

A Study on Speed-Up of Sorting Algorithm
by Precompressing Data嶋田 強志[†] 上土井 陽子[‡] 吉田 典可[†]

Tsuyoshi SHIMADA Yoko KAMIDOI Noriyoshi YOSHIDA

1 はじめに

近年、記憶装置の大容量化に伴いデータベースシステムなど、格納される情報量も莫大なものとなっており、データの格納方法がその後のデータベース上での処理効率に影響を与える。

本研究では単一プロセッサを仮定し、データを事前にハッシュ関数を使用して、バケットへ分配する格納方法でデータを圧縮し、圧縮データを連結し、連結データを操作することで並列操作を実現する高速なソート技法の開発を目指す。

本稿では、提案アルゴリズムについて示し、提案アルゴリズムの部分列ソートについて実験及び考察する。提案アルゴリズムは1) データの事前圧縮、2) 圧縮データの連結、3) 圧縮データを利用したソート手続きから構成される。

2 基本となる従来アルゴリズム

本研究の基本となる従来ソートアルゴリズム Range Reduction[1] と Merging Two Words[2] について説明する。

2.1 Range Reduction アルゴリズム

Range Reduction は b ビット、 n 個の要素のソート問題を $b/2^k$ (k は 1 以上の整数) ビット、 n 個の要素のソート問題に帰着させるアルゴリズムである。入力データ $X = \{x_0, x_1, \dots, x_n\}$ が与えられたとき、ハッシュ関数 $H(x_i) = \lfloor x_i/2^k \rfloor$ で x_i を圧縮し、圧縮要素を $x_i \bmod 2^k$ でインデックス付けられたバケットに格納する。ソートはハッシュ関数で圧縮された要素 $\lfloor x_i/2^k \rfloor$ に対して行なわれる。Range Reduction ソートアルゴリズムにおいて、 $S(n, b)$ を b ビットの要素 n 個のソートに対する計算複雑さとする。このとき、以下の不等式が成り立つ。

$$S(n, b) \leq S(n, \lfloor b/2 \rfloor) + O(n)$$

ここで、 $O(n)$ は Range Reduction ソートのコストである。Range Reduction アルゴリズムは、要素を構成するビット数が小さいと有効なソートに対して効果的であるが、値域が大きい場合、Range Reduction を用いなかった場合とコストは同じである。

2.2 Merging Two Words アルゴリズム

Merging Two Words は 2^k 要素からなる双調列データを連結し、連結データ上で 2^{k-1} の要素対に対する k 回の比較・交換により、バイトニックソート[2] を実行するアルゴリズムである。連結データの各領域には右から左にアドレスが付けられ、各領域の最上位ビットはフラグである。

図 1 に Merging Two Words における比較操作について示す。比較は入力列 Z を各グループごとにアドレス

ビットの最上位から最下位へ、参照するビットを変えながら 1 と 0 で分割し、分割データ A, B を得る。次に A のフラグを 1 にセットし A と B の減算を行なう。減算の結果、フラグが 1 なら A に属すデータが大きい。交換は減算結果のフラグを基に作成された論理マスクを用い、シフトと論理演算により行なわれる。Merging Two Words の計算複雑さは入力要素数を m とおくと $O(\log m)$ である。Merging Two Words は単一プロセッサ上で並列操作を可能とするが、計算機上で実現する場合に、連結データサイズが大きくなり実装が困難である。例えば、int 型データ (4Byte) を 8 個連結した場合、連結データサイズは 256 ビットとなる。また、連結データを構成する要素のみに対応するソートしかできない。連結できる要素数は演算器が使用するレジスタの大きさに依存する。

address	111	110	101	100	011	010	001	000	
Input	0: 7	0: 2	0: 4	0: 9	0: 12	0: 15	0: 10	0: 8	Z
	00111	00010	00100	01001	01100	01111	01010	01000	
					1: 7	1: 2	1: 4	1: 9	A
					10111	10010	10100	11001	
-					0: 12	0: 15	0: 10	0: 8	B
					01100	01111	01010	01000	
=					0: 0:	0: 0:	0: 0:	1:	
					01011	00011	01010	10001	

図 1: 連結データ上での要素の並列比較

3 提案アルゴリズム

Range Reduction と Merging Two Words を基礎とする理論上のソートアルゴリズムとして、Signature Sort[3] が存在する。本稿では、Range Reduction と Merging Two Words を基本とし、計算機上での実現を試みる。

3.1 アルゴリズムの概要

Merging Two Words を計算機上で実現する場合に、連結データサイズが大きくなり実装が困難で、大きな入力データ数に対応できない。また、Range Reduction では、データの取り得る値を減少 (圧縮) させてもソートの十分な高速化が得られない。提案手法では、Range Reduction を複数回適用し、入力データを計算機上で最も効率的にデータアクセス可能だと考えられる要素単位 1Byte に圧縮し、圧縮データを連結する。次に圧縮・連結データを拡張 Merging Two Words によりソートし、Range Reduction の逆操作により元のデータを得る。

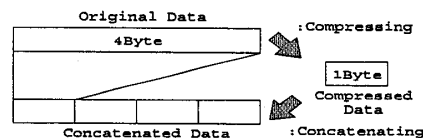


図 2: データ要素の圧縮と連結

演算器が直接操作するデータ (図 2) について考える。ある要素が 4Byte の int 型であるとき、この要素を圧縮し 1Byte の圧縮要素を得る。圧縮要素を 4 個連結し

[†]広島市立大学大学院 情報科学研究科 (Graduate School of Information Sciences, Hiroshima City University)

[‡]広島市立大学 情報科学部 (Faculty of Information Sciences, Hiroshima City University)

た場合、圧縮・連結データは int 型と同じ大きさとなる。圧縮要素に情報の損失がなければ、要素を演算器が直接操作することで 4 要素を並列に処理可能で、演算時間が短縮される。本稿では圧縮・連結データを (m, f) Word と呼ぶ。 m は連結した要素数、 f は連結する各要素のビット数である。圧縮・連結データサイズは $m \times f$ である。図 2 の例では $(4, 8)$ Word である。

3.2 Merging Two Words の拡張

拡張 Merging Two Words は、 $(m, 8)$ Word (ここで m は 2 の巾乗) の部分列を入力要素とする。また、部分列ソートに Merging Two Words を使用し、分割・併合しながらバイトニックソートを行うことで要素数の大きなデータをソートする。

図 3 に $m=16$ の Merging Two Words を部分列ソートとする拡張 Merging Two Words を示す。図 3 の入力データ列はすでに Range Reduction により 1Byte に圧縮されている。最初に入力列を 16 要素ごとに連結する。次に 1st step では入力列がランダムであるため、長さ 2 の双調列を 8 個、昇順、降順にソートし、長さ 4 の双調列を 4 個生成し、昇順、降順にソートする。これを繰り返し、長さ 16 の双調列で Merging Two Words をそれぞれ実行する。 2nd step 以降では通常の Merging Two Words による部分列のソートを行なう。また、各 Merging Two Words の出力は半分に分割・併合され新たな $(16, 8)$ Word の双調列とみなせ、この双調列に対する Merging Two Words によるソートを逐次的に繰り返し、全体をソートする。

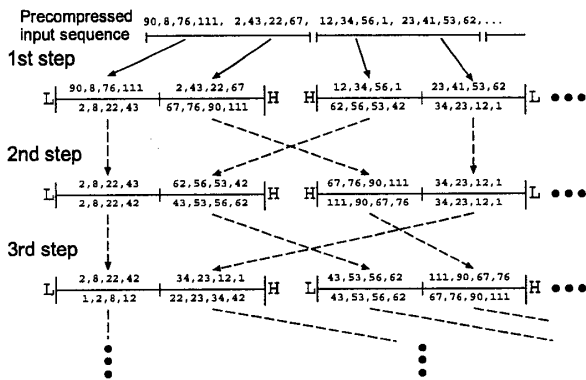


図 3: 分割・併合によるソート

3.3 提案アルゴリズムの計算複雑さ

提案アルゴリズムの計算複雑さについて示す。未圧縮の全入力データ要素数を n 、圧縮データ要素の連結数を m ($n > m \geq 2$) とおく。 Range Reduction の計算複雑さは $O(n)$ である。また、 Merging Two Words による交換コストは 1st step が $\log^2 m$ で 2nd step 以降の Merging Two Words による交換コストは $\log m$ である。 $M(m, n)$ を連結数 m 、要素数 n のときの拡張 Merging Two Words の計算複雑さとする。このとき、以下の関係式が成り立つ。

$$M(m, n) = \log^2\left(\frac{2n}{m}\right) \cdot \frac{n}{m} \cdot \log m + \frac{n}{m} \cdot \log^2 m$$

$$\leq \frac{\log m}{m} \cdot n \cdot (\log^2 n + \log m)$$

$m = 1$ のとき上界、 $U(1, n) = n \log^2 n$ 、
 $m = n$ のとき下界、 $L(n, n) = \log^2 n$ となる。

よって、以下の不等式が成り立つ。

$$n \log^2 n > M(m, n) > \log^2 n$$

Range Reduction の適用回数を C とすると、提案アルゴリズム全体の計算複雑さ $T(m, n)$ は、以下の等式

により表される。

$$T(m, n) = M(m, n) + C \cdot O(n)$$

以上より m が大きい程、ソート時間は短縮される。

4 実験と考察

提案アルゴリズムの部分列ソート Merging Two Words について C 言語を用いて実現し、 2 つのソート列 (双調列) を併合する Merge Sort (計算複雑さは $O(n)$) 及び Quick Sort (計算複雑さは $O(n \log n)$) と比較した。使用計算機は Ultra SPARC-III、要素数 16,8,4 をもつ双調列な入力データに対し、アルゴリズムを 10^6 回繰り返し適用し、CPU 時間を計測した。連結データ数が 16、8 のとき使用言語の制約上、アルゴリズムに忠実な実装ができないため連結データ数 4 の Merging Two Words を再帰的に使用し実現した。入力データは双調列である。実験結果を表 1 に示す。

表 1: 各手法の計算時間 [s]

	入力データの要素数		
	4	8	16
Merging Two Words	0.77	2.05	4.57
Merge Sort	0.57	1.21	2.33
Quick Sort	1.34	3.56	14.24

拡張 Merging Two Words は部分列ソート Merging Two Words の並列操作性を利用したソートアルゴリズムである。そのため、 Merging Two Words の実行時間が Merge Sort より遅ければ、部分列ソートに Merge Sort 適用すればよいということになる。また、 Merge Sort を部分列ソートとする $O(n \log n)$ のソートアルゴリズム [4] が知られているため、 Merging Two Words は Merge Sort より実行が速い場合、高速化の可能性がある。 Merging Two Words は $m = 8, 16$ のとき、言語の制約により忠実な実装ができなかったため、 Merge Sort より全てのデータで遅い結果となった。アルゴリズムに忠実な実装が可能であれば、理論上で実行時間は $m = 8$ のとき 3/5 倍、 $m = 16$ のとき 1/3 倍程度に短縮される。

5 まとめ

実験結果より提案アルゴリズムの部分列ソート Merging Two Words について、連結データ数以内でアルゴリズムに忠実な実装が可能なら Merge Sort より高速であることが分かった。提案アルゴリズムでは、単一プロセッサ上の圧縮・連結データの並列操作によるソートの高速化を試みたが、単一プロセッサのソートの最適計算時間 $O(n \log n)$ を上回ることができなかった。今後の課題として、提案アルゴリズム全体の実現及びその実験的評価、並びに部分列ソート Merging Two Words の効率的な実現のためのハードウェアの構築、部分列ソート Merging Two Words を使用した新たな分割・併合ソートの開発が挙げられる。

参考文献

- [1] D.Kirkpatrick and S.Reisch, "Upper bounds for sorting integers on random access machines", Theoret. Computer Science, Vol. 28, pp. 263-276, 1984.
- [2] A.Albers and T.Hagerup, "Improved parallel integer sorting without concurrent writing", Proc. 3rd Annual ACM-SIAM Symp. on Discrete Algorithms, pp. 463-472, 1992.
- [3] A.Andersson, T.Hagerup, S.Nilsson and R.Raman, "Sorting in Linear Time?", Proc. 27th Annual ACM Symp. on Theory of Computing, pp. 427-436, 1995.
- [4] R.Sedgewick, "Algorithms IN C++ THIRD EDITION", 2000.