

## 連続メディア配信システム:S<sup>3</sup>におけるディスクスケジューリング方式

帆波 幸二<sup>†</sup> 浅見 和男<sup>†</sup> 高野 了成<sup>†</sup> 吉澤 康文<sup>‡</sup>

東京農工大学大学院工学研究科<sup>†</sup> 東京農工大学工学部<sup>‡</sup>

### 1. はじめに

連続メディア配信サーバには、大容量のメディアデータを多数のクライアントに配信するため、高速な I/O 機能が必要である。また、メディアに課せられた時間制約を満たすためにリアルタイム処理の機能も必要となる。そこで、我々は Linux をベースに連続メディアの配信に特化させたオペレーティングシステム:S<sup>3</sup>を開発している。

本報告では S<sup>3</sup> の概要と、実装したリアルタイムディスクスケジューラについて述べる。本スケジューラは、最悪の入出力時間を用いた悲観的なスケジューリングを廃止するために、アクセス時間の予測を用いる。予測はシステム起動時に取得したディスクの性能データをもとに、I/O リクエストに対して動的に見積もり時間を付加することで実現される。

以下、本稿では、2章で S<sup>3</sup> の概要、3章でディスクスケジューラ、4章でスケジューラの評価について述べ、5章でまとめと考察について述べる。

### 2. S<sup>3</sup> の概要

S<sup>3</sup> は通常の Linux に、連続メディア配信支援のためのメモリ管理、ディスク管理、ネットワーク管理の3つの管理モジュールを付加した形で構成される。メモリ管理では、クライアントごとに周期的に行われるシーケンシャルアクセス特性に基づいた新しいメモリ管理手法を提供し、ディスク管理では、見積もりに基づいたリアルタイムディスクスケジューリングを行う。ネットワーク管理はゼロコピープロトコルスタックによるパケット生成時のコピー削減と、周期駆動による割込みの削減を実現する。これらの機能により、配信スループットの向上と I/O レートの保証を目指す。

### 3. S<sup>3</sup> ディスクスケジューラ

S<sup>3</sup> におけるディスクスケジューラの最大の特徴はディスクアクセスの処理時間を見積もり、スケジューリングに利用することである。リアルタイムシステムにおいてはデッドラインミスを防ぐため、ワーストケースの処理時間を想定してスケジューリングを行うのが一般的である。しかし、磁気ディスクの

場合、1回のアクセス時間は、数百マイクロ秒から数十ミリ秒までとばらつきが大きく、常にワーストケースを想定すると、スケジュール可能なリクエストの数が大きく減少してしまう。本スケジューラは、リクエストの前後関係からアクセス時間の見積もりを行うことでこの問題を緩和し、リアルタイム性を維持しつつディスクの利用率を向上させるアプローチをとる。

#### 3.1. S<sup>3</sup> ディスクスケジューラの構成

S<sup>3</sup> ディスクスケジューラはリクエストインタフェース、ディスクスケジューラ本体、性能測定モジュールから構成される(図1参照)。

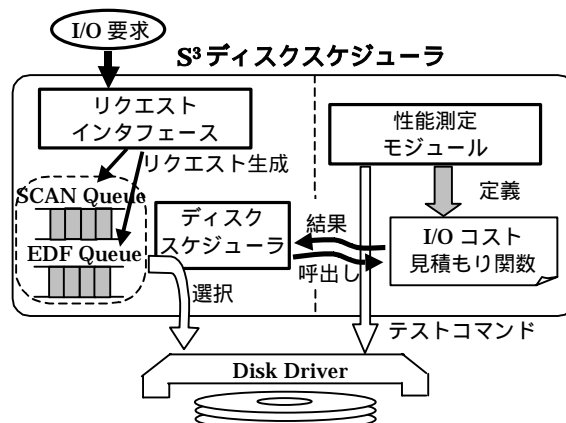


図1 S<sup>3</sup>ディスクスケジューラの構成

#### (1) リクエストインタフェース

Linux のバッファキャッシュからの I/O 要求と、S<sup>3</sup> メモリ管理からのデッドラインを持つ I/O 要求を受け付け、リクエストキューを生成する。デッドラインを持つリクエストは、デッドラインの早い順に並べた EDF(Earliest Deadline First)キューに、そうでないリクエストはセクタ順に並べた SCAN キューにそれぞれつながれる。

#### (2) ディスクスケジューラ

I/O の完了時にたびに実行され、リクエストキューの中から次に実行するリクエストを選択する役割を持つ。スケジューリングアルゴリズムについては次節で述べる。

#### (3) 性能測定ルーチン

システム起動時に実行され、磁気ディスクの性能情報を取得し、アクセス時間の見積もり関数の定義を行う。見積もり関数は、シーク時間、回転待ち時間、転送時間の3つについて定義される。

Disk Scheduling Using Access Time Estimates in S<sup>3</sup> System  
Koji Honami<sup>†</sup>, Kazuo Asami<sup>†</sup>, Ryousei Takano<sup>†</sup> and  
Yasufumi Yoshizawa<sup>‡</sup>

<sup>†</sup>Graduate School of Engineering, Tokyo University of  
Agriculture and Technology

<sup>‡</sup>Faculty of Engineering, Tokyo University of Agriculture  
and Technology

### 3.2. スケジューリングアルゴリズム

本スケジューラは、EDF キューの先頭から一つずつリクエストを取り出し、C-LOOK[1]の並び順になるようにスケジュールキューを生成するアルゴリズムを採用している。デッドラインミスを起こさない範囲でスケジュールキューの長さを最大にすることで、見積りによって得られたスケジューリングにおける余裕をシーク最適化に割り当てることを実現している。アルゴリズムのフローを図2に示す。

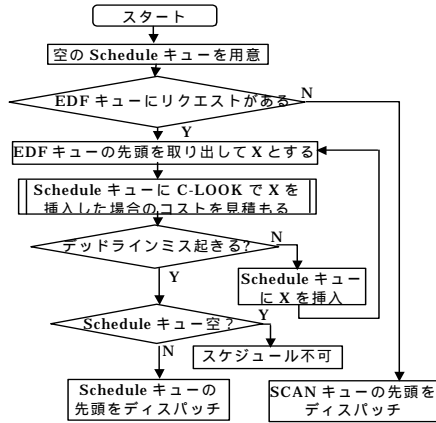


図2 スケジューリングアルゴリズム

### 3.3. アクセス時間の見積もり

ここでは、ディスクアクセス時間の中でもっとも大きな割合を占めるシーク時間の見積もり方法について述べる。シーク時間は、シークの移動距離に対する関数で示すことができる。そこで、システム起動時にシークの移動距離とシーク時間の対応を測定し、表を作成しておくことで、コストの見積もりを行っている。

実際に磁気ディスク(表1参照)にSEEKコマンドを発行して測定した結果を図3に示す。また、アクチュエータの加速度が一定であると仮定して導き出される見積もり関数を図3, 4に示す。

表1 ディスクのスペック

Fujitsu MPE3084AE	
Formatted Capacity	8.45GB
Disks / Heads	1 / 2
Bytes/Sector	512
Logical blocks	16,514,064
Drive Interface	Ultra EIDE (ATA/66)
Average Seek Time	9.5ms(typ)
Record Density	352,728bpi

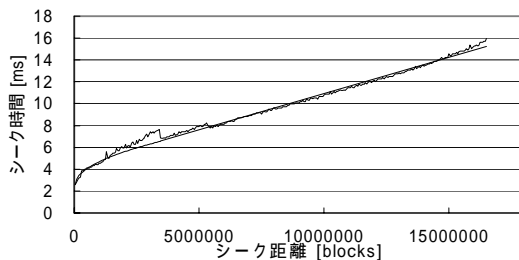


図3 シーク距離とシーク時間の関係

$$seektime(x)[ms] = \begin{cases} 0 & \text{if } x=0 \\ 2.6 + 3.5 \times 10^{-3} \sqrt{x} & \text{if } x \leq 3.4 \times 10^6 \\ 6.5 + 7.3 \times 10^{-7} (x - 3.4 \times 10^6) & \text{if } x > 3.4 \times 10^6 \end{cases}$$

図4 シーク時間の見積もり関数

### 4. スケジューラの評価

本スケジューラの有効性を示すために、ワーストケースを想定した場合との比較に加え、Linux で採用されている C-LOOK, リアルタイムスケジューリングアルゴリズムである SCAN-EDF[2]との比較も行った。実験は、仮想クライアントによりサーバに負荷をかけ、クライアント数とデッドラインミスの関係を測定した。その結果を図5に示す。

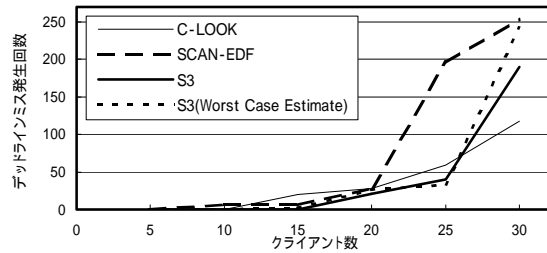


図5 アルゴリズムの比較

この結果から、本スケジューラは、デッドラインミスを起こさずにサポートできるクライアント数がSCAN-EDF 同様に C-LOOK よりも多くすることが可能で、なおかつ、SCAN-EDF よりもデッドラインミスの発生回数を小さく抑えていることがわかる。このことから、本スケジューラがリアルタイムスケジューラとして有効であることが確認できる。また、ワーストケースを想定した場合よりも 18%スケジュール可能性を向上させるという結果が得られた。

本スケジューラの欠点として、デッドラインミスが発生し始めると、ドミノ効果でデッドライン発生件数が増えてしまうという点が挙げられる。これは、負荷が大きくなるとスケジューラの挙動が EDF に近づくことに起因する。この問題はアドミッション制御を行うことで回避できると考える。

### 5. おわりに

S<sup>3</sup>のディスクスケジューラについて述べた。アクセス時間の見積もりを用いることにより、デッドラインミスの発生を防ぎ、ワーストケースを想定した場合よりスケジュール可能性を 18%向上させることを示した。これらの結果からリアルタイムディスクスケジューリングで見積もりを用いることの実効性を示すことができた。

### 参考文献

- [1] Abraham Silberschatz, Peter B. Galvin, "Operating system concepts," Addison Wesley, 4th ed., 1994.
- [2] A. L. N. Reddy and J. Wyllie, "Disk Scheduling in Multimedia I/O System," in Proceeding of ACM Multimedia'93, (Anaheim, CA), pp. 225-234, August 1993.