

外付け I/O エンジン方式を用いたストリーミングサーバの実現

竹内 理[†] レ・モアルダミアン[†] 坂東 忠秋[†]

高速アクセス網の普及とともに、高品質なストリーミングデータの有料配信サービスが注目されている。このサービスの実用化には、高い配信能力、配信品質の保証機能、既存の市販ストリーミングサーバとの機能互換性のすべてを兼ね備えたストリーミングサーバが必要である。本研究では、上記ストリーミングサーバの実現を目指し、外付け I/O エンジン方式と呼ぶストリーミングサーバの構成方式を新規に提案する。外付け I/O エンジン方式では、市販のストリーミングサーバと専用 OS 搭載サーバを連動させることで、上記要件を同時に充足するストリーミングサーバを少ない開発工数で実現する。さらに、Darwin ストリーミングサーバと、HiTactix (著者らが開発したストリーミングサーバ向け専用 OS) を搭載したサーバを実際に連動させ、外付け I/O エンジン方式の定量的な評価を行った。評価の結果、外付け I/O エンジン方式を用いることで、サーバに 9K 行程度のコードを追加するだけで、機能互換性を維持しながらも、HiTactix が提供する配信性能を保持するストリーミングサーバが実現できること等が明らかになった。

Implementation of a Streaming Server Using External I/O Engine Architecture

TADASHI TAKEUCHI,[†] DAMIEN LE MOAL[†] and TADAAKI BANDO[†]

With broadband networks expansion, charged streaming services of high quality stream data are likely to become common. In order to efficiently implement such services, streaming servers providing high data throughput performance, mechanisms to ensure streaming quality and compatibility with existing streaming servers are desirable. In this paper, we propose the external I/O engine architecture, which implements a streaming server as described above using both a conventional streaming server running on a general-purpose OS and an I/O engine server implemented over a special-purpose OS. With this combination, we can lower the necessary coding amount in order to implement the services provided. We implemented the server using the Darwin Streaming Server and HiTactix I/O engine server and evaluated the performances. We confirmed that the modified server can achieve HiTactix I/O performances without losing compatibility, with only 9K lines of code added.

1. はじめに

近年、光ファイバ網、CATV 網等の高速アクセス網が普及しつつある。これら高速アクセス網を提供している各キャリアは、自社が提供するアクセス網への加入者のさらなる増加を促進し、かつ、基本料金以外の追加収入をもたらすサービスの提供を模索している。特に、高品質なストリーミングデータの有料配信サービスは、上記サービスの有力候補として注目されている¹⁾。

しかしながら、高品質なストリーミングデータの有料配信サービスは、現在までのところ実用化に至っていない。その理由の一部として、既存の市販ストリーミングサーバに以下の問題があることがあげられる。

- (1) 既存の市販ストリーミングサーバのストリーミング配信性能が低いため、少ないユーザで 1 つのサーバを共有しなければならない。そのため、サービス提供コストが高くなる²⁾。
- (2) 既存の市販ストリーミングサーバはストリーミング配信の品質保証を行えない³⁾。そのため、有料サービス化が困難である (ユーザは課金されているにもかかわらず、満足いく品質でストリーミングデータの配信を受けられない可能性がある)。

しかし一方で、サービス普及のためには、既存の市販ストリーミングサーバとの機能互換性を維持し、多数のユーザが保有するストリーミングクライアントアプリケーションや、サービス管理者が使いなれている各種管理アプリケーション (ストリーミング配信モニタ等) を従来と同様に利用可能にした方が望ましい。

著者らは次世代ストリーミングサーバ向け専用 OS Hi-

[†] 株式会社日立製作所システム開発研究所
Systems Development Laboratory, Hitachi Ltd.

Tactix の研究開発を行ってきており、HiTactix を用いることで、低オーバーヘッド、かつデータ転送レートを保証した I/O 実行が実現可能であることを実証してきた^{4)~8)}。本研究の目的は、冒頭で述べた課題を解決するストリームサーバの実現、具体的には、既存の市販ストリームサーバとの機能互換性を維持しながらも、HiTactix のような専用 OS を利用した高性能、かつ配信品質を保証したストリーム配信を少ない開発工数で実現することにある。

上記目的の実現のため本研究では、「外付け I/O エンジン方式」という新しいストリームサーバの構成方式を提案する。「外付け I/O エンジン方式」では、既存の市販ストリームサーバと HiTactix のような専用 OS を搭載したサーバ（外付け I/O エンジン）を連動させる。既存の市販ストリームサーバには、外付け I/O エンジンと連動するための改変のみを加える。一方、外付け I/O エンジンは、従来市販ストリームサーバが行っていたストリーム配信を代理実行する。このような連動により、市販ストリームサーバとの機能互換性を維持しながらも、専用 OS を利用して、ストリーム配信の性能向上、および配信品質保証が実現できる。上記実現のために必要な開発工数は、既存の市販ストリームサーバの改変工数と、専用 OS 上で動作するストリーム配信処理の開発工数のみに抑えられる。

さらに本研究では、Darwin ストリームサーバを「外付け I/O エンジン方式」に対応させるべく改変を加え、HiTactix 搭載の外付け I/O エンジンと連動させる。そして、両サーバの連動のために必要となる開発工数が大きくなること、および連動オーバーヘッドに起因するストリーム配信性能の劣化が大きく発生せず、HiTactix が提供する配信性能を引き出せることを定量的に検証する。

以下、2章で外付け I/O エンジン方式の概要について述べ、さらにその実現方法を検討する。次に、3章で HiTactix の概要を示し、HiTactix を利用することにより、高性能かつ配信品質を保証したストリーム配信を行う外付け I/O エンジンの開発が容易なることを示す。さらに、4章で外付け I/O エンジン方式に対応させるために加えた Darwin ストリームサーバの改変内容について説明し、その改変のための開発工数が大きくなることを示す。最後に、5章で外付け I/O エンジン方式の定量的な性能評価結果について述べ、外付け I/O エンジンとの連動オーバーヘッドに起因

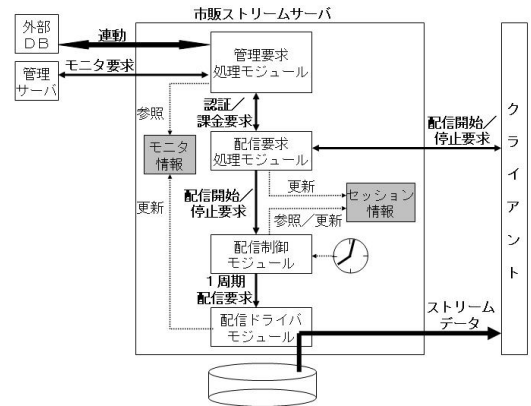


図1 市販ストリームサーバのモジュール構成
Fig. 1 Conventional streaming server modules.

するストリーム配信性能の劣化が大きく発生しないこと等を定量的に明らかにする。

2. 外付け I/O エンジン方式

本章では、著者らが新規に提案した「外付け I/O エンジン方式」、すなわち、既存の市販ストリームサーバにわずかな改変を加えるだけで、当該サーバのストリーム配信性能の向上、およびストリーム配信の品質保証機能の追加を可能にするストリームサーバの構成方式について述べる。

まず、市販のストリームサーバの典型的なモジュール構成について述べる。次に、外付け I/O エンジン方式の概要について述べる。最後にその実現方法の検討結果について説明する。

2.1 市販のストリームサーバのモジュール構成

市販ストリームサーバの一般的なモジュール構成を図1に示す。市販ストリームサーバは、配信要求処理モジュール、管理要求処理モジュール、配信制御モジュール、配信ドライバモジュールからなる。

配信要求処理モジュールは、クライアントから配信開始/停止要求を受信する。配信開始要求受信時には、管理要求処理モジュールに対して認証/課金要求の発行を行う。また、セッション情報（現在ストリーム配信を行っているクライアントとの間で確立されているコネクションに関する情報）の更新後、当該クライアントに対する配信開始/停止要求を配信制御モジュールに対して発行する。

管理要求処理モジュールは、配信要求処理モジュールからの認証/課金要求を受信すると、必要に応じ外部データベースと連動しながら当該要求に応じた処理を実行する。また、外部の管理サーバからモニタ要求を受信すると、配信ドライバモジュールが設定するモ

QuickTime や MPEG4 フォーマットのストリーム配信が可能なストリームサーバ。なお、QuickTime は米国 Apple Computer 社の登録商標である。

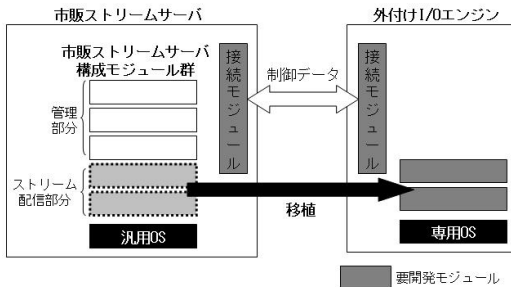


図2 外付け I/O エンジン方式を用いたストリームサーバの構成
Fig. 2 External I/O engine architecture.

ニタ情報 (CPU 負荷情報, 配信レート情報等) を参照しつつ, 当該情報を外部の管理サーバに返送する。

配信制御モジュールは周期的に駆動する (通常 100 ミリ秒程度の周期で駆動する) モジュールで, セッション情報に記載されているクライアントに対して, 自モジュール駆動周期分の配信要求を配信ドライバモジュールに対して発行する。ストリームデータの最後まで配信が完了した等の理由で, 配信すべきストリームデータが存在しない状態が一定時間以上継続している場合, セッション情報を更新する。配信要求処理モジュールは定期的に (通常 30 秒に 1 度程度) セッション情報をチェックし, 上記更新処理を検知したら, 当該クライアントに対する配信処理の強制終了処理, 具体的には, クライアントとの間に形成されている接続の強制切断処理等を実行する。

配信ドライバモジュールは, 配信制御モジュールから配信制御モジュール駆動周期分の配信要求を受信すると, 当該周期分のストリームデータをディスクから読み込み, ネットワーク経由でクライアントに対して読み込んだデータを配信する。ストリームデータのフォーマット (QuickTime, MPEG4 等) を解釈し, 当該周期に相当するデータのバイト数を決定する処理も本モジュールが行う。さらに, 必要に応じてモニタ情報を更新する。

2.2 外付け I/O エンジン方式の概要

外付け I/O エンジン方式は, 市販ストリームサーバの配信性能向上と配信品質保証機能の追加を, 機能互換性を維持しつつも, 少ない開発工数で実現するストリームサーバの構成方式である。

外付け I/O エンジン方式を用いたストリームサーバの構成を図 2 に示す。本方式では, 市販ストリームサーバと機能互換なストリームサーバを, 汎用 OS を搭載した市販ストリームサーバと, HiTactix のような専用 OS を搭載した外付け I/O エンジンとを連動させることにより実現する。

本方式では, 市販ストリームサーバの構成モジュールを, 管理部分とストリーム配信部分に分割し, 専用 OS 上にはストリーム配信部分のみを移植する。さらに, 両サーバ間で制御データを送受するために, 接続モジュールを両サーバ上に新規実装する。

本方式では, 既存の市販ストリームサーバで行っていたストリーム配信処理を外付け I/O エンジンに代理実行させる。外付け I/O エンジンは, 専用 OS を利用して, 高性能かつ配信品質を保証したストリーム配信を実現する。

また, 本方式実現のために必要となる開発は, 接続モジュールの新規実装と, ストリーム配信部分の専用 OS への移植のみに抑えられ, 専用 OS 上に市販ストリームサーバの全構成モジュールを移植する従来方式と比して, 開発工数の低減を期待できる。

2.3 実現方法の検討

HiTactix を用いて外付け I/O エンジンを実現する方法は, 2.1 節で述べた構成モジュールのうちどのモジュールを HiTactix に移植するかにより, 図 3 に示す 3 方法が考えられる。

方法 1 は, 配信要求処理モジュール, 配信制御モジュール, 配信ドライバモジュールを HiTactix に移植する方法である。この方法では, 外付け I/O エンジンがクライアントからの配信開始/停止要求を受信し, かつ, 当該要求に応じたストリーム配信を実行する。認証/課金処理や管理サーバからのモニタ要求の処理のみを市販ストリームサーバが行う。モニタ情報は外付け I/O エンジン側で更新するため, 外付け I/O エンジンが保持するモニタ情報の内容を定期的 (通常 1 秒に 1 度程度) に市販ストリームサーバ側に反映させ, 同期をとる必要がある。

方法 2 は, 配信制御モジュール, 配信ドライバモジュールを HiTactix に移植する方法である。この方法では, 市販ストリームサーバがクライアントからの配信開始/停止要求を受信する。当該要求を受信すると, 市販ストリームサーバは外付け I/O エンジンに要求を転送する。外付け I/O エンジンは, 要求に応じたストリーム配信を代理実行する。この方法では, モニタ情報以外にセッション情報も, 市販ストリームサーバ, 外付け I/O エンジンの双方で保持する。そのため, 上記モニタ情報の同期のほかに, セッション情報の同期が必要になる。セッション情報の同期のため, 定期的に (30 秒に 1 度程度), 各セッションにおけるストリーム配信が正常に行われているか否かを, 市販ストリームサーバから外付け I/O エンジンに問い合わせる。そして, 外付け I/O エンジン側で配信完了等

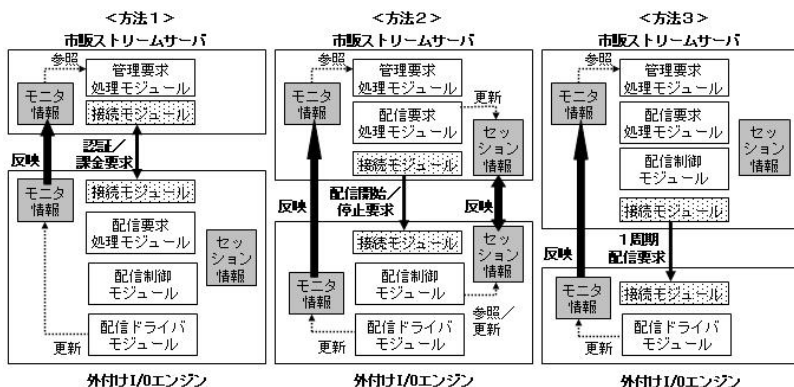


図 3 外付け I/O エンジンの実現方法
Fig. 3 External I/O engine implementation schemes.

の理由でセッション情報の更新を行ってれば、市販ストリーム側にもその更新を反映させる。

方法 3 は、配信ドライバモジュールのみを HiTactix に移植する方法である。この方法では、クライアントからの配信開始/停止要求の受信のほか、配信タイミングの制御も市販ストリームサーバが行う。外付け I/O エンジン、このタイミング制御に応じてストリームデータの I/O 処理を実行する。モニタ情報のみ上記方法で同期をとる必要がある。

上記 3 方法の優劣を、以下の観点から比較した。

- 連動オーバーヘッド (ストリーム配信性能劣化の可能性)
- 配信品質保証の容易性
- 既存システムへの影響
- 開発工数

比較結果を表 1 にまとめる。

まず、連動オーバーヘッドによるストリーム配信性能劣化の可能性についての比較結果を示す。連動オーバーヘッドは、連動にともなう市販ストリームサーバと外付け I/O エンジンとの間で発生しうる通信回数により比較した。発生しうる通信としては、

- モニタ情報同期のための通信
- セッション情報同期のための通信
- モジュール間での要求受け渡しのための通信

が考えられる。

先述したとおり、モニタ情報同期のための通信はどの方法でも発生し、その頻度は 1 秒につき 1 回程度である。セッション情報同期のための通信は方法 2 でのみ発生し、その頻度は、たかだか 30 秒につき N 回程度 (ただし N は、外付け I/O エンジンの同時配信ストリーム数を表す) である。

モジュール間での要求受け渡しのための通信は、方法 1 では、配信要求処理モジュールと管理要求処理モ

ジュール間における認証/課金処理要求にともない発生する。認証/課金処理要求はクライアントからの配信開始要求の到達にともない発生する。配信開始要求の到達頻度は、たかだかストリームデータの長さにつき N 回である。この比較では、ストリームデータの平均の長さを 5 分 (300 秒) とし、たかだか 300 秒につき N 回の頻度で発生すると仮定した。方法 2 では、配信要求処理モジュールと配信制御モジュール間における配信開始/停止要求にともない発生する。4 章で後述するように、Darwin ストリームサーバでは、配信開始/停止要求にともない 6 回の制御メッセージの通信を行う。すなわち、配信要求処理モジュールと配信制御モジュール間における配信開始/停止要求の通信頻度はたかだか 300 秒につき $6N$ 回である。方法 3 では、配信制御モジュールと配信ドライバモジュールとの間の 1 周期分の配信要求にともない発生する。この頻度は、たかだか 100 ミリ秒につき N 回程度である。

以上より、1 秒あたりに発生しうる最大の通信回数は、たかだか表 1 に記載した回数程度になると予測できる。 N が 1000 程度であっても、方法 1、方法 2 では、通信回数は 1 秒あたり 5 ~ 55 回程度で抑えられ、このオーバーヘッドはストリームデータの配信にともなう通信 (Ethernet 経由で 1.5 KB のデータ送信を 1 Gbps 程度の速度を行っている) と比べると無視しうる。それに対して方式 3 では、通信回数が 10,000 回を超え、ストリームデータの配信にともなう通信の 10% 以上を占めうる。

次に、配信品質保証の容易性の比較結果を示す。一般に、QuickTime や MPEG4 等のストリームデータには配信時刻情報 (タイムスタンプ) が格納されており、ストリームサーバは、この情報に厳密に従ったタイミングで配信を行う。方法 1、方法 2 では、配信ド

表 1 外付け I/O エンジンの実現方法の比較
Table 1 External I/O engine schemes comparison.

	方法 1	方法 2	方法 3
運動オーバーヘッド (通信回数/秒)	$(1 + N/300)$	$(1 + N/30 + N/50)$	\times $(1 + N/30 + 10N)$
配信品質保証	(HiTactix による配信制御)	(HiTactix による配信制御)	\times (汎用 OS による配信制御)
既存システム への影響	\times (サーバ IP 変更要)	(サーバ IP 変更不要)	(サーバ IP 変更不要)
開発工数	\times		

ライバモジュールの駆動を HiTactix 上に実装された配信制御モジュールが実行する。それに対して、方法 3 ではこの制御を汎用 OS 上に実装された配信制御モジュールが行う。次章で述べるとおり、HiTactix には高精度 (8 ミリ秒) で配信タイミングの制御を行う GDSL (Generic Data Streaming Library) と呼ぶライブラリが存在する。一方、汎用 OS には通常このような高精度な配信タイミングの制御機能がない。すなわち、配信品質保証に関しては、方法 1、方法 2 が方法 3 より優れているといえる。

次に、既存システムへの影響についての比較結果を示す。ここでいう「既存システムへの影響」とは、市販ストリームサーバを用いて構築された配信システムに、外付け I/O エンジンを追加し、ストリーム配信性能等の増強を図る際に必要となるシステム変更の大きさを意味する。方法 1 では、クライアントからの配信開始/停止要求の送信先となるサーバと、管理サーバからのモニタ要求の送信先となるサーバが異なる。そのため、クライアントまたは管理サーバで設定されているストリームサーバの IP アドレスの変更が不可欠となる。また、外付け I/O エンジンを複数台導入した場合、クライアントにシングルサーバとして見せるためには、負荷分散装置等の導入が不可欠となる。一方、方法 2、方法 3 ではこれらの必要はない。

最後に、外付け I/O エンジンを実現するために必要となる開発工数についての比較結果を示す。接続モジュールの開発工数はどの方法でも大きな差はないため、HiTactix に移植するモジュール数が少ない方法 3、方法 2、方法 1 の順に開発工数は小さくなる。

以上の考察から、著者らは、方法 2 を用いて HiTactix 搭載の外付け I/O エンジンを実現することにした。開発工数についてのみ方法 3 より劣るが、HiTactix において、配信制御モジュール、配信ドライバモジュールを少ない開発工数で実装することを支援する GDSL というライブラリを提供することにより、その課題を一部解決した。GDSL の詳細については次章で述べる。

3. HiTactix の概要

本章では、HiTactix の概要について示す。HiTactix は、外付け I/O エンジンにおける高性能かつ配信品質を保証したストリーム配信を実現するために必要となる基本機構を保持する。この基本機構を応用し、2 章で述べた「配信制御モジュール」「配信ドライバモジュール」を HiTactix 上に少ない開発工数で実現可能とするために、著者らは、GDSL と呼ぶライブラリを新規に HiTactix 内部に実装した。以下、HiTactix が保持する基本機構と GDSL について説明し、HiTactix を用いることで外付け I/O エンジンの開発が容易になることを明らかにする。

3.1 基本機構

HiTactix は、高性能かつ配信品質を保証したストリーム配信を実現するために必要不可欠となる以下の基本機能を持つ。

高速 I/O 機構 HiTactix は、他の高速 I/O を指向する専用 OS と同様、ゼロコピーでネットワーク I/O およびディスク I/O を実現する機構を持つ。さらに、アプリケーションが複数の非同期 I/O 要求を 1 つのシステムコールで一括して発行可能とする機能も持ち、I/O 完了通知オーバーヘッドも低く抑えられる⁸⁾。

アイソクロナススケジューラ HiTactix は、アイソクロナススケジューラと呼ぶタイマ駆動型のスケジューラを持つ。アイソクロナススケジューラは、図 4 に示すように、周期駆動が必要なスレッドの CPU 割当て時間をあらかじめテーブルを用いて予約する。予約した時間にわたり確実に当該スレッドが CPU を占有できるように、HiTactix はカーネル内部の排他制御や割り込み処理の管理を行っている⁴⁾。

サイクリックディスクスケジューリング HiTactix は、サイクリックディスクスケジューリングと呼ぶディスクスケジューラを持ち、指定したストリームのディスク I/O レートを保証できる。この I/O

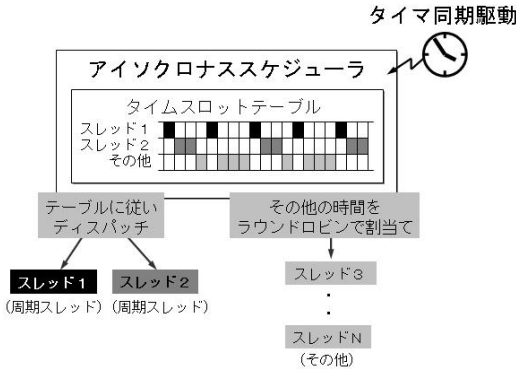


図4 アイソクロナススケジューラ Fig. 4 Isochronous scheduler.

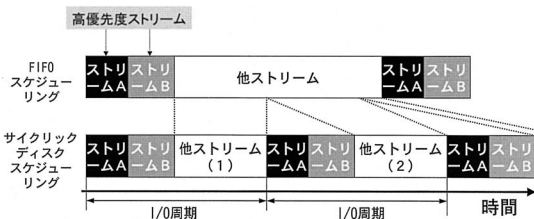


図5 サイクリックディスクスケジューラ Fig. 5 Cyclic disk scheduler.

レート保証は、図5に示すように、ディスク I/O レートを保証すべき高優先度ストリームのディスク I/O の実行時間をあらかじめ予約し、他の低優先度ストリームのディスク I/O による実行遅延が大きく発生しないことを保証する。このスケジューリングを応用すると、たとえばディスクが高負荷になった場合にも、低優先度ストリームのディスク I/O レートを落とすことで、高優先度ストリームのディスク I/O レートを保証し続けることが可能になる⁶⁾。

3.2 GDSL

GDSL は、2 章で述べた配信制御モジュール、配信ドライバモジュールを少ない開発工数で実現することを目的に HiTactix 内部に新規実装したライブラリである。以下、本ライブラリの概要と実装の詳細について説明する。

3.2.1 GDSL の概要

GDSL は、デバイス間のストリームデータ転送処理を実行する際に必要となる I/O スケジューリング、I/O 要求発行処理、バッファリング処理をアプリケーションに代わって実行するライブラリである。アプリケーションは I/O の実行内容とその実行タイミングを GDSL に対して与える処理のみを行う。GDSL は自動的に指定されたストリームデータの転送処理を、指

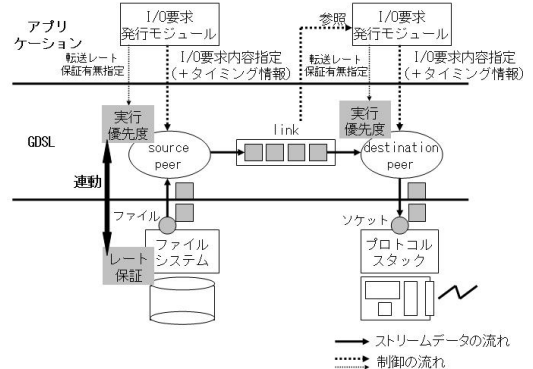


図6 GDSLを用いたストリームデータの転送処理 Fig. 6 Data streaming over GDSL.

定されたタイミングに従い、低オーバーヘッドで実行する。また、必要に応じてストリームデータの転送レート保証の有り/無しを制御することもできる。

GDSL を用いることで、高精度なタイミング制御機能を持つ「配信制御モジュール」の実現が容易になる。アプリケーションで I/O の実行タイミング指示を与えるだけで、GDSL が自動的に高精度な I/O スケジューリングを実現する。また、GDSL を用いることで、様々なストリームデータフォーマットや配信形態（ディスクからのオンデマンド配信、ネットワークから到達するストリームデータのライブ配信等）に対応した「配信ドライバモジュール」を少ない開発工数で実現できる。ストリームデータフォーマットや配信形態に非依存で必要となる処理は GDSL で、ストリームデータフォーマットや配信形態に依存する処理はアプリケーションで実現するように GDSL は設計されており、「配信ドライバモジュール」開発者がこのレイヤ設計を行う必要をなくしている。

GDSL を用いたストリームデータの転送処理の概要を図6に示す。

アプリケーションは、GDSL の提供する source peer/destination peer/link と呼ぶカーネル資源を利用しながら、ディスクやネットワーク等のデバイス間におけるストリームデータ転送を実現する。source peer/destination peer はストリームデータの読み出し先/書き込み先となるファイルまたはソケットごとに生成し、当該ファイルやソケットに対する I/O 要求を発行する。link は、peer 間でのストリームデータ受け渡しのためのバッファキュー領域を保持する。

デバイス間におけるストリームデータの転送のために必要となる I/O の内容は、アプリケーションが peer に対応させて登録する「I/O 要求発行モジュール」により制御する。I/O 要求発行モジュールは、各 peer

がファイルシステムやプロトコルスタックに発行すべき I/O 要求の内容 (I/O 先となるファイル名, ファイルオフセット, 送信先アドレス, I/O サイズ, I/O すべきデータの指定等) を指定する. destination peer に登録された「I/O 要求発行モジュール」は, この要求内容を決定する際に, link にキューイングされているストリームデータを参照することもできる.

また, peer による I/O 要求発行タイミングの制御は「I/O 要求発行モジュール」が I/O 要求内容を指定する際に, あわせて I/O 要求発行タイミング情報を与えることにより行う. また, アプリケーションが peer を生成する際に, 転送レート保証の有/無しの指定を可能にすることにより, I/O 要求発行の優先度の制御も可能にしている. 転送レート保証有りの peer による I/O の要求発行は, 転送レート保証無し of peer による I/O の要求発行より優先して実行される. また, ファイルに対応する peer に転送レート保証有りを指定した場合, その読み出し/書き込みレートもサイクリックディスクスケジューリングを用いて保証される.

GDSL を用いることにより, I/O 要求発行モジュールの開発のみで, HiTactix 上の「配信制御モジュール」「配信ドライバモジュール」を実現可能になる「配信ドライバモジュール」におけるサポートストリームデータフォーマットや配信形態の追加も, GDSL の提供機能を追加せずに実現できる.

3.2.2 GDSL の実装

GDSL では, I/O 要求発行モジュールからの I/O 要求発行内容の受取りを, 当該モジュールをアップコールすることにより実現する. また, I/O 要求発行モジュールからの I/O 要求発行タイミング情報の受取りは, 当該モジュールのアップコールからリターンの際に, 次回アップコール時刻をリターン値として受け取ることにより実現している.

GDSL では, 上記 I/O 要求発行モジュールのアップコールの実行, および当該モジュールから指示された I/O 要求の発行を, 以下の 2 つのスレッドを用いて実現している.

- アイソクロナススケジューラにより, 周期駆動と 1 周期あたりの CPU 割当て時間が保証されたリアルタイム I/O 実行スレッド
- インターバルタイマを用いて周期駆動する非リアルタイム I/O 実行スレッド

リアルタイム I/O 実行スレッドは転送レート保証有りの peer を用いた I/O 要求発行を, 非リアルタイム I/O 実行スレッドは転送レート保証無し of peer を用

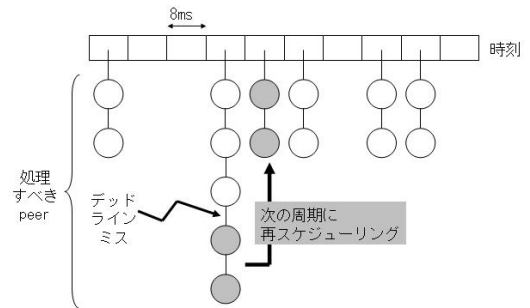


図 7 I/O 実行スケジューリングテーブル
Fig. 7 I/O scheduling table.

いた I/O 要求発行を行う. 両者は I/O 要求発行のための CPU 時間の割当て優先順位のみが異なる. 両スレッドは, とともに以下に示す処理を周期的に繰り返す.

- (1) 各 peer に登録されている I/O 要求発行モジュールを指定された時刻にアップコールする. I/O 実行スレッドは, アップコール時刻管理のため, 「I/O 実行スケジューリングテーブル」を保持している. この詳細はすぐ後で述べる.
- (2) アップコールの結果, I/O 要求発行内容の指示を受け取る. また, 次のアップコール時刻もあわせて受け取る.
- (3) ファイルまたはソケットに対して, (2) で指示された内容に基づき, I/O 要求を HiTactix の提供する非同期 I/O インタフェースを用いて一括発行する.
- (4) 前の周期で発行した I/O 完了をチェックする. source peer に対して発行した I/O が完了していたら, 当該 I/O の結果読み込んだストリームデータを link にバッファリングする.
- (5) (2) で指示されたアップコール時刻情報に基づき「I/O 実行スケジューリングテーブル」を更新する.

「I/O 実行スケジューリングテーブル」の概要を図 7 に示す. 「I/O 実行スケジューリングテーブル」は I/O 実行スレッドの駆動周期 (現在の実装では 8 ミリ秒) ごとに, 上記処理を行うべき peer が登録されている. ある周期に処理すべき peer が集中し, 1 周期内ですべての peer の処理が完了しなかった場合, 残った peer の処理は次の周期に実行する. この結果, I/O 実行要求モジュールが要求した I/O 要求発行タイミングと実際の I/O 要求発行タイミングとの間に誤差が発生しうる. この誤差の大きさがどの程度になるかについては, 5 章で定量的に評価する.

表2 接続プロトコルの制御メッセージ一覧
Table 2 Control messages.

メッセージ名	対応する RTSP メッセージ	I/O エンジンへの 要求内容
CreateSession	Describe	新規セッション情報の生成 . ストリームデータの情報 (ストリームデータフォーマット等) の返送 .
SetupSession	Setup	ストリーム配信に必要なカーネル資源 (peer, link 等) の生成 .
StartSession	Play	ストリーム配信の開始準備 (ファイルオフセット現在値等の配信制御データ初期化) .
RunSession	Play	ストリーム配信の開始 .
StopSession	Pause	ストリーム配信の停止 .
DestroySession	Teardown	ストリーム配信に使用したカーネル資源の解放 . セッション情報の削除 .

イルモジュールに相当するモジュールを新規に追加し、ファイルからストリームデータを読み込むのではなく、ネットワークからストリームデータを読み込む。

4.2 変更内容の概要

今回の Darwin ストリームサーバの変更では、通常の Darwin ストリームサーバのバージョンアップ作業と同様、基本モジュールにはほとんど手を加えていない。図9に示すとおり、上位層モジュールとして I/O エンジン駆動モジュールの新規実装を行ったのが主要な変更内容である。

I/O エンジン駆動モジュールは、接続モジュール、およびポーリングモジュールを含む。接続モジュールは、クライアントからの配信開始/停止要求の外付け I/O エンジンへの転送を行う。この転送のための接続プロトコルの詳細は次節で述べる。ポーリングモジュールはセッション情報、およびモニタ情報の同期を行うモジュールで、従来の配信ドライバモジュール(ファイル読み込み部分)を流用して実現している。

従来の配信ドライバモジュール(ファイル読み込み部分)と同様、ポーリングモジュールは配信制御モジュールから駆動される。しかし、配信ドライバモジュール(ファイル読み込み部分)とは異なり、ファイルからストリームデータを読み込まない。代わりに、セッションごとに外付け I/O エンジンが正常にストリームデータの配信を実行しているか否かのチェックを行う。このチェックは、接続モジュールに次節で述べる接続プロトコルを用いたポーリングを要求することにより実現している。正常に配信していなければ、セッション強制終了処理を実行する。配信制御モジュールに通知する次回駆動時刻を 30 秒後とすることで、ポーリングモジュールの駆動頻度を抑えている。同様の手法で、1 秒間隔でモニタ情報の同期も行っている。

外付け I/O エンジンには、接続モジュールのほか I/O エンジンモジュールを新規実装した。I/O エンジンモジュールとは、配信制御モジュール、配信ドラ

イバモジュール相当部分を GDSDL 上に実装したモジュールである。

以上述べたとおり、Darwin ストリームサーバを外付け I/O エンジン方式に対応させる際にも、Darwin ストリームサーバが管理するデータ構造や、QTSS ファイルモジュール以外のモジュールに変更を加える必要はほとんどなく、運動のために必要となる開発工数は大きくならない。より詳細な変更量の解析は 4.4 節で述べる。

4.3 接続プロトコルの概要

本節では、Darwin ストリームサーバと HiTactix 搭載の外付け I/O エンジンとの運動の際に使用する接続プロトコルの概要について述べる。

接続プロトコルは、

- 配信開始/停止要求を転送する機能
 - セッション情報/モニタ情報の同期をとる機能
- を持つ。これらの機能はすべて、Darwin ストリームサーバから外付け I/O エンジンに対して要求制御メッセージを送信し、外付け I/O エンジンが Darwin ストリームサーバに応答制御メッセージを返送することにより実現している。

配信開始/停止要求を転送する際に使用する制御メッセージの一覧を表2に示す。本制御メッセージの種類は、サーバ-クライアント間で RTSP⁹⁾プロトコルに従い送受される制御メッセージに対応して存在する。

配信要求処理モジュールは、Play 以外の RTSP 制御メッセージ(配信開始/停止要求)を受信すると、接続モジュールに対して、対応する接続プロトコルの制御メッセージの外付け I/O エンジンへの送信を要求する。接続モジュールは外付け I/O エンジンからその要求の実行結果を応答として受信する。当該実行結果は配信要求処理モジュールを介してクライアントに返送される。

配信要求処理モジュールが、Play を RTSP 制御メッセージとして受信すると、まず、StartSession の送信

表 3 接続プロトコルの制御メッセージ一覧
Table 3 Control messages.

メッセージ名	I/O エンジンへの要求内容
GetEngineStatus	モニタ情報 (CPU 負荷, 配信レート等) の返送
GetSessionStatus	セッションの配信状態 (ファイルオフセットの現在値等) の返送

表 4 改変ソースコード量
Table 4 Modified source code amount.

モジュール名	改変対象量 (K 行, 使用言語)	改変量 (K 行, 使用言語)
接続モジュール/ ポーリングモジュール (Darwin)	1.0 (C++)	3.0 (C++, 改変)
接続モジュール (HiTactix)	-	1.0 (C, 新規実装)
I/O エンジン (HiTactix)	-	5.0 (C, 新規実装)
その他 (Darwin)	-	0.1 (C++, データ構造追加)
合計	1.0	9.1

を接続モジュールに要求する。上記と同様な手順でクライアントに StartSession の実行結果を返送した後、接続モジュールに対して RunSession の送信を要求し、外付け I/O エンジンによるストリーム配信を開始させる。この制御により、クライアントへの Play 要求の実行結果の返送が、ストリーム配信開始よりも前に行われることを保証する。

Darwin ストリームサーバと外付け I/O エンジンとの間では、クライアントからの配信開始/停止要求到達にともない、6 回の制御メッセージ送受信 (外付け I/O エンジンへの要求送信と外付け I/O エンジンからの応答の受信) を行う。各制御メッセージの要求/応答には、要求/応答の種類の識別フィールド、各種要求実行時に必要となる情報 (URL, クライアント IP アドレス等)、実行結果が格納されている。その要求/応答のメッセージの大きさはそれぞれ最大で 296 バイト、272 バイトである。

セッション情報やモニタ情報の同期の際に使用する制御メッセージの一覧を表 3 に示す。GetEngineStatus は 1 秒につき 1 回、GetSessionStatus は 30 秒につき配信実行中のストリーム数だけ、Darwin ストリームサーバと外付け I/O エンジンとの間で送受される。GetEngineStatus により、Darwin ストリームサーバは外付け I/O エンジンからモニタ情報 (CPU 負荷, 配信スループット等) を取得する。現在の実装では、本制御メッセージの要求は 12 バイトで、応答は 428 バイトである。また、GetSessionStatus により、セッションごとの配信状態 (ファイルオフセットの現在値等) を取得する。本制御メッセージにより、配信が一定時間以上進行していないセッションを検知可能になる。現在の実装では、本制御メッセージの要求は 12 バイトで、応答は 44 バイトである。

本節で述べた制御メッセージの送受が、どの程度の

外付け I/O エンジンの配信性能の劣化を招くかについての定量的な評価は、5 章で述べる。

4.4 改変ソースコード量の考察

本章で示した改変で、Darwin ストリームサーバ/HiTactix 上に新規実装、もしくは改変したソースコード量を表 4 に示す。表 4 に示すとおり、今回の改変は、合計で 9.1 K 行のソースコードの新規実装、改変、データ構造の追加を行った。改変対象となった Darwin ストリームサーバのコード量は 1.0 K 行 (QTSS ファイルモジュールのみ) であり、Darwin ストリームサーバ全体 (66 K 行) の 2% 以下に抑えられた。

従来方式に従い、専用 OS 上に Darwin ストリームサーバの全構成モジュールを移植する場合、以下の開発工数が必要になると予想できる。

- (1) OS 依存モジュール (15.0 K 行) の HiTactix への移植
- (2) パスワード管理機能等の HiTactix に不足している機能の新規実装
- (3) 同期 I/O を用いて実現されているストリーム配信処理の高性能化

(3) の工数は HiTactix 上の I/O エンジンモジュールの新規実装とほぼ同等の工数が必要になる。(1) の工数を接続モジュール/ポーリングモジュールの実装と同等の工数であると仮定しても、(2) の工数の分だけ外付け I/O エンジン方式は従来方式と比して開発工数を低減していると考えられる。

また、外付け I/O エンジン方式を用いると、将来、Darwin ストリームサーバがバージョンアップした場合においても、Darwin ストリームサーバや HiTactix 上のモジュールに大きな改変を必要としない。

Darwin ストリームサーバは、他 OS への移植を容易にするため、OS 依存モジュールの提供インタフェースを定義し、このインタフェース上に他モジュールを構築している。

表5 ハードウェア仕様
Table 5 Hardware specification.

サーバ種類	サーバ名	ハードウェア仕様
汎用 OS サーバ	Darwin ストリームサーバ (外付け I/O エンジン方式非適用版)	PC/AT 互換機 (Pentium III 1.26 GHz 搭載) Gigabit Ethernet NIC (x2) (Alteon AceNIC チップ搭載) Ultra160 SCSI カード (x4) (LSILogic 53C1010 チップ搭載) Ultra160 HDD (x16)
専用 OS サーバ	簡易 Darwin ストリームサーバ	Darwin ストリームサーバ (外付け I/O エンジン方式非適用版) と同一構成
外付け I/O エンジン方式サーバ	Darwin ストリームサーバ (外付け I/O エンジン方式適用版)	PC/AT 互換機 (Pentium III 1.26 GHz 搭載) Gigabit Ethernet NIC (x1) (Alteon AceNIC チップ搭載)
	外付け I/O エンジン	Darwin ストリームサーバ (外付け I/O エンジン方式非適用版) と同一構成

たとえば、管理機構の機能向上がなされても、接続モジュールや I/O エンジンモジュールの変更は不要である。サポートするストリームデータフォーマットの追加やライブ配信サポートがなされた場合には、HiTactix 上の I/O エンジンモジュールをバージョンアップする必要がある。しかし、HiTactix には GDSL が搭載されており、このバージョンアップに要する開発工数を低減できる。バージョンアップの際に追加実装すべきコード量は、たかだか現在の I/O エンジンモジュールのコード量、すなわち C 言語で 5.0 K 行程度だと予測できる。

一方、従来方式では、Darwin ストリームサーバのバージョンアップ、特に管理機構の機能向上により、HiTactix 不足機能の新規実装がさらに必要になる可能性がある。

5. 性能評価実験

本章では、外付け I/O エンジン方式の定量的な性能評価のために行った実験内容とその結果について述べる。性能評価は、配信性能、配信タイミング保証性能、配信レート保証性能の 3 つについて行った。以下、それぞれの概要を説明する。

5.1 配信性能の評価実験

配信性能の評価実験では、外付け I/O エンジン方式を用いることで、従来方式で構築した専用 OS サーバ (専用 OS 上に市販ストリームサーバ全体を移植して実現したサーバ) と比してどの程度の配信性能の劣化が発生しうるのか、また、汎用 OS サーバ (汎用 OS 上に構築されたストリームサーバ) と比してどの程度の配信性能向上が期待できるのか、定量的に評価した。

配信性能の評価実験で用いた実験システムの構成を図 10 に示す。本実験では、汎用 OS サーバとして、

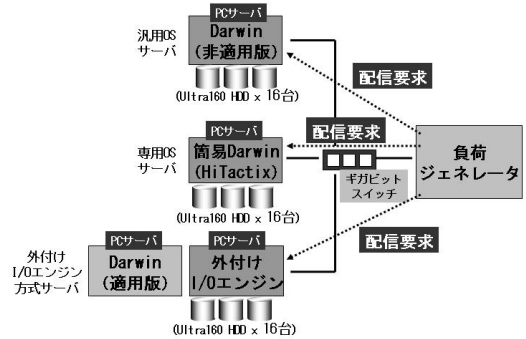


図 10 配信性能評価実験システム

Fig. 10 Experimental system for streaming performance evaluation.

い) Darwin ストリームサーバを使用した。また、専用 OS サーバとして、HiTactix を搭載した簡易 Darwin ストリームサーバを使用した。簡易 Darwin ストリームサーバとは、外付け I/O エンジンに、Darwin ストリームサーバの配信要求処理モジュールの一部を HiTactix に追加移植することで実現したサーバである。外付け I/O エンジン方式サーバとして、前章で加えた変更を加えた (外付け I/O エンジン方式に対応した) Darwin ストリームサーバと HiTactix 搭載の外付け I/O エンジンを使用した。

本実験で用いたハードウェアの仕様を表 5 に示す。本実験では、汎用 OS サーバ、専用 OS サーバは同一のハードウェア構成である。外付け I/O エンジン方式サーバは、上記ハードウェア構成に加えて、Darwin ストリームサーバ (外付け I/O エンジン方式適用版) 用 PC サーバをもう 1 台追加して用いる。

実験は、負荷ジェネレータから各サーバに対して可変のストリーム数 (各ストリームはビットレート 1.42 Mbps の QuickTime フォーマットのデータ) の配信要求を送信し、各サーバのストリーム配信レートおよび配信実行中の CPU 負荷がどのように変動するかを測定した。ただし、外付け I/O エンジン方式サー

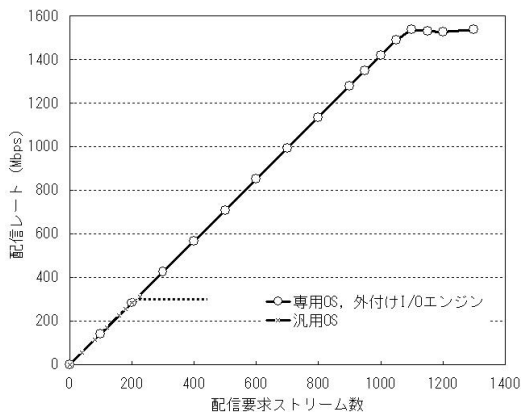


図 11 ストリーム配信レートの比較
Fig. 11 Streaming rate.

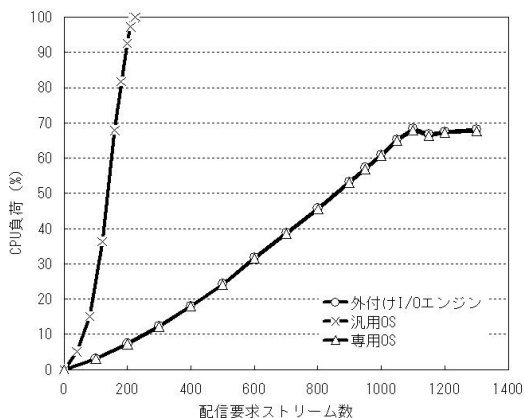


図 12 CPU 負荷の比較
Fig. 12 CPU load comparison.

は、外付け I/O エンジンの CPU 負荷を当該サーバの CPU 負荷として測定した。配信要求は、各サーバが保持する 16 台のハードディスクから均等に配信するように発行した。また、ディスクキャッシュを使用しないように、各配信要求はすべて異なるストリームデータの配信を要求した。

測定結果を図 11 および図 12 に示す。各グラフの横軸は負荷ジェネレータが配信を要求したストリーム数を、縦軸は各サーバのストリーム配信レートおよび CPU 負荷を表している。汎用 OS サーバは、配信要求ストリーム数が 220 を超えると、配信の異常終了が発生する（一定時間以上配信が実行できないストリームが生じはじめる）ため、測定ができなかった。そのため、配信要求ストリーム数が 220 ストリームを超えたときは、配信要求ストリーム数が 220 ストリームのとおり配信レートであると仮定した（図 11 の破線部分）。また、専用 OS サーバと外付け I/O エンジン方式サーバのストリーム配信レートの変動は完全に

一致するため、図 11 では両者を同一の実線で表している。

測定の結果、以下のことが明らかになった。

- 専用 OS サーバと比しても、外付け I/O エンジン方式サーバの CPU 負荷の増大は、平均で 0.49% に抑えられている。本実験では配信実行中の CPU 負荷を測定しているため、この配信性能の劣化は、2 章で述べたセッション情報の同期、モニタ情報の同期のために必要となる通信処理のみに起因している。短いストリームデータを配信した場合、配信開始/停止要求受け渡しのための通信処理による CPU 負荷増大も発生しうる。しかし、2 章で述べたとおり、5 分程度の短いストリームデータを配信した場合でも、この通信の発生回数は 1 秒あたり $N/50$ 程度 (N は同時配信ストリーム数) であり、セッション情報/モニタ情報の同期のために必要となる通信の発生回数 $1 + N/30$ と比して、同程度もしくはそれ以下である。そのため、外付け I/O エンジン方式における運動オーバーヘッドに起因する配信性能の劣化は 1% 以下に抑えられると予想できる。

- 汎用 OS サーバは、1 台のサーバを用いて、300 Mbps 程度のストリーム配信性能を達成している。一方、外付け I/O エンジン方式サーバは、2 台のサーバを用いて 1.5 Gbps 程度のストリーム配信を達成している。すなわち、外付け I/O エンジン方式の適用による 1 サーバあたりのストリーム配信性能向上の効果は 2.5 倍程度である。なお、外付け I/O エンジン方式においては、Darwin ストリームサーバ（外付け I/O エンジン適用版）にはほとんど負荷がかからないため、同時に稼働する外付け I/O エンジン方式を増やした場合、外付け I/O エンジンの台数に応じてスケラブルにストリーム配信性能が向上すると期待できる。このため、外付け I/O エンジンの同時稼働台数を増やすと、1 サーバあたりのストリーム配信性能の向上の効果は、さらに増大すると予想できる。

なお、専用 OS サーバや外付け I/O エンジン方式サーバの場合、CPU 負荷が 100% に満たないにもかかわらず、ストリーム配信レートが増えなくなる。これは、サーバの PCI バスの帯域が飽和しているためだと考えられ、さらに広帯域な PCI バスを持つサーバを使用した場合、上記の 1 サーバあたりの配信性能

実験で使用した PC サーバが保持する PCI バスの理論性能は 532 MB/秒である。1.5 Gbps のストリーム配信を行った場合、この 70% を使用している。

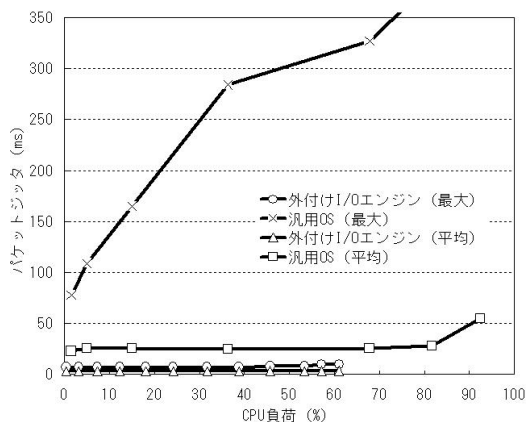


図 13 パケットジッタの比較

Fig. 13 Packet jitter comparison.

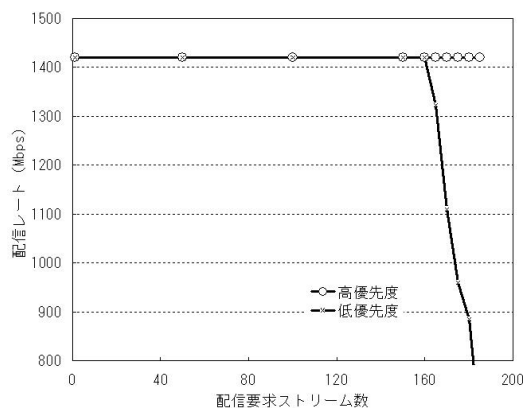


図 14 配信優先度付けの効果

Fig. 14 Prioritized streaming effects.

向上の効果はさらに増大すると期待できる。

5.2 配信タイミング保証性能の評価実験

3章で述べたとおり、HiTactix搭載の外付けI/Oエンジンでは、GDSLを用いてストリーム配信のタイミング制御を行う。GDSLは、アプリケーション(I/O実行要求発行モジュール)が指定するI/O実行要求発行タイミングがある時刻に集中した場合、その一部のI/O実行要求発行を次の周期に持ち越す可能性がある。そのため、厳密な配信タイミング保証ができない。この持ち越しによる配信タイミングの誤差がどの程度発生するか、図10に示した実験システムを用いて評価した。

本実験では、負荷ジェネレータから外付けI/Oエンジン方式サーバに可変ストリーム数の配信要求を送信し、外付けI/Oエンジンのパケットジッタがどのように変動するかを測定した。ここでいうパケットジッタとは、ストリームデータを格納する各パケットを送信すべき時刻と、実際の送信時刻との差をいう。送信すべき時刻は、QuickTimeフォーマットのストリームデータに格納されている時間情報から算出でき、この時刻に基づきアプリケーションは配信タイミングを指定する。本実験では、I/Oエンジンノードに7分間のストリームデータを配信させ、この7分間の最大と平均のパケットジッタを測定した。また、比較のため、同様の実験を汎用OSサーバに対しても行った。

測定結果を図13に示す。グラフの横軸は、外付けI/Oエンジン、またはDarwinストリームサーバ(外付けI/Oエンジン非連動版)のCPU負荷(配信要求ストリーム数を可変にした結果変動する)を示す。縦軸は、パケットジッタの大きさを示す。

測定の結果、以下のことが明らかになった。

GDSLを用いている外付けI/Oエンジン方式サーバは、少なくとも、CPU負荷が60%以下であるなら

ば、パケットジッタは平均で3.5ミリ秒、最大でも7ミリ秒~13ミリ秒で抑えられる。現在のGDSLの実装では、8ミリ秒を1周期として管理しているため、2周期より大きいI/O要求発行の持ち越しは発生していない。この最大パケットジッタの大きさは、汎用OSサーバにおける平均のパケットジッタ(23ミリ秒~55ミリ秒)の1/2以下、最大パケットジッタ(78ミリ秒以上)の1/11以下である。

5.3 配信レート保証性能の評価実験

3章で述べたとおり、GDSLは、ストリームごとにデータの転送レート保証の有り/無しを制御できる。その結果、ディスクが過負荷状態になった場合においても、優先度の低い(転送レート保証のない)ストリームのディスク読み込みレートを落とすことにより、優先度の高い(転送レート保証のある)ストリームのデータを期待されるレートでディスクから読み込み、クライアントに配信できる。

外付けI/Oエンジン方式を適用したDarwinストリームサーバが、上記に示すストリーム配信制御を行うことを確認するため、以下に示す実験を図10の実験システムを用いて行った。

まず負荷ジェネレータから、1ストリームの配信要求を外付けI/Oエンジン方式サーバに対して送る。この1ストリームは、低い優先度での配信を要求する。次に、可変ストリーム数の配信要求を負荷ジェネレータから外付けI/Oエンジン方式サーバに対して送る。これらのストリームは、高い優先度で配信するように要求する。なおこれらのストリームは、優先度の低いストリームと同じハードディスクからデータを読み出すことを要求した。優先度の高い配信を要求するストリーム数を変動させた場合、各優先度のストリームの配信レートがどのように変動するかを測定した。

測定結果を図 14 に示す。図 14 のグラフの横軸は優先度の高い配信を要求したストリーム数を、縦軸はストリーム配信レートの平均を表す。ストリーム配信レートは、優先度の高いストリームと優先度の低いストリームのそれぞれの平均値を測定した。

測定の結果、配信要求ストリーム数が 160 を超えてから期待どおりに優先度の低いストリームの配信レートが落ち込むこと、および優先度の高いストリームの配信レートは要求配信ストリーム数が 160 を超えても厳密に一定に保てる事が確認できた。

6. 関連研究

本稿で提案した外付け I/O エンジン方式のように、市販ストリームサーバとの機能互換性維持、少ない開発工数、高性能ストリーム配信、配信品質保証のすべてを同時に実現するストリームサーバを構築しようとする研究は、著者らの知る限りでは存在しない。

従来のストリームサーバは、

- 汎用 OS 上にストリームサーバを構築する。ストリーム配信性能を向上させる必要がある場合には、汎用 OS とインタフェースの互換性を維持しながらも I/O 性能の向上をはかる機能を使用する。
- 専用 OS 上にストリームサーバを構築する。専用 OS 内部に、高性能ストリーム配信機能や配信品質保証機能を組み込む。

のどちらかのアプローチで構築されていた。

前者のアプローチとしては、たとえば Linux 等の汎用 OS と、Zero-Copy TCP¹⁰⁾ や I/O-Lite¹¹⁾ 等の配信性能の向上機能を組み合わせることが考えられる。しかし、このアプローチでは配信品質保証の実現が難しい。なぜならば、配信品質の保証のためには、配信に必要な資源の予約宣言等の既存 OS のインタフェースでは提供されていない機能をアプリケーションが利用する必要があるためである。

一方、後者のアプローチとしては、splice¹²⁾、RoadRunner¹³⁾、Nemesis¹⁴⁾、Tiger Video File Server¹⁵⁾ 等が知られている。しかし、これらのアプローチは、開発工数の低減や市販ストリームサーバとの機能互換性維持について考慮していない。さらに、以下の点で本研究と差がある。

splice、RoadRunner は、デバイス間のデータ転送を高速に行う特別なパスを作成することにより、高性能なストリーム配信を実現する。しかし、このパスを経由してストリーム配信を行う場合、アプリケーションレベルのヘッダ（たとえば、RTP プロトコルを使用する場合は RTP ヘッダ）を挿入したり、ディスク

から読み込んだデータをネットワークに送信するタイミングを制御したりすることができない。HiTactix 搭載の外付け I/O エンジンでは、アプリケーションが自由に上記制御を行える。

Nemesis は、デバイスドライバ層に I/O スケジューリング機能、すなわち I/O 要求の実行順序を入れ換える機能を組み込むことにより、配信品質保証を実現する。しかし、I/O 要求発行のタイミングの制御については考慮されていないため、I/O 要求を発行するスレッドが厳密にスケジューリングされないことによる配信品質の劣化が発生しうる。HiTactix 搭載の外付け I/O エンジンは、「I/O 実行スケジューリングテーブル」を保持することにより、この配信品質の劣化を防いでいる。

Tiger Video File Server は、ディスク I/O やネットワーク I/O のスケジューリングテーブルを保持することにより、配信品質保証を実現する。しかし、ストリーム配信以外の負荷（たとえば、ストリームデータ登録にともなうディスク負荷）がかかった場合を想定しておらず、この場合に配信品質の劣化が発生しうる。HiTactix 搭載の外付け I/O エンジンは、3 章で述べたとおり、このような場合にも配信品質を保証する。

7. まとめ

本稿では、既存の市販ストリームサーバの機能互換性を維持しつつ、当該サーバの配信性能の向上、配信品質保証機能追加を少ない開発工数で実現するストリームサーバの構成方式である外付け I/O エンジン方式を提案した。本方式では、既存の市販ストリームサーバと HiTactix のような専用 OS を搭載した外付け I/O エンジンを連動させ、既存の市販ストリームサーバのストリーム配信処理のみを外付け I/O エンジンに代行させる。

さらに本稿では、Darwin ストリームサーバと HiTactix 搭載の外付け I/O エンジンを連動させた。この連動のために、HiTactix に、GDSDL と呼ぶ低オーバーヘッドでストリームデータの転送処理を行うライブラリを実装した。加えて、外付け I/O エンジン方式の定量的な評価を行った。評価の結果、上記連動は 9 K 行程度のコード追加で実現できること、また、連動オーバーヘッドによるストリーム配信性能の劣化は 1% 以下に抑えられること、等を確認した。

参考文献

- 1) 阿蘇和人：ブロードバンドに命を吹き込み、日経コミュニケーション，No.343, pp.94-111 (2001).

- 2) nCube: n4 Streaming Media System, *nCube White Paper*, No.PN109786 (2000).
- 3) Venkatasubramanian, N. and Nahrstedt, K.: An Integrated Metric for Video QoS, *Proc. 5th ACM international conference on Multimedia*, No.208, pp.371-380 (1997).
- 4) 竹内 理ほか: 連続メディア処理向き OS の周期駆動保証機構の設計と実装, *情報処理学会論文誌*, Vol.40, No.3, pp.1204-1215 (1998).
- 5) 中野隆裕ほか: Ethernet 上で QoS を保証する通信方法の設計と実装, *情報処理学会論文誌*, Vol.40, No.2, pp.322-332 (2000).
- 6) 竹内 理ほか: OS 接続モジュール Symbiose を用いた BSD-HiTactix 運動システムの設計と実装, *情報処理学会研究報告*, Vol.2000-OS-83, pp.31-36 (2000).
- 7) Iwasaki, M., Takeuchi, T., Nakano, T. and Nakahara, M.: Isochronous Scheduling and its Application to Traffic Control, *19th IEEE Real-Time System Symposium*, pp.14-25 (1998).
- 8) 岩崎正明ほか: 連続メディア処理向きマイクロカーネルの開発(1)~(5), 第53回全国大会講演論文集(1), pp.141-150 (1996).
- 9) Schulzrinne, H.: Real Time Streaming Protocol (RTSP), RFC-2326 (1998).
- 10) Chu, H.K.J.: Zero-Copy TCP in Solaris, *Proc. USENIX Annual Technical Conference*, pp.253-264 (1996).
- 11) Pai, V.S., Druschel, P. and Zwanepoel, W.: I/O-Lite: A Unified I/O Buffering and Caching System, *ACM Trans. Comput. Syst.*, Vol.18, No.1, pp.37-66 (2000).
- 12) Fall, K. and Pasquale, J.: Improving Continuous Media Playback Performance with In-Kernel Data Paths, *International Conference on Multimedia Computing and Systems*, pp.100-109 (1994).
- 13) Miller, F.W. and Tripathi, S.K.: An Integrated Input/Output System for Kernel Data Streaming, *SPIE/ACM Multimedia Computing and Networking*, pp.57-68 (1998).
- 14) Barham, P.R.: A Fresh Approach to File System Quality of Service, *7th International Workshop on Network and Operating System Support for Digital Audio and Video*, pp.119-128 (1997).
- 15) Bolosky, W.J., Fitzgerald, R.P. and Douceur, J.R.: Distributed Schedule Management in the Tiger Video Fileserver, *16th ACM Symposium*

on Operating Systems Principles, pp.212-223 (1997).

(平成 14 年 9 月 30 日受付)

(平成 15 年 5 月 6 日採録)



竹内 理 (正会員)

1969 年生。1992 年東京大学理学部情報科学科卒業。1994 年同大学院理学系研究科情報科学専攻修士課程修了。同年(株)日立製作所システム開発研究所入社。連続メディア処理向きマイクロカーネルの研究, 特にリアルタイムスケジューリング方式, リアルタイム通信方式, 異種 OS 共存技術, ストリーミングサービスアーキテクチャの研究に従事。



レ・モアル ダミエン

1972 年生。1995 年 ENSEEIHT (National Engineering School of Electrotechnics, Electronics, Informatics and Hydraulics of Toulouse, France) 卒業。2000 年京都大学大学院情報学研究所通信情報システム専攻修士課程修了。同年(株)日立製作所システム開発研究所入社。連続メディア処理向きマイクロカーネルの研究, 特にリアルタイムスケジューリング方式, サービス品質保証可能なシステムアーキテクチャ, ストリーミングサーバ向けミドルウェアの研究に従事。



坂東 忠秋 (正会員)

1968 年東京大学工学部卒業, 同年(株)日立製作所入社。日立研究所部長, 中央研究所部長, システム開発研究所副所長を経て, 現在システム開発研究所主管研究長。制御用計算機システム, 画像処理, マイクロプロセッサ, 並列型スパコン, マルチメディア, ブロードバンドコンテンツ配信等の研究開発に従事。1976 年スタンフォード大学 Electrical Engineering 修了。1984 年工学博士。1991 年~1993 年千葉大学非常勤講師。2002 年より電通大学客員教授。