

パイプラインマージソータに於ける 4U-4 String Length Tuning用フラグ自動生成機構

楊 維康 喜連川 優 高木 幹雄

東京大学生産技術研究所

1. はじめに

我々はパイプラインマージソータのアルゴリズム及びその実装法についての研究を行い、既に18段のソートプロセッサから構成される大容量ハードウェアソータを試作し[1]、その試作システムでソータの動作確認と性能評価を行った。又、そのソートプロセッサのLSI化を行い、1チップのソートプロセッサを実現した[2]。そのLSIソートチップを用いて、大容量(最大64MB)、高速(最大8MB/Sec.)のハードウェアソータをコンパクトに実現することができる。本ソートチップの最大の特長は我々が開発したSLT機構を始め、幾つかの拡張機能を有しており、それらの機能により、本チップで構成されたソータは、ソートするファイルのレコード長、レコード数等のパラメータに対して、非常に高い柔軟性を有している。それらの拡張機能を利用する為には、ソータに入力されるデータストリームの各レコードの先頭に制御フラグを付加し、必要な制御情報をソートプロセッサに知らせる必要がある。又、本ソータの動作はデータストリーム指向で、ディスク等からのデータを上位プロセッサの主記憶等に転送する時間を利用して、On-the-flyにソート処理を行うことが可能である。そのような処理ができるようにする為には、フラグの生成をOn-the-flyに行う自動生成機構が不可欠である。

本稿では、そのような制御フラグをハードウェアで自動的に生成するアルゴリズムを紹介し、その実装法について述べる。

2. SLT 及び制御フラグの種類

SLTはソートするファイルのレコード長がソータの設計レコード長と異なる時、ソータの最初の数段のソートプロセッサに於いてダイナミックに出力するストリング(ソートされた部分列)の長さを調整して、ある中間プロセッサのメモリを最大限に利用できるようなストリングを生成し、それによって、ソータ全体のメモリ利用効率をほぼデータのレコード長に影響されないように最適化するアルゴリズムである。そのアルゴリズムによって、レコード長に対する柔軟性を実現した。ストリング長の

調整は図1のマージ木で示すアルゴリズムに従って行い、その実現は、初段のソートプロセッサに於いてレコードを適宜マージせずにバイパスすることによって実現されている(図1の終端で兄弟のないノード)。図の中では、 d はチューニングの次数と言い、 N_d は P_d が出力するストリングのレコード数である。

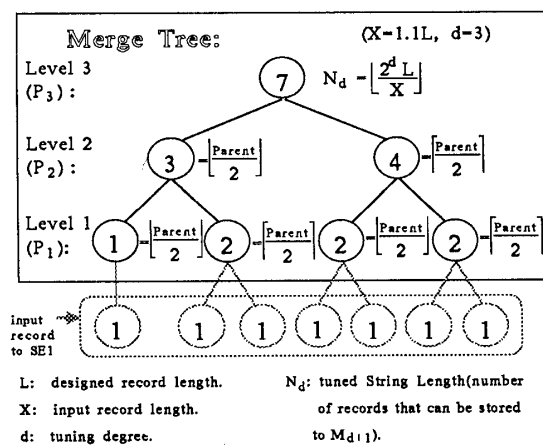


図1. マージ木の例

本ソータに於いて、その他の拡張機能の実現に合わせて、次のようなフラグを導入している。

- BOS ——各ストリングの先頭レコードを示す。
- LS ——データストリーム中、最終ストリングの先頭レコードを示す。
- EOS ——データストリームの終了を示す。
- NMP ——SLT動作時、初段のプロセッサに於いてマージペアのないレコードを示す。
- SBP ——このフラグ以降LSフラグまでのデータをバイパスすることを指示する。
- DR ——SLTの動作中に内部的に生じたダミーレコードを示す。

SLTを行わない場合、ソータに入力するデータの各レコードにBOSフラグを付加し、ソートプロセッサはBOSフラグのレコードをマージする。SLT動作時、初段でバイパスするレコードにはNMPフラグを付加する。BOSとNMPのシーケンスによって、SLTが実現されている。

3. SLT用フラグの自動生成アルゴリズム

次数 d のSLTを行う場合、 N_d レコードからなるストリングを生成するマージ木は、深さ $d+1$ 、 N_d 個の葉ノード

Flag Generator for String Length
Tuning in Pipeline Merge Sorter

W. Ynag, M. Kitsuregawa, M. Takagi
Institute of Industrial Science,
University of Tokyo

ドを持つバランスした2進木である。一般的に、 $2^{d-1} < N_d < 2^d$ であり、言い換えれば、そのマージ木は深さd+1の完全2進木の 2^d 個の葉ノードの内、 $2^d - N_d$ 個をバランスして切り取ったものである。但し、そのマージ木の左部分木の葉ノード数は右部分木のそれと等しいか、又は1個少ないことがある。この条件はマージ木に対して再帰的に適用できる。その切り取られた葉ノードは次のように決定することができる。

深さd+1の完全2進木の 2^d 個の葉ノードを左から $(2^{d-1}-1), (2^{d-1}-2), \dots, 0$ と番号をつける。番号iを2進数で表し、そのiの2進表現のMSBをLSBとして得られた値を*i'*とする。*i'* $\geq N_d$ の葉ノードの右兄弟を切り取る。図2はその例を示す。

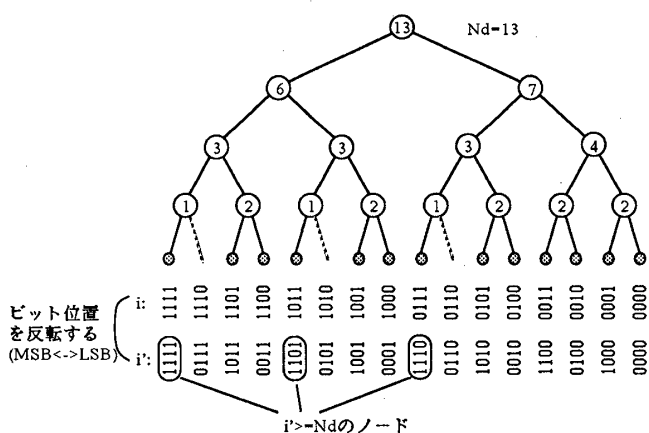


図2. SLT用フラグ生成アルゴリズムの例

i' $\geq N_d$ の葉ノードは初段ではマージペアがなく、バイパスされるレコードに相当する。故に、そのノードにはBOSフラグの代わりにNMPフラグが付加される。

以上をまとめると、BOSとNMPフラグのシーケンスは次のようなアルゴリズムで簡単に生成することができる。

BOS, NMPフラグを生成するアルゴリズム:

RNC: 入力レコード数をカウントするカウンター

RNC': RNCのMSBとLSBを反転した値

BEGIN

DO

RNCを'11..1'にリセットする;

LAB1: IF (RNC' < N_d) THEN

BOSフラグを出力する;

ELSE /* RNC' $\geq N_d$ の場合 */

NMPフラグを出力し、RNCを1回カウントダウンする;

RNCを1回カウントダウンする。

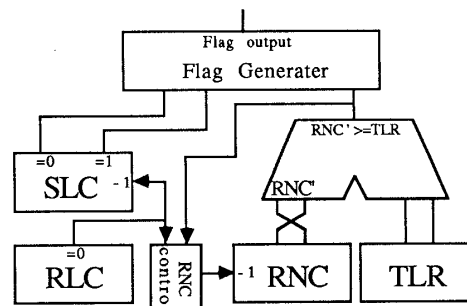
IF (RNC ≥ 0) GOTO LAB1;

UNTIL(入力終了)

END

4. フラグ自動生成機構のハードウェア構成

図3にフラグ自動生成機構のハードウェア構成を示す。



RLC: Record Length Counter
RNC: Record Number Counter
SLC: Stream Length Counter
TLR: Tuning Length Register

The Function of Flag Generator:

SLC=0	SLC=1	RNC' \geq TLR	Flag
0	0	1	NMP
0	0	0	BOS
0	1	X	LS
1	0	X	EOS

図3. フラグ自動生成機構のハードウェア構成

この図では、RNCはフラグ生成用のレコード数カウンタで、SLCはストリームの残りレコード数、RLCは処理中のレコードの残りバイト数を示すカウンタである。これらは皆ダウンカウンタである。RLCの初期値はレコード長で、1バイト入力する度にカウントダウンする。RLC=0の時、データストリームにフラグを付加する処理を行い、RLCをレコード長にリセットする。TLRは N_d の値を保持するレジスタで、RNCのMSBとLSBを反転した値RNC'と比較される。その比較結果によってSLT用のフラグを生成する。

図3でわかるように以上のアルゴリズムを実現するハードウェアは非常に簡単で、データストリーム入力と同時に、フラグを生成することができる。簡単にする為に、各レジスタをセット、リセットする部分や、SBPフラグを生成する部分はこの図に示していない。

5. むすび

本稿では、SLTを始め、幾つかの拡張機能を有するパイプラインマージソータ用の制御フラグを自動的に生成するハードウェアのアルゴリズム及びその構成について述べた。本ハードウェアを用いて、ソータ用のフラグをOn-the-flyに生成でき、ディスク等からのデータ流に沿ってソート処理を行うことを可能にした。このような回路をソータ駆動系に組み込み、実用可能なシステムを実現するのが今後の課題である。

<参考文献>

[1] 楊, 鈴木, 喜連川『高速(4MB/Sec.)大容量(8MB)ハードウェアソータの実装』信学技報 CPSY 86-26, 1986
[2] 楊, 喜連川『LSIソータチップの試作』第37回情報処理学会全国大会 7Q-4