

Map Sort: マルチコアプロセッサに向けた スケーラブルなソートアルゴリズム

枝廣 正人[†] 山下 康子[‡]

[†] NEC システムデバイス研究所 〒229-1198 神奈川県相模原市下九沢 1120

[‡] NEC ソリューション開発研究本部 〒305-8501 茨城県つくば市御幸が丘 34

E-mail: [†]eda@bp.jp.nec.com

あらまし マルチコア向けの並列ソートアルゴリズムMap Sortを提案する。今後単体CPUの性能向上が鈍化し、プロセッサがマルチコアによって性能向上する時代では、並列対応されていないソフトウェアは計算機が進歩しても性能は向上しない。従って単体CPUでは従来と同等処理時間で、かつ並列CPUではスケーラブルに性能向上するようなアルゴリズムが必須となるが、我々はそれをスケーラブルアルゴリズムとよんでいる。本論文ではソート問題を取り上げ、新しいスケーラブルアルゴリズムMap Sortを提案する。Map Sortの時間に関する計算複雑度は M 個のデータ、 P 台のCPUで $O((N/P) \log N)$ であり、単体CPU上での下界値 $O(N \log N)$ の $(1/P)$ である。また計算機実験の結果、単体CPU上のクイックソートと比較し、単体CPUでは同等性能、4CPUでは3倍の性能向上であることが示された。

キーワード マルチコア、並列、ソート、アルゴリズム

Map Sort: A Scalable Sorting Algorithm for Multi-Core Processors

Masato EDAHIRO[†] and Yoshiko YAMASHITA[‡]

[†] System Devices Res. Labs, NEC Corporation, 1120 Shimokuzawa, Sagamihara, 229-1198 Japan

[‡] Solutions Dev. Labs, NEC Corporation, 34 Miyukigaoka, Tsukuba, 305-8501 Japan

E-mail: [†]eda@bp.jp.nec.com

Abstract A sorting algorithm, called *Map Sort*, is proposed for multi-core processors. The proposed algorithm uses a map from subsets of input data to intervals in data range. Using the map, our algorithm is executed in parallel with multiple CPUs and has $O((N/P) \log N)$ time complexity for N data and P CPUs. Experimental results show that, in comparison with quick sort on a single CPU, processing time of Map Sort is comparable on a single CPU and three times faster on four CPUs.

Keyword Multi-Core, Parallelism, Sorting, Algorithm

1. はじめに

今後の半導体デバイスプロセスにおいては、微細化は進むものの動作周波数を上げつつ電力を下げる事がますます難しくなってきており、装置の電力制約を考慮すると単体CPUの動作周波数向上はあまり期待できない。そのため現在パソコンや高性能組込システムにおいて始まっているように[1, 2]、単体CPUの動作周波数よりはマルチコア化など並列により性能を高めていく傾向が強まると考えられている。従って並列化されていないソフトウェアは装置の進化に合わせた性能向上は従来ほど期待できず、並列化されたソフトウェアのみが性能向上していくことになる。

このような時代において、我々は以下の三点の性質

を持つアルゴリズムが重要であると考え、スケーラブルなアルゴリズムとよんでいる。

- 1) P 台のCPUを用いたときに時間に関する計算複雑度が $(1/P)$
- 2) 単体CPUでの実行時間は、単体CPU向けに最適化された従来のアルゴリズムと同等
- 3) 個数CPUでの並列性能向上は高いことが望まれる

メディア処理や科学技術計算においては本質的に並列性を含むものが多く、並列アルゴリズムの研究も進んでいるが、いわゆる基本アルゴリズムや、経路探索、辞書サーチなどのグラフ・ネットワークアルゴリズム[3]は並列性が抽出しにくいものが多い。

ソート処理はそのような基本アルゴリズムの中で代表的な問題である。これまで多くの並列ソートアルゴリズムが提案されてきているが[4]、上記の意味で優れたスケーブルアルゴリズムは提案されていない。クイックソートは単体プロセッサ向けでは最も効率の良いアルゴリズムと言われているが、並列化すると再帰的な分割の最初の数回は十分な並列性能が出せない。マージソートは理論上最も良いアルゴリズムであるが、近年の計算機システムで用いるとメモリ上のコピー回数が多いためクイックソートよりも実行時間性能が良くない。またソート済み配列のマージ処理を複数CPUで行うアルゴリズムが複雑であることも実行時間を遅くする原因となっている。並列計算機上での計算複雑度が優れたアルゴリズムについてもバイトニックソート、Odd-Evenマージソートなど多くの提案があるが、 $O(N)$ 台のCPUを仮定しており、現実的なCPU数ではCPU数に見合った並列実行性能は見込めない。

本論文では、スケーラブルなソートアルゴリズム *Map Sort*を提案する。*Map Sort*では入力データ配列を複数部分配列に分割し、データが分布する範囲を複数区間に分割する。そして複数部分配列から複数区間へのマップ(Map)を考える。マップの計算、複数部分配列から複数区間へのデータコピー、複数区間ごとのソート処理をそれぞれ並列に行うことができ、結果として $O((N/P) \log N)$ の時間複雑度を達成する。

計算機実験の結果、*Map Sort*は単体CPUではクイックソートと同等性能を達成し、4 CPUでは単体CPU上のクイックソートと比較して3倍の性能向上を達成できることがわかった。

2. 用語

N 個のデータからなる全順序集合 S があたえられたとき、ソートとは昇順（または降順）に並べ替える処理である。本論文では、ソート処理を P 個のコアを持つSMP (Symmetric Parallel Processing)、共有メモリ型のマルチ（コア）プロセッサで解くことを考える。また P は N よりもはるかに小さいと仮定する。

提案するアルゴリズムでは、与えられたデータ集合 S を M 個の部分集合 $\{S_1, \dots, S_M\}$ に分割し、すべてのデータを含む範囲を L 個の区間に分割する。 M および L は m 、 l を小さい定数とするとき $M=mP$ 、 $L=lP$ のように選ぶことが多い。実際のプログラムでは入力データ配列を M 個の部分配列に分割すると都合がよい。また、出力データ配列は、ソート後の状態を考えると、 L 個の区間にに対応して L 個の部分配列に分割されることになる。提案アルゴリズムでは出力データ配列の L 個の

部分配列のそれぞれがさらに M 個の部分配列に分割されるが、それにより入力データ配列において異なる部分配列に属するデータは、たとえ同じ区间に属していたとしても並行して転送することができる。この並行データ転送のために提案アルゴリズムでは二つの2次元($M\times L$)配列を用いる。まず入力部分配列 S_i に含まれ、 j 番目の区间に属するデータの集合を D_{ij} とする。このときカウント2次元配列 $C=\{c_{ij}\}$ を $c_{ij}=|D_{ij}|$ で定義する。またマップ2次元配列 $Q=\{q_{ij}\}$ を D_{ij} を保持すべき出力データ配列上の位置として定義する。

3. Map Sort

*Map Sort*の概要は下記のようになる。

Algorithm *MapSort*

Input: $S[0..N-1]$

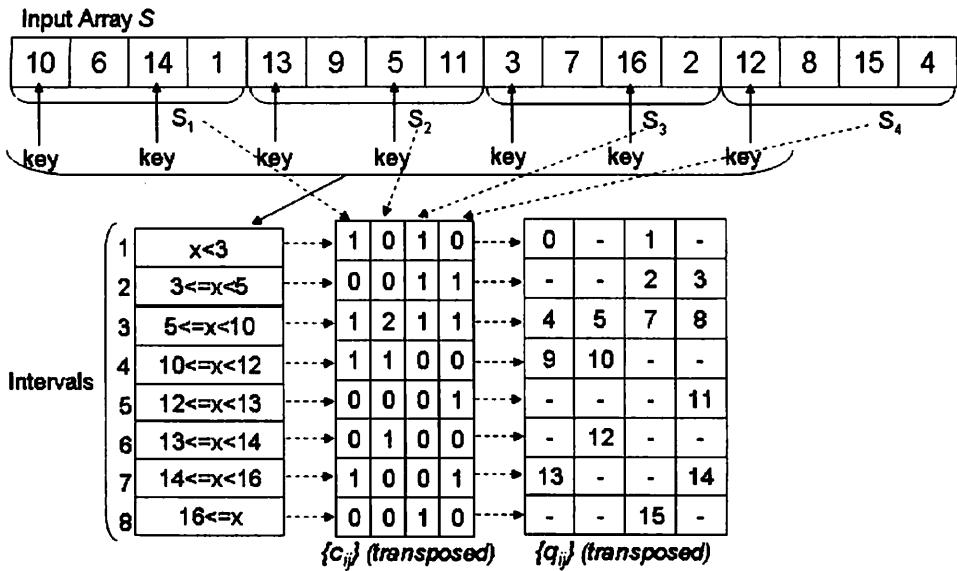
Output: $S'[0..N-1]$ (sorted array)

- Step 1: 入力データ S から $L-1$ 個の基準値を計算し、それらの基準値を境界とする L 個の区間を計算
- Step 2: データ $d \in S_i$ に対し、 d を含む区間（ j 番目とする）を求め、カウント2次元配列 c_{ij} の値を更新（部分集合 S_i ごとに並行処理可能）
- Step 3: P 個のCPUにより Q を計算
- Step 4: データ $d \in S_i$ に対し、 d を含む区間（ j 番目とする）を求め、出力データ配列 S' 上の位置 $q_{ij} \in Q$ に d をコピーし、 q_{ij} に1を加算（部分集合 S_i ごとに並行処理可能）
- Step 5: 出力データ配列 S' 上の各区间に対応する部分配列をソート（各区間のソートはクイックソートなど単体CPU向けに最適化されたソートアルゴリズムを用い、区間ごとに並行処理可能）

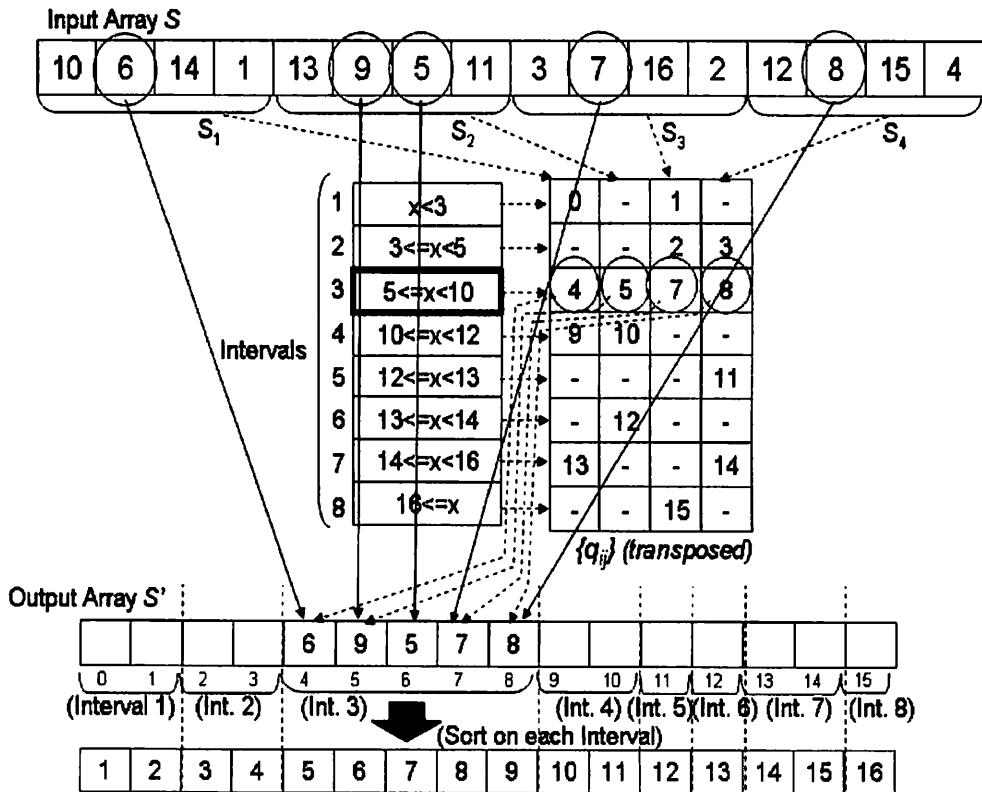
Fig. 1 (a) に示すような16データの例を用いて提案アルゴリズムを説明する。ここで $M=4$ 、 $L=8$ とする。Step 1では S から7(= $L-1$)個の基準値(key)を選び8個の区間(Interval)を作る。Step 2では入力データ S を4(= M)分割し、各部分配列 S_i に対し各区間にデータ数 c_{ij} を数え上げる。Step 2の処理は各部分配列 S_i ごとに並行処理可能である。

そしてStep 3では q_{ij} を次のように計算する。 c_{ij} は部分集合 S_i に属し j 番目の区間に含まれるデータ数である

ため、 j 番目の区間に含まれるデータの総和は $\sum_{i=1}^M c_{ij}$



(a) Input data S , $(L-1)$ keys, L intervals, c_{ij} and q_{ij} for $M=4$, $L=8$. No value for q_{ij} when $c_{ij}=0$.



(b) Data in S_i within j -th interval is moved to the position of output array S' indicated by q_{ij} . Then, each sub-array of the output array associated with intervals is sorted.

Fig. 1. Proposed Algorithm.

となる。Fig. 1の例では、1番目の区間に2個、2番目の区間に2個のデータが存在する。従って出力データ配列

上で j 番目の区間に含まれるデータは $\sum_{j'=1}^{j-1} \sum_{l'=1}^M c_{l'j'}$ から

開始されることになる。Fig. 1の例では1番目の区間を出力データ配列の位置0、2番目の区間を位置2、3番目の区間を位置4から開始すればよい。これより j 番目の区間に含まれる部分集合 S_j に属するデータを

$\sum_{j'=1}^{j-1} \sum_{l'=1}^M c_{l'j'} + \sum_{l'=1}^{i-1} c_{lj'}$ から開始すると重なりなくデータ

転送でき、これを q_{ij} とする。Fig. 1の例では上述のように3番目の区間は位置4から開始するため $q_{1,3}=4$ となり、 S_1 に含まれるデータはここから開始する。そして $c_{1,3}=1$ であるため $q_{2,3}=5$ となり、 S_2 に含まれるデータは位置5から開始すればよい。なお図では、 $c_{ij}=0$ のときには S_i に含まれる j 番目の区間に属するデータが存在しないため q_{ij} を記載していない。Step 3の処理は $M \times L$ 行列上の集計処理であり、並行処理が可能である。ただし M 、 L の値が小さいため単体プロセッサで実行した方が高速になる場合も多い。

次に Step 4において、入力データ配列 S 上のデータをすべて出力データ配列に転送する。例えば部分配列 S_i に属するデータが j 番目の区間に含まれるとき、出力データ配列上の位置 q_{ij} にコピーすると共に (Fig. 1(b))、 q_{ij} に 1 を加算する。このとき Step 3 の計算方法により、 q_{ij} と $q_{i+1,j}$ (または $q_{1,j-1}$) の間のデータ数が c_{ij} であることから、すべてのデータは重なることなくコピーできることが重要である。Fig. 1の例では、上述のように3番目の区間に含まれるデータは、 S_1 は位置4から、 S_2 は位置5から、 S_3 は位置7から、 S_4 は位置8からデータ転送することにより重なりなく転送できる。このため Step 4 についても各部分配列 S_i ごとに並行処理可能である。

さて、Step 4 の終了後は、出力データ配列は各区間に応じて部分配列に分けられることがわかる。よって Step 5においてそれぞれの区間に応じて部分配列に対して複数プロセッサを用いて並行してソートすることにより、全体のソートが完了する (Fig. 1 (b))。

このアルゴリズムにおいては基準値の選び方が重要である。なぜならば基準値の選び方により Step 4 の結果生成される L 個の部分配列に偏りができ、Step 5 を並列実行した場合の負荷バランスが悪くなるためである。そのため、例えば L 個よりも多くのサンプルを取得し、簡単な統計処理を行った上で基準値を計算する方が結果的に高速になる場合が多い。また、特に元のデータ

が一様分布、正規分布など既知の場合には効果が大きくなる。

次に計算複雑度について考える。 L 、 M を $O(P)$ とする。 P は実際に定数であるが $O(\sqrt{N})$ とする。空間複雑度は、出力データ配列に $O(N)$ 、マップ情報に $O(P^2)$ であるため、全体で $O(N)$ となる。

時間複雑度は、Step 1 は選択アルゴリズムにも依存するが $O(P)$ 個のデータを選びソートして区間を設定するため $O(P \log P)$ となる。Step 2 において各部分配列は $O(N/P)$ の長さであり、各データに対して区間を計算する処理に $O(\log P)$ 必要であるため、 P 台の CPU を用いて $O((N/P) \log P)$ となる。Step 3 におけるマップデータ構造の作成は、 P 台の CPU を用いて $O(P)$ ができる。Step 4 も Step 2 と同様の処理になるため、 P 台の CPU を用いて $O((N/P) \log P)$ となる。最後に Step 5 は、各区間に応じて部分配列が平均的に $O(N/P)$ の長さであるため、ソート処理は $O((N/P) \log (N/P))$ となるが、 P のオーダーを考えると $O((N/P) \log N)$ となる。従って全体の時間複雑度は $O((N/P) \log N)$ となり、ソート処理の時間複雑度 $O(N \log N)$ の $(1/P)$ となる。

4. 計算機実験

提案アルゴリズムを評価するため、OpenMP ディレクティブを用いた C 言語を用いてプログラム化し、Intel クアッドコアプロセッサ QX6700 (2.66GHz) を用いたパソコン上で計算機実験を行った。各実験においては乱数により生成された整数データ集合を 10 通り用意し、平均を取ることにより実行時間を測定した。今回の実験ではデータ数は 10^4 ~ 10^7 とした。また比較のためクイックソートもプログラム化した。単体 CPU 上のクイックソートは [3] に基づくプログラムとし、複数 CPU の場合には再帰的な配列分割をスレッド化することにより並列化を行った。なおスレッドを数多く発生させるとオーバーヘッドが大きくなり、かえって実行時間が遅くなるため、最大スレッド数を制限するようにした。

Fig. 2 は $M=L=32$ における提案アルゴリズムとクイックソートの実行時間を示している。Fig. 3 では P を変化させたときの Map Sort の実行時間を計測した。Map Sort は、2PE (Processing Element) の時に約 1.8 倍、4PE では約 3.5 倍の並列性能を達成した。Fig. 4 ではクイックソートとの比較を示した。単体 CPU 上のクイックソートと比較して、並列クイックソートは 4PE で 1.6 倍であるのに対し、提案アルゴリズムは 4PE で約 3 倍の性能向上を達成した。また Map Sort では、バラメタ M 、 L に対し、単体 CPU での実行速度と並列加速性能がトレ

ードオフの関係にあるが、Fig. 5は L の値に対する変化を示している。 L の値が小さいとき、単体CPUでの実行時間は短いが並列加速は悪くなる。 L の値が大きいと並列加速は良くなってくるが、オーバーヘッドが大きくなるため単体CPUでの実行時間は遅くなる。また、さらに L の値を大きくすると、単体CPU、複数CPUのどちらの性能も悪くなる。この結果より、 P の値に応じて M 、 L の値を変えることにより優れた性能が達成できることがわかる。例えば $P=1$ のとき $L=1$ 、 $P=2$ のとき $L=16$ 、 $P=4$ のとき $L=32$ とすることにより、単体CPUではクイックソートと同等性能、4 CPUでは3倍の性能を達成し、スケーラブルなアルゴリズムとなることがわかる。

5. まとめ

本論文では、新しい並列ソートアルゴリズムMap Sortを提案した。このアルゴリズムはマルチコア上でスケーラブルな性能を達成する。 P 台のCPUを用いたとき単体CPU上のソートアルゴリズムと比較して時間複雑度が (I/P) となり $O((N/P) \log N)$ となる。また計算機実験の結果、単体CPU上のクイックソートと比較し、単体CPUでの性能は同等で、複数CPUでは4 CPUで3倍の性能を達成できることがわかった。

文 献

- [1] Intel, *Multi-Core Processors. Making the Move to Quad-Core and Beyond*, White Paper, 2006.
- [2] P. Middleton, "A New ARMv6 Symmetric Multiprocessing Core", *Embedded Processor Forum*, 2004.
- [3] R. Sedgewick, *Algorithms in C*, Addison-Wesley, Reading, MA, 1990.
- [4] D. Bitton, D. J. DeWitt, D. K. Hsiao, and J. Menon, "A Taxonomy of Parallel Sorting," *ACM Computing Surveys*, Vol. 16 (1984), No. 3, pp.287-318.

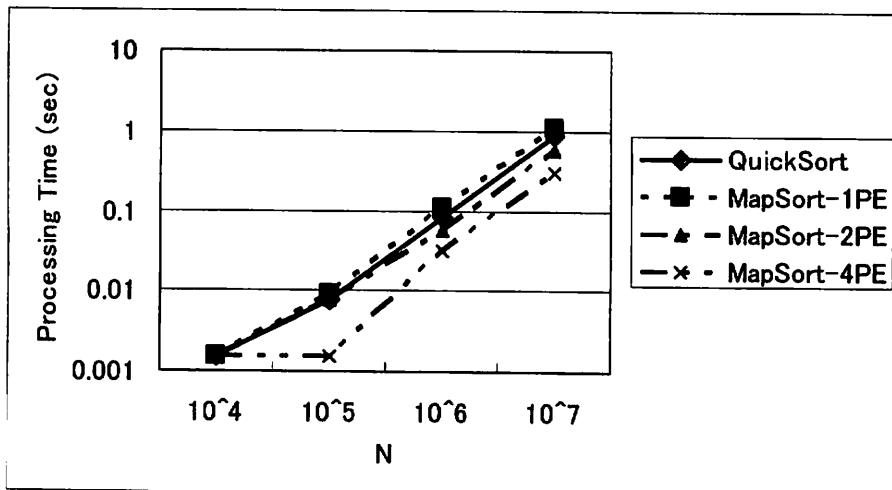


Fig. 2. Processing Time for Quick Sort on 1 PE and Map Sort

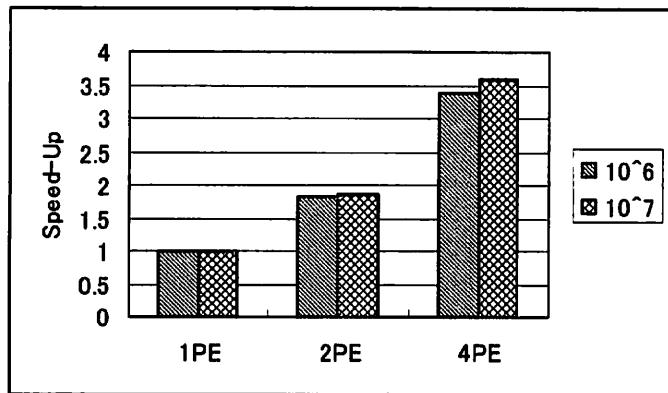


Fig. 3. Speed-Up for Map Sort (Normalized by Map Sort-1PE)

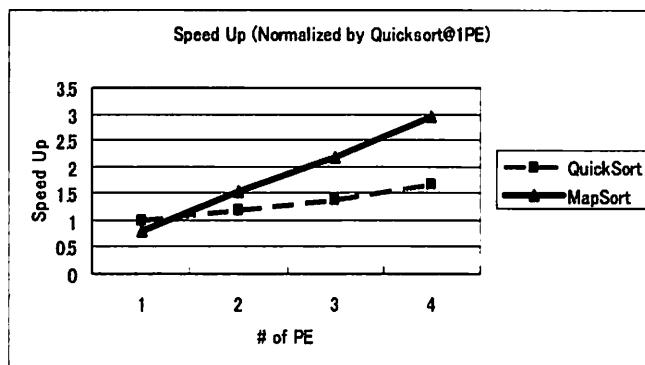


Fig. 4. Speed-Up for Quick Sort and Map Sort
(Normalized by Quick Sort-1PE, $N=10^7$)

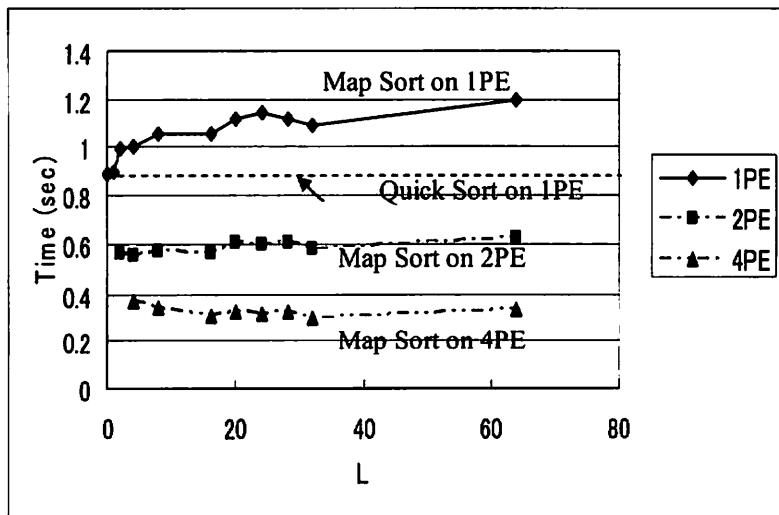


Fig. 5. Processing Time for several L values ($N=10^7$)