑すぎることも分かった．そこで，通信路の混雑をより単純にモデル化した「通信負荷累積モデル」を作成し，解析的手法による計算量の見積りを可能にしようと試みた．

通信負荷累積モデルでも，通信を担うパケットの各時点での速度，およびそのパケットが影響を与える範囲は複層の通信路モデルと同じである．結果，距離 x だけ離れた 2 地点間の通信は少なくとも $f(x)$ だけの時間がかかる．2 つのモデルは，通信の集中による通信路の混雑の表現についてのみ異なっている．

パケットは現在占めている通信路の領域にある負荷を与える．複数のパケットが同じ領域に存在するとき，その領域にはそれぞれの負荷を加算した負荷が与えられる．通信路には決まった容量が定められており，ある領域の負荷が容量を超えたとき，その時刻の通信路全体にわたって混雑による通信遅延が起きるものとする．

パケットは，速度に応じて負荷を与える範囲が変わる．単位量の通信パケットは全体としてはある一定の負荷を持ち，ある時刻での地点 x での負荷を $l(x)$ とすると，

$$\int l(x)dx = constant \quad (1)$$

が成り立つ．つまり，速度が上がると広い範囲に負荷を与えるが，その場合ある地点に与える負荷の値は小さくなる．このような $l(x)$ の形はさまざまなものが考えられるが，以降の議論では，ある時刻のパケットの速度を v としたとき， v の範囲に均一に $1/v$ の負荷が与えられる，という単純なモデルを採用する．

図 6 に，1 つのパケットのみが存在する時各時刻において通信路に与えられる負荷を示す．パケットの速度，存在するメモリ上の位置は図 3 と同様に変化しているが，通信路に層は存在せず，1 つの通信路に与え

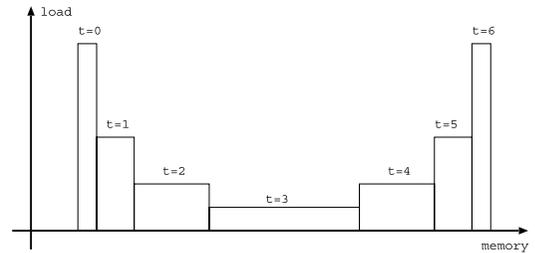


図 6 通信負荷累積モデルでの 1 つのパケットの負荷

Fig. 6 Transition of communication load caused by a packet in load sum-up model.

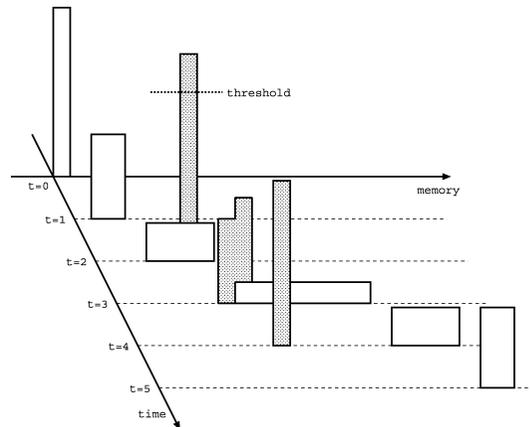


図 7 通信負荷累積モデルでの複数のパケットの負荷

Fig. 7 Transition of communication load caused by multiple packets in load sum-up model.

る負荷 ($load$) が速度に応じて変化している．

図 7 は，2 つのパケットが衝突する様子を示している． $t = 0$ で 1 つ目， $t = 2$ で 2 つ目のパケットが発生し，ともに右に向かって進んでいる． $t = 2, 3$ では 2 つのパケットが通信路の同じ領域に存在するため，2 つの負荷を加算したものがその時刻の通信路負荷となる． $t = 2$ で，2 つの負荷の合計がある定められた通信路容量 $threshold$ を超えるため，この時刻の通信路中すべての通信はペナルティを受けて遅くなる．ペナルティは，負荷 l が容量 T を超えたとき，通信は通常の l/T 倍の時間がかかるというモデルを想定している．

このような通信路は，「複層になった通信路モデル」と比較して，現実の通信路の特性をいくつか無視している．たとえば，複数のパケットの負荷を単純に加算してある地点の通信路負荷を求めているため，LAN 内の通信が輻輳を起こすとローカルメモリの読み書きが遅くなる，というモデルであるといえる．しかし，このような単純化によって，複数の通信が同時に行われるときの通信路の振舞いを簡単に表現でき，解析的

手法で計算量を求めることが可能になると期待される。

3. 仮想機械

これまで述べてきた計算モデルは，次のような命令セットを持つ仮想機械であると定義できる．以下の説明で， $\langle input \rangle$ ， $\langle output \rangle$ は即値およびダイレクト，インダイレクトのアドレッシングモードをとるメモリ内容である．なお，アドレスは現在の計算実体が存在している場所からの相対位置を指定する．各々の命令は現在の実行場所で実行され，実行後 PC を 1 つ進める．

3.1 演算・メモリ

$\langle operation-name \rangle \langle input_1 \rangle \langle input_2 \rangle \langle output \rangle$
 $input_1 \cdot input_2 \rightarrow output$ という演算を行う．

$\langle operation-name \rangle = \text{add, sub, などの基本演算}$
 $\langle copy \langle src \rangle \langle dest \rangle \rangle \quad src \rightarrow dest$ というメモリ内容コピーを行う．メモリ上に定数を書き込みたいときもこの命令を使う (src が即値になる)．ブロック転送も可能．

3.2 実行場所の制御

$\langle next_place \langle next-place \rangle \rangle$ 次の実行場所を現在の場所からの相対位置で指示．

3.3 制御

$\langle jump \langle program-point \rangle \rangle$ 無条件ジャンプ．

$\langle program-point \rangle$ の命令を実行．

$\langle branch \langle input \rangle \langle program-point \rangle \rangle \quad \langle input \rangle$ の場所のメモリの値によって分岐．

$\langle fork \langle program-point \rangle \rangle$ 別の計算実体を作成し， $\langle program-point \rangle$ の命令を，現在の実行場所で実行させる．自らは次の命令を現在の実行場所で実行．

$\langle compare_and_swap \langle src \rangle \langle org \rangle \langle update \rangle \langle pp \rangle \rangle$

$\langle src \rangle$ の内容を compare and swap ．

$\langle vanish \rangle$ 計算実体が消える．

これらの命令セットを用いてアルゴリズムを記述することで，すべての計算に必要なデータの配置とその間の通信を明記することになり，アクセス計算量モデルに基づいた計算量を求めることができる．

我々は，この仮想機械プログラムを解釈実行し，その計算の振舞いを把握できるシミュレータを実装した．また，この仮想機械上でやや高水準なプログラミングが行える言語 CEMA⁹⁾ を設計，実装した．複雑な並列アルゴリズムの計算量を考察する場合においては，このようなシミュレーションを用いることができる．

4. 並列アルゴリズムの解析

ここでは，よく知られた並列アルゴリズムである，

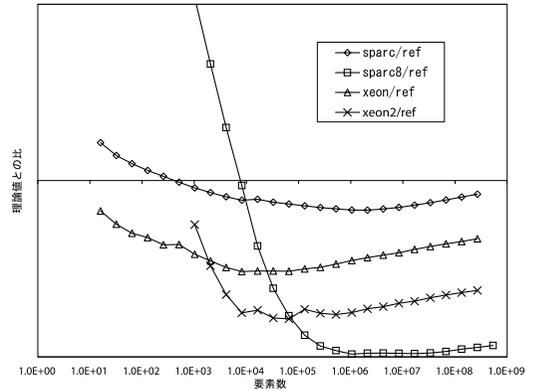


図 8 bitonic sort における実計算機と RAM モデルの計算量比較
 Fig. 8 Ratio of the measured and predicted performances of the RAM model in bitonic sort.

bitonic sort, merge sort, FFT を例にアクセス計算量モデルでの計算量解析を行う．また，bitonic sort を実計算機上で実行し，(P)RAM モデルとアクセス計算量モデルの解析結果と比較する．

4.1 bitonic sort

bitonic sort アルゴリズムは，sorting network の一種である．規則的な比較演算を全データに対して繰り返すアルゴリズムであり，比較を行う要素の対がデータにかかわらず固定されているため，並列計算の適用も比較的容易である．

n 要素の bitonic sort は $(n/2) \times (\log n (\log n + 1) / 2)$ 回の比較演算を含むため，RAM モデルにおいて計算量は $O(n(\log n)^2)$ と表される．また，並列計算の際には $n/2$ 回の比較操作を並列に実行できるため，並列度 P のとき $O(n(\log n)^2 / P)$ ，最高で $O((\log n)^2)$ の計算量となる．

実計算機上での実験結果を図 8 に示す．計算環境は

- Ultra SPARC III Cu 1.2 GHz, 8 CPU
- Pentium4 Xeon 2.4 GHz, 2 CPU

である．いずれも pthread ライブラリを用い，共有メモリを使って並列化を行っている．4つのグラフ系列は，2つのアーキテクチャそれぞれで single thread 実行時と 8(2)thread 実行時の 2通りの計算時間を示している．横軸はデータ数 n ，縦軸は RAM モデルでの理論計算量と実際の計算時間の比（分母が理論計算量）である．ただし，縦軸は適当に倍率を決めたものであり，SPARC どうし，Xeon どうしの系列の値の比は計算時間の比を示しているが，SPARC と Xeon の系列の縦軸の値の比，ならびにその絶対値は何の意味も持たない．

図 8 に示されるとおり，データの量が大きくなるに

つれて理論値との比が単調に増加している。つまり、RAM による単純なモデル化では把握できないメモリ階層などの要素が現実の計算時間を押し上げており、RAM モデルの計算量と現実との乖離が起きているのである。

次に、「アクセス計算量モデル」による解析を行う。

アクセス計算量モデルでは、データの配置と計算実体の位置が計算量を決定する。実験に用いたプログラムでは、データは最初に連続領域に配置され、そのまま配置を変更することなく同じ連続領域でソートされる。このとき、 n 要素 bitonic sort に対して一度呼ばれる n 要素 bitonic merge は、 $i = 0$ から $n/2 - 1$ まで i を順に変えつつ、 i 番目と $i + n/2$ 番目要素の比較を行う。

これをアクセス計算量モデルで解析するとき、データ配置は元プログラムの配置に従えばよいが、計算実体の場所は通常のプログラムには概念が存在しないため、解析者が適切に定める必要がある。今回の場合、要素に順番にアクセスするため、キャッシュメモリにデータが乗っていることが多いと考えられる。解析をより正確にするためには、次に比較に使う 2 つの要素付近のメモリを計算実体そばにブロック転送し、キャッシュの効果を再現する方がよいが、今回は比較に使う片方の要素の上に計算実体が移動することで、簡易的にキャッシュの効果を再現することにした。この場合、もう片方の要素はやや遠い位置に存在することになるが、アクセスの際の距離がつねに一定になるため、それほどひどいメモリアクセスコストを生じず「順番に行うアクセス」の振舞いを正しく把握できるのではないかと考えた。また、メモリの距離 x に対するアクセスコストは HMM の結果より $\log x$ と定めた。

このとき、 n 要素 bitonic sort に対して一度呼ばれる n 要素 bitonic merge 内での比較演算は距離 $n/2$ のメモリアクセスを行っているので、 $O(\log n)$ の遅延が起きることになる。これを再帰的に適用していくと、bitonic sort 全体の計算量は $O(n(\log n)^3)$ となる。

また、 P プロセッサでの並列実行時には $n/2P$ だけ離れた箇所でも同じ距離の通信が並行して行われる。最も通信路が混雑するのは $P = n/2$ のときであり、連続した P 個のメモリ上の点から、 P だけ離れた点への通信がいっせいに行われることになる (図 9 参照)。

このような場合について、通信路の「通信負荷累積モデル」を用いて解析する。図 9 の状況において、任意の地点 x の通信路の負荷は、同一の負荷の形をしたパケットが距離 1 ずつずれながら加算されたものであり、1 つのパケットのある時刻の負荷を $l(x)$ として、

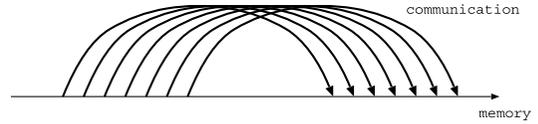


図 9 連続した点からの同一方向へのいっせい通信

Fig. 9 Concurrent communication from and to continuous locations.

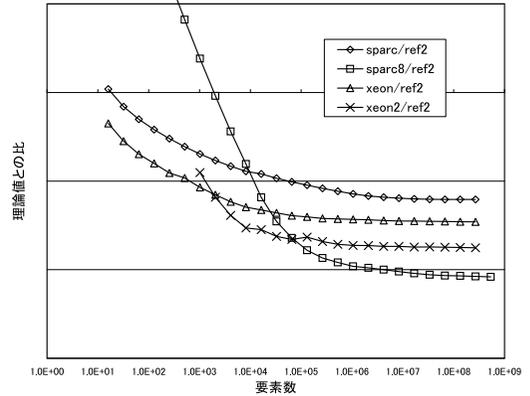


図 10 bitonic sort における実計算機とアクセス計算量モデルの計算量比較

Fig. 10 Ratio of the measured and predicted performances of the access complexity model in bitonic sort.

$$\sum^d l(x-d)$$

で表される。これは、1 つのパケットの持つ負荷の総量が一定であるという制限 (式 (1)) より、一定値であるといえる。つまり、複数のパケットが生まれていても、あらゆる地点の負荷は単一のパケットのみが存在したときの負荷を超えないことになり、このような通信パターンは通信路を混雑させない、と解析されることになる。

よって、並列計算時でも通信路の混雑は起きず、計算量は $O(n(\log n)^3/P)$ 、最高で $O((\log n)^3)$ となる。

先の実験の実行時間を、アクセス計算量モデルの理論計算量と比較したものを図 10 に示す。縦軸・横軸などは図 8 と同じ意味である。図 8 の RAM モデルでは表現できていなかった、要素数の大きい領域において、計算時間とアクセス計算量モデルの理論値との比は一定に収束していている。つまり、アクセス計算量モデルは bitonic sort の計算量をうまく見積もれたといえる。

4.2 merge sort

merge sort は並列化手法を決定的に決めることができるアルゴリズムであり、解析は容易である。RAM モデルでの計算量は $O(n \log n)$ 、並列化すると最も多くのプロセッサを使った場合で $O(n)$ となる。

次に、アクセス計算量モデルでの計算量を考える。まず、サイズ n の連続した入力データに対し、隣に同じ大きさのバッファを用意し、その2つのメモリ領域を交互に使ってソートを行うものとする。通信は2つのソートされた部分領域の要素（2領域合わせた要素数を k とする）を順々に比較するとき、および n だけ離れたバッファに比較後の結果をコピーするときの両方で起きる。部分領域間の通信より結果コピーときの距離 n の通信コストが支配的であるので、結果コピー時の通信が最も遅くなるような場合を考えると、これは k 要素すべてがソートされていて、すべての要素が図9のような距離 n の通信を行うときだと考えられる。もし、図9の通信のどこかの要素の順序が入れ替わっていると、その部分の通信は「距離 n を2回」から「 $n-1$ と $n+1$ の距離の通信を1回ずつ」へと変化する。このとき、 $2f(n)$ と $f(n-1) + f(n+1)$ の通信遅延を比較すると、 $f(x)$ が上に凸な関数であるので $2f(n) > f(n-1) + f(n+1)$ であり、順序を入れ替えた方が通信時間が短くなる。よって、最も遅くなるようなアクセスパターンは図9のようにすべての要素が距離 n の通信をするときで、そのときの通信コストは $k \log n$ である。

このとき、ソートされた領域の比較に必要な通信回数は $2 \sum_{i=1}^{k/2} \log i = k \log k - 2(k-1)$ となり、 k 要素の merge sort の時間 $T(k)$ についての漸化式（全体の要素数は n ）を書くと

$$T(k) = 2T(k/2) + k \log n + k \log k - 2(k-1)$$

となる。これを解くと計算量は $O(n(\log n)^2)$ となる。

異なったデータ配置でのアルゴリズム、たとえば入力データと作業用バッファを1要素ずつ交互に配置する方式も考えられるが、この配置の場合でもアクセス計算量モデルでの計算量は変わらなかった。

また、並列化について考えると、bitonic sort と同様の考え方で通信路は混雑しないと結論づけられる。そのため、十分多いプロセッサが存在する場合は $O(n \log n)$ の計算量となる。

実計算機上での実験結果を各モデルと比較したものを図11、図12に示す。merge sort の場合、SPARC 1 CPU, Xeon 2 CPU の構成では RAM モデルとの理論値比は要素数が大きい領域でじわじわと上昇しているが、Xeon 1 CPU の構成では RAM モデルとよく一致しているように見える。一方、アクセス計算量モデルでの予測との比を見ると、SPARC 1 CPU や Xeon 2 CPU では要素数の大きい領域で予測と一致する傾向にあるが、一般的に理論値比が低下する方向にある。つまり、本実験のプログラムは、RAM モデルの

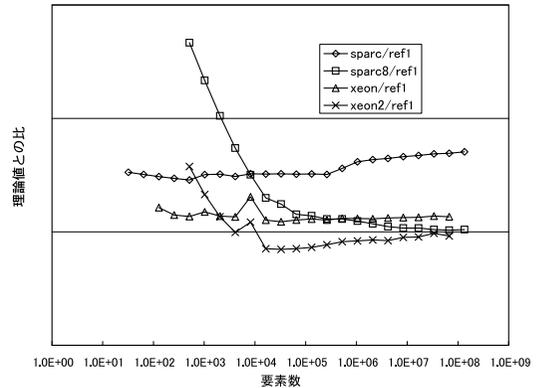


図 11 merge sort における実計算機と RAM モデルの計算量比較

Fig. 11 Ratio of the measured and predicted performances of the RAM model in merge sort.

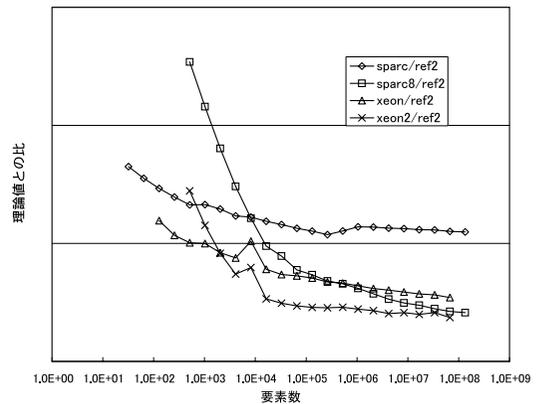


図 12 merge sort における実計算機とアクセス計算量モデルの計算量比較

Fig. 12 Ratio of the measured and predicted performances of the access complexity model in merge sort.

$O(n \log n)$ と、アクセス計算量モデルの $O(n(\log n)^2)$ の間のあたりの計算量を持ち、アクセス計算量モデルによる計算量と完全には一致していない。

これは、2つの領域をそれぞれ連続アクセスし、その後でその領域2つを連続したものとしてアクセスする、という merge sort のメモリアクセスパターンが、実計算機のキャッシュの仕組みにうまく当たるのに対し、アクセス計算量モデルで解析を行う際、明示的にキャッシュを使うようなアルゴリズムを採用しなかったためであると考えられる。

4.3 FFT

FFT については、メモリアクセスのローカリティを高めるためのさまざまな手法が知られている。また、演算数を減らすための手法も多い。しかし、工夫を行っても演算数のオーダそのものは変化しないこと、また

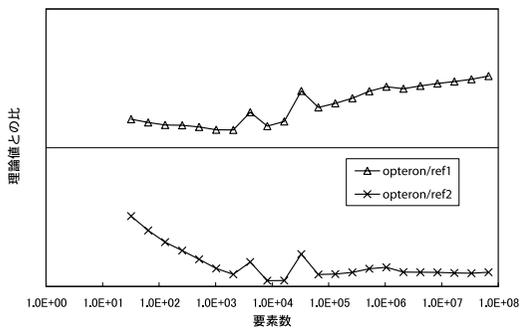


図 13 実計算機と各モデルの計算量比較

Fig. 13 Ratio of the measured and predicted performances of two models in pseudo FFT.

演算数を減らす手法は並列化したときの解析が面倒になることを考慮し、今回は一番単純な形の FFT、つまり n 入力 FFT を「バタフライ回路 → 2 個の $n/2$ 入力 FFT → シャッフル回路」と計算する方式について解析を行う。RAM モデルでの計算量は $O(n \log n)$ 、十分にプロセッサが存在するときは $O(\log n)$ である。

アクセス計算量モデルでは、 n 個の入力データは連続したメモリ上に置かれ、隣接した同じ大きさの作業領域を使って書き換えられるものとする。全体が n 要素の FFT で、 k 要素の部分 FFT を考えるとき、バタフライ回路の通信は $\log(k/2)$ の遅延が k 回、作業領域への書き込みが遅延 $\log(n)$ を k 回、シャッフル回路の通信コストは作業領域から入力データ領域への移動なので $k \log(n)$ である。よって、計算時間の漸化式は

$$T(k) = 2T(k/2) + k \log(k/2) + 2k \log(n)$$

となり、これを解いて全体の計算量は $O(n(\log n)^2)$ となる。また、並列計算時は bitonic sort と同様の議論で通信路が混雑しないことが示され、十分にプロセッサが存在する場合で $O((\log n)^2)$ となる。

簡単な実験を行ったものとモデルとの比較を図 13 に示す。これは厳密な FFT ではなく、バタフライで三角関数を掛けるべきところを定数を掛けている。これは、通常のライブラリの三角関数は重いので、三角関数計算が律速になってしまうからである。実用的な FFT の場合は三角関数も表を引くなどして高速化する必要があるのだが、どの程度のメモリを使って高速化するか、もアクセス計算量モデルでの解析では重要な要素となってくるため、今回の解析ではその部分を無視することにした。参考程度の実験であるので、並列化も行っておらず、実験環境も少し異なって Opteron 1.4GHz である。opteron/ref1 が RAM モデル理論値との比、opteron/ref2 がアクセス計算量との比であ

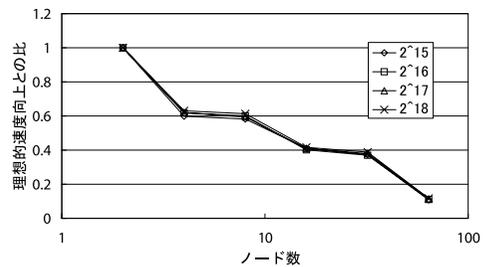


図 14 同時通信数と通信速度

Fig. 14 Ratio of measured and predicted performances in parallel communication.

る。図から読めるとおり、RAM モデルでは理論値との比が増大してしまっているが、アクセス計算量モデルでは正しく実行時間を把握できている。

merge sort と異なり、今回の単純な FFT の場合はシャッフル操作の部分がほとんどランダムアクセスで、実計算機ではキャッシュをうまく使えない。そのため、今回のキャッシュを意識しない解析方法でもうまく計算量が見積もれたと考えられる。

4.4 通信路の輻輳実験

4.1 節の bitonic sort の共有メモリによる並列計算では、効率はそれほど落ちることなく並列化が行われている。これは連続体モデルの計算量とも合致しているが、通信路のバンド幅がもっと狭い、ネットワーク接続された分散計算機のような環境では同様の結果となるか、についても確認する必要がある。

そこで、bitonic sort で起きる通信パターン、同じ距離の通信が複数密集して起きる図 9 のような状態を、PC クラスタ上で再現した実験を行った。実験環境は Xeon 2.4 GHz dual, 64 node のクラスタであり、ネットワークは Gbit-ether, 通信ライブラリに phoenix¹⁰⁾ を用いた。全体のデータサイズを一定とし、通信を行うノードの数を変化させて通信時間を測定した。モデルでは台数に比例した速度向上が得られると期待される。結果を図 14 に示す。グラフの系列は全体のデータサイズの違いを示しているが、結果は 4 系列ともほとんど変わらない。横軸はノード数、縦軸は期待される理想速度向上と実際の通信時間との比を示している。図 14 から分かる通り、ノード数が増えるに従って速度向上率が落ちる結果となっている。

つまり、ネットワークの帯域幅が十分に大きい環境では現在の通信負荷累積モデルで精度良く振舞いを予測できるが、帯域幅が小さいところは現在のモデルではうまく説明できていないといえる。通信路をすべて統一的に表現することを目指し、通信路のモデル化についてはさらに検討を続ける必要がある。

5. ま と め

我々は、計算のコストの本質は演算ではなく通信にこそ存在するという基本概念から、アルゴリズムの計算量を、均一に広がったメモリ上に存在するデータ間の通信の遅延の総和であるとする、アクセス計算量モデルを提案する。このモデルを用いることで、単一計算機内のメモリ階層から計算機間のネットワーク遅延の差異までを統一的に、かつ簡潔に記述することができる。このモデルは十分簡潔なものであり、並列アルゴリズムの計算量を解析的に求めることができることを、merge sort, FFT のアルゴリズムで検証した。また、bitonic sort アルゴリズムの実計算機上での実行と比較することで、従来の PRAM モデルが表現しきれなかった振舞いを、このモデルがよく表現できていることが示された。FFT の簡単なプログラムでもこのモデルの予測が実験とよく一致した。

ただし、merge sort の解析結果が実計算機と必ずしも一致しなかったことから分かれるとおり、キャッシュを明示的に操作するアルゴリズムで解析しなければならぬ場合も多いと考えられる。さらに、通信路のモデル化には、帯域幅が十分でない環境において、通信の輻輳を必ずしも十分に表現しきれていないという側面もある。これらを今後の検討課題としたい。また、より多くのアルゴリズムやより大規模なアルゴリズムに適用することで、このモデルによる計算量解析の妥当性と適用可能性を検証していくことも、今後取り組むべき課題である。

謝辞 本研究の一部は文部科学省科学研究費特定領域研究「IT の深化の基盤を拓く情報学研究」の一環として行った。

参 考 文 献

- 1) Aho, A.V., Hopcroft, J.E. and Ullman, J.E.: *The Design and Analysis of Computer Algorithms*, Addison-Wesley (1974).
- 2) Aggarwal, A., Alpern, B., Chandra, A. and Snir, M.: A Model for Hierarchical Memory, *Proc. 19th Annual ACM Symposium on Theory of Computing*, pp.305–314 (1987).
- 3) Alpern, B., Carter, L., Feig, E. and Selker, T.: The Uniform Memory Hierarchy Model of Computation, *Algorithmica*, Vol.12, No.2/3, pp.72–109 (1994).
- 4) Culler, D.E., Karp, R.M., Patterson, D.A., Sahay, A., Schauer, K.E., Santos, E., Subramonian, R. and von Eicken, T.: LogP: Towards a Realistic Model of Parallel Compu-

tation, *PPoPP*, pp.1–12 (1993).

- 5) Dehne, F.K.H.A., Fabri, A. and Rau-Chaplin, A.: Scalable parallel computational geometry for coarse grained multicomputers, *Intl. Journal of Computational Geometry and Applications*, Vol.6, No.3, pp.379–400 (1996).
- 6) Valiant, L.G.: A bridging model for parallel computation, *Comm. ACM archive*, Vol.33, No.8, pp.103–111 (1990).
- 7) Williams, T.L. and Parsons, R.J.: The Heterogeneous Bulk Synchronous Parallel Model, *LNCS*, Vol.1800, pp.102–108 (2000).
- 8) Rauber, T. and Rünger, G.: Program-Based Locality Measures for Scientific Computing, *IPDPS*, p.164 (2003).
- 9) 渡邊誠也, 横山大作, 近山 隆, 小宮常康, 湯浅太一: メモリ上の配置を意識する並列処理向き高水準機械語の設計と実装, ソフトウェア学会第 20 回大会 (2003).
- 10) Taura, K., Endo, T., Kaneda, K. and Yonezawa, A.: Phoenix: a Parallel Programming Model for Accommodating Dynamically Joining/Leaving Resources, *PPoPP* (2003).

(平成 17 年 1 月 24 日受付)

(平成 17 年 5 月 16 日採録)



横山 大作 (正会員)

昭和 49 年生。平成 12 年東京大学大学院工学研究科情報工学専攻修士課程修了。平成 14 年同博士課程中退。同年より東京大学新領域創成科学研究科助手。並列計算量理論、並列プログラミングライブラリ、およびゲームプログラミングに関する研究に従事。ソフトウェア学会、IEEE-CS 各会員。



近山 隆 (正会員)

昭和 28 年生。昭和 52 年東京大学工学部計数工学科卒業。昭和 57 年同大学院工学系研究科情報工学専門課程博士課程修了。工学博士。同年(株)富士通入社その後(財)新世代コンピュータ技術開発機構に出向し、第五世代コンピュータプロジェクトに参加。平成 7 年東京大学助教授。平成 8 年同教授。現在同新領域創成科学研究科基盤情報学専攻教授。プログラミング言語と処理系・開発環境、並列分散処理、機械学習と応用等の研究に従事。