

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4426589号  
(P4426589)

(45) 発行日 平成22年3月3日(2010.3.3)

(24) 登録日 平成21年12月18日(2009.12.18)

(51) Int. Cl.	F I		
HO4N 7/173 (2006.01)	HO4N	7/173	610A
GO6F 3/06 (2006.01)	GO6F	3/06	540
HO4L 12/40 (2006.01)	HO4L	12/40	Z

請求項の数 24 (全 16 頁)

(21) 出願番号	特願2006-542714 (P2006-542714)	(73) 特許権者	504199079
(86) (22) 出願日	平成16年12月2日(2004.12.2)		インタラクティブ コンテンツ エンジン
(65) 公表番号	特表2007-513582 (P2007-513582A)		ズ、エルエルシー
(43) 公表日	平成19年5月24日(2007.5.24)		アメリカ合衆国 ハワイ州 96813
(86) 国際出願番号	PCT/US2004/040235		ホノルル ビショップ・ストリート 10
(87) 国際公開番号	W02005/057828		88 4100号
(87) 国際公開日	平成17年6月23日(2005.6.23)	(74) 代理人	100070150
審査請求日	平成18年6月1日(2006.6.1)		弁理士 伊東 忠彦
(31) 優先権主張番号	60/526, 437	(74) 代理人	100091214
(32) 優先日	平成15年12月2日(2003.12.2)		弁理士 大貫 進介
(33) 優先権主張国	米国 (US)	(74) 代理人	100107766
(31) 優先権主張番号	10/999, 661		弁理士 伊東 忠重
(32) 優先日	平成16年11月30日(2004.11.30)	(72) 発明者	ローズ, スティーヴン, ダブリュ
(33) 優先権主張国	米国 (US)		アメリカ合衆国 ハワイ州 96768
			ハリイメール マイカイ 866
			最終頁に続く

(54) 【発明の名称】 同期データ転送システム

(57) 【特許請求の範囲】

【請求項1】

同期データ転送システムであって：

複数のプロセッサ・ノードと、

前記複数のプロセッサ・ノードに結合された、該複数のプロセッサ・ノード間の通信を可能にするためのバックボーンネットワークスイッチと、

前記複数のプロセッサ・ノードに結合され前記複数のプロセッサ・ノードにまたがって分散されており、複数のタイトルを保存する複数の記憶装置であって、各タイトルが複数のサブチャンクに分割され、該サブチャンクが前記複数の記憶装置にまたがって分散されるような複数の記憶装置と、

前記複数のプロセッサ・ノードのうちの対応するものの上でそれぞれ実行され、ソース・プロセッサ・ノードから宛先のプロセッサ・ノードに転送されるべき各サブチャンクについてのソース・ノード識別子および宛先ノード識別子を含むメッセージを送るよう動作する複数の転送プロセスと、

複数の逐次的な送信周期のそれぞれを開始させる送信コマンドを周期的にブロードキャストし、複数のメッセージを受け取り、各プロセッサ・ノードが各送信周期の間に高々一つのサブチャンクを送信して高々一つのサブチャンクを受信することを保証するよう各送信周期に先立って前記複数のメッセージの間から選択を行い、選択されたメッセージに対応する複数の送信要求を送る、前記複数のプロセッサ・ノードのうちの少なくとも一つの上で実行される同期スイッチマネージャ・プロセス、

とを有しており、少なくとも一つのメッセージを送り、前記同期スイッチマネージャ・プロセスから対応するサブチャUNKを同定する送信要求を受信した各転送プロセスが、前記対応するサブチャUNKをブロードキャストされた送信コマンドによって開始される次の送信周期の間に送ることを特徴とするシステム。

【請求項2】

請求項1記載の同期データ転送システムであって、前記複数のメッセージのそれぞれがタイムスタンプを有し、前記同期スイッチマネージャ・プロセスが前記複数のメッセージに対してタイムスタンプ順に基づいて優先度付けをし、前記複数の送信要求をタイムスタンプ順に送ることを特徴とするシステム。

【請求項3】

請求項2記載の同期データ転送システムであって、さらに：

前記複数のプロセッサ・ノードのうちの対応するものの上でそれぞれ実行され、複数のタイムスタンプ付き読み取り要求を送るよう動作する複数のユーザプロセスを有しており、

各転送プロセスが対応するタイムスタンプ付き読み取り要求からタイムスタンプを対応するメッセージに組み込むことを特徴とするシステム。

【請求項4】

請求項3記載の同期データ転送システムであって、前記同期スイッチマネージャ・プロセスが、前記複数のメッセージを準備完了メッセージリストにタイムスタンプ順に整理し、前記複数の逐次的送信周期の各周期の直前に前記準備完了メッセージリストをタイムスタンプ順にスキャンし、タイムスタンプの優先度に基づいてメッセージを選択することを特徴とするシステム。

【請求項5】

請求項4記載の同期データ転送システムであって、前記同期スイッチマネージャ・プロセスがあるメッセージを、同定されたソース・プロセッサ・ノードがすでに次の送信周期の間にサブチャUNKを送信するために選択されているのでない場合であり、かつ同定された宛先プロセッサ・ノードがすでに前記次の送信周期の間にサブチャUNKを受信するために選択されているのでない場合に選択することを特徴とするシステム。

【請求項6】

請求項1記載の同期データ転送システムであって、さらに：

前記複数の転送プロセスのそれぞれが、受け取ったサブチャUNK読み取り要求を読み取り要求待ち行列に保存し、各サブチャUNK読み取り要求がローカルに保存されたサブチャUNKを示していることと、

前記複数の記憶装置のそれぞれが、ローカルな読み取り要求待ち行列で同定されているサブチャUNKを物理的な順序で読み取ることと、

前記複数のプロセッサ・ノードのそれぞれが、対応する記憶装置による読み取りが成功したサブチャUNKを成功読み取り待ち行列にリストすることと、

前記複数の転送プロセスのそれぞれが、対応する成功読み取り待ち行列中の各項目についてのメッセージを前記同期スイッチマネージャ・プロセスに送ること、

とを有することを特徴とするシステム。

【請求項7】

請求項6記載の同期データ転送システムであって、前記サブチャUNK読み取り要求のそれぞれがタイムスタンプ付き読み取り要求を含んでおり、前記成功読み取り待ち行列のそれぞれにおける項目がタイムスタンプ順にリストされており、各転送プロセスが対応する成功読み取り待ち行列における各項目についてのメッセージをタイムスタンプ順に送ることを特徴とするシステム。

【請求項8】

請求項6記載の同期データ転送システムであって、さらに：

前記複数の転送プロセスのそれぞれが、ある成功要求待ち行列から、対応する送信要求によって同定されたあるサブチャUNKと関連付けられている項目を除去することと、

10

20

30

40

50

前記複数のプロセッサ・ノードのうちの対応するものの上でそれぞれ実行されており、それぞれ同定されたサブチャUNKを送信コマンドに反応して宛先プロセッサ・ノードに転送するために使われるネットワークパケットを構築するよう動作する複数のネットワーク転送プロセス、  
とを有することを特徴とするシステム。

【請求項 9】

請求項 1 記載の同期データ転送システムであって、前記ネットワークスイッチが複数のポートをもつギガビットイーサネット（登録商標）スイッチであり、前記複数のプロセッサ・ノードのそれぞれが前記ネットワークスイッチの対応するポートに結合されていることを特徴とするシステム。

10

【請求項 10】

請求項 1 記載の同期データ転送システムであって、前記複数のプロセッサ・ノードのうちに、前記同期スイッチマネージャ・プロセスを実行する管理ノードが含まれることを特徴とするシステム。

【請求項 11】

請求項 1 記載の同期データ転送システムであって、前記複数のプロセッサ・ノードのうちに、前記同期スイッチマネージャ・プロセスを実行する第一の管理ノードと、ミラーされた同期スイッチマネージャ・プロセスを実行する第二の管理ノードとが含まれることを特徴とするシステム。

【請求項 12】

ネットワークスイッチ結合された複数のプロセッサ・ノードの間でデータの分散されたサブチャUNKを同期的に転送する方法であって：

20

前記プロセッサ・ノードのうちの少なくとも一つの上で実行される管理プロセスによって、複数の逐次的な送信周期のそれぞれを開始させる送信コマンドを周期的にブロードキャストし、

前記管理プロセスに対して、送るべき少なくとも一つのサブチャUNKをもつ各プロセッサによって、送られるべき各サブチャUNKについてのソース・プロセッサ・ノードおよび宛先プロセッサ・ノードを同定するメッセージを送り、

前記管理プロセスによって、各プロセッサ・ノードが次の送信周期の間に高々一つのサブチャUNKを送信し、各プロセッサ・ノードが前記次の送信周期の間に高々一つのサブチャUNKを受信することを保証するよう、前記プロセッサ・ノードから受け取られたメッセージを選択し、

30

前記管理プロセスによって、選択された対応するメッセージを送ったプロセッサ・ノードにそれぞれ送られる複数の送信要求を送り、

送信要求を受け取る各プロセッサ・ノードによって、前記受け取られた送信要求によって同定されたサブチャUNKを次の送信コマンドに反応して宛先プロセッサ・ノードに送信する、

ことを含むことを特徴とする方法。

【請求項 13】

請求項 12 記載の方法であって、さらに：

40

前記した送られるべき各サブチャUNKについてのメッセージを送ることに先立って、各メッセージにタイムスタンプを付けることを含んでおり、

前記した選択がタイムスタンプ順に基づいて優先度を付けることであり、

前記した複数の送信要求を送ることがタイムスタンプ順に送信要求を送ることであることを特徴とする方法。

【請求項 14】

請求項 13 記載の方法であって、さらに：

少なくとも一つのプロセッサ・ノードによって、複数のタイムスタンプ付き読み取り要求を送ることを含んでおり、

前記した各メッセージにタイムスタンプを付けることが、受け取られたタイムスタンプ付

50

き読み取り要求からのタイムスタンプを対応するメッセージに組み込むことを含んでいることを特徴とする方法。

【請求項 15】

請求項 14 記載の方法であって、さらに：

前記管理プロセスによって、受け取られたメッセージを準備完了メッセージリスト中でタイムスタンプ順に整理し、

前記管理プロセスによって、各送信周期の直前に、前記準備完了メッセージリストをタイムスタンプ順にスキャンする、

ことを含むことを特徴とする方法。

【請求項 16】

請求項 15 記載の方法であって、前記したスキャンすることが、あるメッセージを、同定されたソース・プロセッサ・ノードがすでに次の送信周期の間にサブチャックを送信するために選択されているのでない場合であり、かつ同定された宛先プロセッサ・ノードがすでに前記次の送信周期の間にサブチャックを受信するために選択されているのでない場合に選択することを含むことを特徴とする方法。

【請求項 17】

請求項 16 記載の方法であって、前記したスキャンすることが、準備完了メッセージリスト全体がスキャンされたとき、あるいは全部のプロセッサ・ノードがサブチャックを送信するために選択されてしまっている場合、あるいは全部のプロセッサ・ノードがサブチャックを受信するために選択されてしまっている場合に、完了することを特徴とする方法。

【請求項 18】

請求項 12 記載の方法であって、さらに：

受け取られたサブチャック読み取り要求を読み取り要求待ち行列に保存し、各サブチャック読み取り要求がローカルに保存されたサブチャックを示しており、

ローカルなディスクドライブによって、読み取り要求待ち行列で同定されているサブチャックを物理的な順序で読み取り、

読み取りが成功したサブチャックの項目を成功読み取り待ち行列にリストする、

ことを有しており、前記した送られるべき各サブチャックについてのメッセージを送ることが、前記成功読み取り待ち行列中に各項目についてのメッセージを送ることであることを特徴とする方法。

【請求項 19】

請求項 18 記載の方法であって、各サブチャック読み取り要求がタイムスタンプ付き読み取り要求であり、前記した読み取りに成功したサブチャックの項目を成功読み取り待ち行列にリストすることが項目をタイムスタンプ順にリストすることを含んでおり、前記した前記構成読み取り待ち行列中の各項目についてのメッセージを送ることがタイムスタンプ順にメッセージを送ることを含んでいることを特徴とする方法。

【請求項 20】

請求項 18 記載の方法であって、さらに：

前記成功要求待ち行列から、対応する送信要求によって同定されたあるサブチャックと関連付けられている項目を除去し、

同定されたサブチャックを送信コマンドに反応して宛先プロセッサ・ノードに転送するために使われるネットワークパケットを構築する、

ことを含むことを特徴とする方法。

【請求項 21】

請求項 12 記載の方法であって、さらに、第一の管理ノード上で前記管理プロセスを実行し、前記第一の管理ノードをミラーするミラーされた管理ノード上でミラーされた管理プロセスを実行することを含むことを特徴とする方法。

【請求項 22】

同期データ転送システムであって：

第一および第二のユーザーノードおよび管理ノードを含む複数のストレージ・プロセッサ・ノードと、

前記複数のストレージ・プロセッサ・ノードに結合されたバックボーン通信スイッチと、

前記複数のストレージ・プロセッサ・ノードにまたがって分散される複数のサブチャUNKにそれぞれ細分された複数のタイトルと、

それぞれが対応するサブチャUNKを要求するための複数のタイムスタンプ付き読み取り要求を送る、前記第一のユーザーノード上で実行されるユーザープロセスと、

ローカルに保存されているサブチャUNKを要求する受け取ったタイムスタンプ付き読み取り要求それぞれについてのソース・ノード識別子および宛先ノード識別子を含むメッセージを送る、前記第二のユーザーノード上で実行される転送プロセスと、

前記スイッチを通じて複数の逐次的な送信周期のそれぞれを開始させる送信コマンドを周期的にブロードキャストし、複数のメッセージを受け取り、各プロセッサ・ノードが各送信周期の間に高々一つのサブチャUNKを送信して高々一つのサブチャUNKを受信することを保証するよう各送信周期に先立って前記複数のメッセージの間から選択を行い、選択されたメッセージに対応する複数の送信要求を送る、前記管理ノード上で実行される管理プロセス、

とを有しており、前記転送プロセスが、前記管理プロセスから送信要求を受け取るのに反応して、対応するサブチャUNKを次のブロードキャストされる送信コマンドによって開始される次の送信周期の間に送ることを特徴とするシステム。

#### 【請求項 2 3】

請求項 2 2 記載の同期データ転送システムであって、前記管理プロセスが前記複数のメッセージからタイムスタンプ優先度に基づいて選択を行うことを特徴とするシステム。

#### 【請求項 2 4】

請求項 2 3 記載の同期データ転送システムであって、前記管理プロセスがまず最高優先度のタイムスタンプをもつメッセージを選択し、次いでその後の各メッセージを、同定されたソース・ノードがすでに次の送信周期の間にサブチャUNKを送信するために選択されているのでない場合であり、かつ同定された宛先ノードがすでに前記次の送信周期の間にサブチャUNKを受信するために選択されているのでない場合に選択することを特徴とするシステム。

#### 【発明の詳細な説明】

#### 【技術分野】

#### 【0 0 0 1】

関連出願への相互参照

本出願は、2003年12月2日に出願された米国仮出願第60/526,437号の恩恵を主張するものであり、2002年11月26日に出願された“Interactive Broadband Server System”と題する係属中の米国特許出願、通し番号第10/304,378号の部分継続出願であり、この後者自身も2001年11月28日に出願された米国仮出願第60/333,856号の恩恵を主張している。これらすべては共通の発明者を有し、本出願の出願人に譲渡されており、ここにあらゆる面において参照によって組み込まれる。

本発明の技術分野

本発明は、対話式ブロードバンドサーバシステムに、より詳細には、多重同時アイソクロナス・データストリーム (multiple simultaneous isochronous data stream) の高速配送を容易にするための同期 (synchronized) データ転送システムを利用する対話式コンテンツエンジンに関するものである。

#### 【背景技術】

#### 【0 0 0 2】

ストリーミングメディアコンテンツの記憶および配送のためのソリューションを提供することが望まれている。スケーラビリティについての初期の目標は、ストリームあたり4メガビット毎秒 (Mbps) での100ないし100万の同時の個別的なアイソクロナス・コンテ

10

20

30

40

50

ンツストリームである。ただし、異なるデータレートも考えられている。利用可能な全通信帯域幅は利用可能な最大のバックプレーンスイッチによって制約される。現在最大のスイッチはテラビット毎秒の範囲、すなわち20万同時出力ストリーム程度である。出力ストリームの数は一般にストリームあたりのビットレートに反比例する。

【発明の開示】

【発明が解決しようとする課題】

【0003】

コンテンツ記憶装置の最も単純なモデルは、単一のネットワークコネクタを有する単一のプロセッサに接続された単一のディスクドライブである。データはディスクから読み出され、メモリに入れられ、パケットとしてネットワークを通じて各ユーザーに配布される。ウェブページなどのような伝統的なデータは非同期的に配送できる。つまり、ランダムな時間的遅延をもったランダムな量のデータがある。低ボリューム、低解像度のビデオはウェブサーバーから配送できる。ビデオおよびオーディオのようなリアルタイムのメディアコンテンツはアイソクロナス伝送、すなわち配送時間を保証された伝送を必要とする。このシナリオでは、通信帯域幅の制約はディスクドライブに存在する。ディスクはアームの動きや回転による遅延と競合しなければならない。システムが任意の所与の時点においてドライブからプロセッサへの連続的なコンテンツの同時ストリームを6つしかサポートできないとすると、7番目のユーザーの要求は先の6人のユーザーの一人がコンテンツストリームを放棄するのを待たねばならない。この設計の長所は単純さである。短所は、当該設計中で唯一の機械的装置であるディスクがその速さでしかデータのアクセス・転送ができないということである。

【0004】

一つまたは複数の別のドライブを追加したり、ドライブアクセスをインターリーブしたりすることによって改良することもできる。また、重複コンテンツを各ドライブに保存して冗長性およびパフォーマンスを稼ぐこともできる。これは前進ではあるが、それでもいくつかが問題がある。ローカルな単数または複数のドライブ上には一定量のコンテンツしか置くことができない。ディスクドライブ、CPUおよびメモリはそれぞれ破綻につながりかねない単一障害点である。このシステムは、ディスクコントローラが扱うことのできるドライブ数までしか拡大できない。ユニット数がたくさんあったとしても、タイトルの配布に問題がある。実世界では、誰もが最新の映画を見たいと思う。大まかにいって、コンテンツ要求の80%はほんの20%のタイトルについてのものである。機械の通信帯域幅のすべてを一つのタイトルにつぎ込むわけにはいかない。それではその機械にしか保存されていない人気の劣るタイトルへのアクセスが遮断されてしまう。結果として、「高需要」タイトルはほとんど、あるいはすべての機械に登録されなければならなくなる。端的に言えば、ユーザーが古い映画を見たいとしたら、そのユーザーは不運なことになりうる。たとえシステムに登録されていたとしても。大きなライブラリでは、先の比はこの例で用いた80/20の割合よりずっと大きくなりうる。

【0005】

システムがデータ処理で使われる標準的な構内ネットワーク(LAN: Local Area Network)に基づいているとしたら、他の不十分な点もあるだろう。現代のイーサネット(登録商標)ベースのTCP/IPシステムは配送が保証される驚くべきものであるが、パケットの衝突や部分的に紛失したパケットの再送信によって引き起こされる時間的な代償ならびにいっさいを機能させるために必要とされる管理を伴っている。コンテンツストリームのセットが適時に利用可能である保証はない。また、各ユーザーが一つのスイッチポートを消費し、各コンテンツサーバーが一つのスイッチポートを消費するので、スイッチポート数はサーバー数の2倍でなければならず、全体的なオンライン通信帯域幅を制限することになる。

【0006】

本発明の恩恵、特徴および利点は以下の記述および付属の図面に関してよりよく理解されるであろう。

## 【発明を実施するための最良の形態】

## 【0007】

以下の記述は、特定の用途と要件のコンテキストにおいて与えられる本発明を通常の当業者が作成し、利用できるようにするために提示される。ただし、好ましい実施形態に対するさまざまな修正が当業者には明白であろうし、ここに定義される一般的な原理は他の実施形態にも適用されうる。したがって、本発明はここに示され、説明される特定の実施形態に限定されることを意図したのではなく、ここに開示される原理および新規の特徴と一致する最も広範な範囲を与えられることを意図したものである。

## 【0008】

ここに記載されるアーキテクチャはさまざまな性能の個々の構成要素を受け入れられるので、導入が最初のシステム購入時点に限定されるのを避けることができる。利便性 (commodity) 構成要素の使用によって、近年の十分に実証済みの技術、単一ソースの回避およびストリームあたりのコスト最小化を保証する。個別の構成要素の障害は耐えられる。多くの場合には、ユーザーの観点から目につく動作の変化はない。その他の場合には、短い「自己修復」サイクルがある。多くの場合には多重障害にも耐えうる。また、すべてではないまでもほとんどの場合には、システムはすぐに対処してやらなくても復旧できるので、「無人」(lights-out) 運転にとっては理想的となる。

## 【0009】

コンテンツ記憶割り当ておよび内部通信帯域幅は、最長未使用時間 (LRU: Least Recently Used) アルゴリズムによって自動的にやりくりされる。このアルゴリズムは、RAM キャッシュおよびハードディスクアレキッシュにある内容が現在の需要に対して適切であり、バックプレーンスイッチ通信帯域幅が最も効率的な仕方で使用されることを保証する。システム内の通信帯域幅が申し込み過多になることは、よしあるとしてもまれなので、パケットの伝送を破棄または遅延することは必要でない。本アーキテクチャは、各構成要素の複合的な通信帯域幅を最大限に利用する機能を提供し、よって保証が満たされることができ、当該ネットワークはプライベートでありかつ完全なコントロール下にあり、予期しないピーク需要の状況であってもどのデータ経路も過負荷にならない。どんなビットレートのストリームでも受け入れることができるが、典型的なストリームは1ないし20M bpsの範囲に留まると期待される。非同期コンテンツは利用可能な通信帯域幅に基づいて受け入れられる。アプリケーションが要求すればその目的のために通信帯域幅をリザーブしてもよい。ファイルはいかなるサイズでもよく、記憶の不効率が最小限となる。

## 【0010】

図1は、本発明の例示的な実施形態に基づいて実装された対話的コンテンツエンジン (ICE: Interactive Content Engine) 100の一部分の単純化されたブロック図である。本発明の十分にして完全な理解のために適用できない部分は明快のため示されていない。ICE 100は、適切な多重ポート (またはマルチポート) ギガビットイーサネット (登録商標) (GbE) スwitch 101を、いくつかのストレージ・プロセッサ・ノード (SPN: Storage Processor Node) 103に結合された複数のイーサネット (登録商標) ポートを有するバックプレーンファブリックとして含んでいる。各SPN 103は簡略化したサーバーであり、2つのギガビットイーサネット (登録商標) ポート、一つまたは複数のプロセッサ 107、メモリ 109 (たとえばランダムアクセスメモリ (RAM)) および適切な数 (たとえば4ないし8) のディスクドライブ 111を含んでいる。各SPN 103上の第一のGbポート 105は全二重動作 (各SPN / ポート接続における同時の送信および受信) のためにSwitch 101の対応するポートに接続されており、ICE 100内でデータを移動させるのに使う。もう一方のGbポート (図示せず) はコンテンツ出力を下流のユーザー (図示せず) に配送する。

## 【0011】

各SPN 103はそのローカルディスクドライブへの、および5つのSPNからなる各グループ内の他の4つのSPNの他のディスクドライブへの高速アクセスを有する。Switch 101は単なるSPN 103間の通信デバイスではなく、ICE 100のためのバックプレーンであ

10

20

30

40

50

る。図解のためにSPN 1 0 3は5つしか示されていないが、ICE 1 0 0は典型的にはより多くのサーバーを含んでいることは理解されている。各SPN 1 0 3はコンテンツの記憶装置、処理器および送信機の役をする。図示した構成では、各SPN 1 0 3は、市販の構成要素を使って構成され、通常の意味におけるコンピュータではない。標準的なオペレーティングシステムが考えられているが、そのような割り込み駆動型のオペレーティングシステムは無用なボトルネックを生じることもある。

#### 【 0 0 1 2 】

各タイトル（たとえばビデオ、映画または他のメディアコンテンツ）はどの単一のディスクドライブ 1 1 1にも全体としては保存されていない。そうではなく、インターリーブ・アクセスによるスピード上の恩恵を実現するため、各タイトルのデータは分割されてICE 1 0 0内のいくつかのディスクドライブにまたがって保存されている。単一のタイトルの内容は複数のSPN 1 0 3の複数のディスクドライブにまたがっている。タイトルコンテンツの短い「時間フレーム」がラウンドロビン式に各SPN 1 0 3内の各ドライブから収集される。この仕方では、物理的な負荷は拡散されてSCSIおよびIDEのドライブ数限界を免れ、ある形のフェイルセーフ動作が得られ、タイトルの大きな集合が組織・管理される。

#### 【 0 0 1 3 】

図示した特定の構成では、各コンテンツタイトルは固定サイズのばらばらなチャンク（かたまり）（典型的にはチャンク1つあたり2メガバイト（MB）程度）に分割される。各チャンクはラウンドロビン式に異なるSPN 1 0 3のセットに保存される。各チャンクは4つのサブチャンクに分割され、パリティを表す5番目のサブチャンクが生成される。各サブチャンクは異なるSPN 1 0 3のディスクドライブ上に保存される。図示して説明される構成では、約512キロバイト（KB）（ここで「K」は1024である）のサブチャンクのサイズはディスクドライブ 1 1 1のそれぞれのデータの名目的な単位に一致する。SPN 1 0 3は5つずつグループ化され、各グループすなわちSPNセットがタイトルのデータの一つのチャンクを保存する。図示したように、前記5つのSPN 1 0 3は1～4、そして「パリティ」とラベル付けされ、これらが集団的にチャンク 1 1 3を、それぞれSPN 1、2、3、4および「パリティ」に保存される5つの別個のサブチャンク 1 1 3 a、1 1 3 b、1 1 3 c、1 1 3 d、1 1 3 eとして保存する。サブチャンク 1 1 3 a～1 1 3 eは各異なるSPNのための異なるドライブ上に（たとえばSPN1/ドライブ1、SPN2/ドライブ2、SPN3/ドライブ3など）分散的に保存されて示されているが、他のいかなる可能な組み合わせで保存されてもよい（たとえばSPN1/ドライブ1、SPN2/ドライブ1、SPN3/ドライブ3など）。サブチャンク 1～4はデータを含み、サブチャンク「パリティ」はデータサブチャンクのためのパリティ情報を含む。各SPNセットのサイズは典型的には5であるが任意であり、たとえば2つのSPNから10のSPNまでといったその他いかなる好適な数であっても全く同じようにできる。2つのSPNはその記憶の50%を冗長性のために使用することになり、10個ならそれは10%である。5は記憶の効率と障害の確率との間の妥協点である。

#### 【 0 0 1 4 】

このようにしてコンテンツを分散させることによって、少なくとも二つの目標が達成される。まず、単一のタイトルを視聴できるユーザーの数が単一のSPNセットによってサービスを受けられる数に限定されず、その限界はSPNセットすべてを合わせた通信帯域幅によって決まるようになる。したがって、各コンテンツタイトルのコピーは一つしか必要とされない。その代償は、毎秒立ち上げることのできる所与のタイトルのための新規視聴者の数の制限であるが、これは冗長記憶による無駄なスペースおよび管理上のオーバーヘッドに比べればはるかに取るに足りない制約である。第二の目標とは、ICE 1 0 0の全体的な信頼性の向上である。単一ドライブの障害はパリティドライブを使ってのその内容のリアルタイムの再生成によって隠蔽される。独立したディスクの冗長なアレイ（RAID: redundant array of independent disks）と同様である。SPN 1 0 3の障害は、それが、それぞれは動作し続けるいくつかのRAIDセットのうちの一つのドライブを含んでいるという事実によって隠蔽される。障害のあったSPNに接続していたユーザーは非常にすばやく他のSPNで走っている影のプロセスによって引き継がれる。ディスクドライ



ブの、あるいはあるSPN全体の障害の場合には、障害のあった装備を修理または交換するよう運用者に通知される。紛失したサブチャックがユーザープロセスによって再構築されると、それはそれを提供したはずだったSPNに送り返され、そこでRAM内にキャッシュされる(あたかもそのローカルディスクから読まれた場合のように)。これにより、人気のあるタイトルについて同じ再構築をすることにおいて他のユーザープロセスの時間を無駄にすることが回避される。その人気サブチャックがキャッシュされ続けるのに十分なほどでありさえすれば、その後の要求はRAMから充填されるからである。

#### 【0015】

各「ユーザー」SPN103で走っているユーザープロセス(UP: user process)の目標は、自分のディスクからのサブチャックと他のユーザーSPNからの対応する4つのサブチャックとを収集して、配送のためにビデオコンテンツのチャックを組み立てることである。ユーザーSPNは一つまたは複数の管理用の(management)MGMT SPNからは区別される。後者は同じように構成されるが、後述するように異なる機能を実行する。信頼性およびパフォーマンスを向上させるために、冗長なMGMT SPNのペアが考えられる。各UPによって実行される収集および組み立て機能は、各ユーザーSPN103上の多くのユーザーのために何度もなされる。結果として、ユーザーSPN103どうしの間にはかなりの量のデータトラフィックが行き交うことになる。そうでなければパケット衝突検出および再試行をもつ典型的なイーサネット(登録商標)プロトコルは圧倒されるだろう。典型的なプロトコルはランダムな送信のために設計されており、それらのイベントの間の不活発な時間をあてにしている。よって、このアプローチは使用されない。ICE100では、衝突は全二重の、フルスイッチ式(fully switched)アーキテクチャを使うことによって、そして通信帯域幅を注意深く管理することによって回避される。ほとんどの通信は同期的になされる。のちにさらに説明するように、スイッチ101そのものは同期的な仕方管理され、伝送の調整がされる。どのSPN103がいつ送信できるかが決定されているので、ポートが所与の期間の間にさばける以上のデータで圧倒されることはない。実際、データはまずユーザーSPN103のメモリ109において収集され、次いでその転送が同期的に管理される。調和の一環として、ユーザーSPN103どうしの中に状態信号がある。エンドユーザーに向かう実際のコンテンツとは異なり、ユーザーSPNユニットどうしの中の信号のためのデータサイズはきわめて小さい。

#### 【0016】

各サブチャックの長さ(約512Kバイト;ここで「K」は1024)は、もしそうではなくサブチャックの送信がランダムまたは非同期的に行われることが許されるとしたら、GbEスイッチ101において利用可能ないかなるバッファリングをも圧倒してしまうだろう。これだけの情報を送信するための期間は約4ミリ秒(ms)であり、複数のポートが単一のポートに同時に送信しようとしないうことを保証することが望まれる。したがって、のちにさらに説明するように、スイッチ101は、すべてのポートが最大負荷条件のもとでフルに利用されて同期的に動作するように管理される。

#### 【0017】

ファイルシステム(あるいは仮想ファイルシステム(virtual file system)すなわちVFS)を管理する冗長ディレクトリプロセスは、ユーザーによる要求があったときに所与のコンテンツタイトルがどこに保存されているかを報告することを受け持っている。それはまた、新たなタイトルをロードすべき時に必要になる記憶スペースを割り当てることも受け持っている。すべての割り当ては一体の諸チャックにおいてなされ、その各チャックは5つのサブチャックからなる。各ディスクドライブ上のスペースは当該ドライブ内では論理ブロックアドレス(LBA: Logical Block Address)によって管理される。一つのサブチャックはあるディスクドライブ上で連続的なセクタまたはLBAアドレスに保存される。ICE100における各ディスクドライブの容量は、その最大LBAアドレスをサブチャックあたりのセクタ数で割ったもので表される。

#### 【0018】

各タイトルマップまたは「ディレクトリ項目」は、そのタイトルの諸チャックがどこに

10

20

30

40

50

保存されているか、そしてより特定的には各チャンクの各サブチャンクがどこに位置しているかを示すリストを含んでいる。図示した実施例では、リスト中で一つのサブチャンクを表す各項目には、特定のユーザーSPN 1 0 3を同定するSPNID、同定されたユーザーSPN 1 0 3の特定のディスクドライブ1 1 1を同定するディスクドライブ番号(DD#)およびサブチャンクポインタ(または論理ブロックアドレスすなわちLBA)が64ビット値としてパックされて含まれている。各ディレクトリ項目は、公称4Mbpsでの約30分のコンテンツのためのサブチャンクリストを含む。これは450チャンク、すなわち2250サブチャンクに等しい。各ディレクトリ項目は補助データを含めて約20KBである。SPN上で実行されているUPがディレクトリ項目を要求すると、その項目全体が送られ、対応するユーザーのためにローカルに保存される。たとえSPNが1000ユーザーをサポートしていたとしても、ローカルなリストまたはディレクトリ項目のために消費されるメモリは20MBでしかない。

10

#### 【0019】

ICE 1 0 0はあるユーザーに利用可能な全タイトルのデータベースを維持している。このリストはローカルな光ディスクライブラリ、リアルタイムネットワークプログラミングおよび使用許諾および転送の手配がなされているところではリモート位置にあるタイトルを含む。このデータベースは各タイトルについての全メタデータを含む。それには管理情報(使用許諾期間、ビットレート、解像度など)とともにユーザーにとって関心のある情報(プロデューサー、監督、キャスト、スタッフ、原作者など)も含まれる。ユーザーが選択をすると、仮想ファイルシステム(VFS) 2 0 9(図2)のディレクトリに問い合わせがされて、そのタイトルがすでにディスクアレイにロードされているかどうか判別される。まだであれば、ロードプロセス(図示せず)がそのコンテンツ作品について開始され、必要ならいつ視聴のために利用可能になるかについてUPに通知される。たいていの場合、遅延は光ディスク取得ロボット(図示せず)の機械的遅延、すなわち約30秒を超えない。

20

#### 【0020】

光ディスク(図示せず)に保存されている情報は全メタデータ(ディスクが最初にライブラリにロードされたときにデータベースに読み込まれる)とともに、当該タイトルを表す圧縮されたデジタルビデオおよびオーディオならびにそれらのデータストリームについて事前に知得できる全情報を含む。たとえば、クロック値およびタイムスタンプのようなデータストリーム中のあらゆる有意な情報へのポインタを含んでいる。それはすでにサブチャンクに分割されており、パリティサブチャンクも事前計算されてディスク上に保存されている。一般に、ロード時間と処理オーバーヘッドを節約するために事前になしうことは何でも光ディスク上に含められる。

30

#### 【0021】

リソース管理システムに含まれるものに、ディスパッチャー(図示せず)がある。UPはこのディスパッチャーを参照してそのストリームについての開始時刻を受け取る(通例要求から数ミリ秒以内に)。ディスパッチャーはシステムへの負荷が均一のままであり、遅延が最小化され、いかなる時点でもICE 1 0 0内で要求される通信帯域幅が利用可能な値を超えることのないことを保証する。ユーザーが停止、一時停止、早送り、巻き戻しあるいはその他ストリームの流れを中断する動作を要求したときには常に、その通信帯域幅は割り当てを解除され、新たに要求される何らかのサービス(たとえば早送りストリーム)のために新たな割り当てが行われる。

40

#### 【0022】

図2は、本発明のある実施形態に基づいて実装された同期データ転送システム200を図解する、ICE 1 0 0の一部分の論理ブロック図である。スイッチ101は、いくつかの例示的なSPN 1 0 3に結合して示されている。SPN 1 0 3には第一のユーザーSPN 2 0 1、第二のユーザーSPN 2 0 3および管理(MGMT)SPN 2 0 5が含まれる。先に注記したように、多くのSPN 1 0 3がスイッチ101に結合されており、本発明を解説するために二つのユーザーSPN 2 0 1、2 0 3のみが示されており、先に述べたようにどのSPN 1 0 3として物理的に実装されてもよい。MGMT SPN 2 0 5は物理的実装は他のどのSPN 1 0 3とも同様

50

であるが、一般に特定のユーザー機能よりむしろ管理機能を実行する。SPN 2 0 1 は各ユーザー-SPN 1 0 3 のある種の機能を図解し、SPN 2 0 3 は他の機能を図解する。しかし、各ユーザー-SPN 1 0 3 が同様の機能を実行するよう構成され、SPN 2 0 1 について述べた機能（およびプロセス）はSPN 2 0 3 で与えることもでき、その逆も言えることは理解されるものである。

#### 【 0 0 2 3 】

先に述べたように、スイッチ 1 0 1 はポートあたり1Gbpsで動作し、よって各サブチャUNK（約512KB）があるSPNから別のSPNに渡されるのに約4msかかる。各ユーザー-SPN 1 0 3 は一つまたは複数のユーザープロセス（UP）を実行するが、そのそれぞれが下流のユーザーをサポートするためのものである。あるタイトルの新しいチャUNKがユーザー出力バッファ（図示せず）を再充填するために必要とされるとき、リストからの次の5つのサブチャUNKが、それらのサブチャUNKを保存している他のユーザー-SPNに対して要求される。多くのUPが複数サブチャUNKを実質同時に要求する可能性があるため、サブチャUNK伝送期間はほうっておいたら単一ポートのためのGbEスイッチでさえほとんどどんなものでもバッファリング容量を圧倒してしまうだろう。ましてや全体スイッチの場合はそうである。これは図示したスイッチ 1 0 1 について成り立つ。もしサブチャUNK伝送が管理されていなければ、それは可能性として各UPのための5つのサブチャUNKすべてが同時に返されて出力ポートの通信帯域幅を圧倒する結果につながる。ICE 1 0 0 のSPNの伝送のタイミングを緊密にし、最も重要なデータが最初に無傷で伝送されるようにすることが望まれる。

#### 【 0 0 2 4 】

SPN 2 0 1 は、対応する下流のユーザーにサービスを提供するためにUP 2 0 7 を実行しているものとして示されている。ユーザーがタイトル（たとえば映画）を要求し、その要求はUP 2 0 7 に転送される。UP 2 0 7 はタイトル要求（TR: title request）をMGMT SPN 2 0 5 に位置しているVFS 2 0 9（のちにさらに説明する）に送信する。VFS 2 0 9 はディレクトリ項目（DE: directory entry）をUP 2 0 7 に返し、該UP 2 0 7 は2 1 1で示されるDEをローカルに保存する。DE 2 1 1 は当該タイトルの各サブチャUNK（SC1、SC2など）の位置を示すリストを含んでおり、各項目には、特定のユーザー-SPN 1 0 3 を同定するSPNID、同定されたSPN 1 0 3 の特定のディスクドライブ 1 1 1 を同定するディスクドライブ番号（DD#）および同定されたディスクドライブ上のサブチャUNKの特定の位置を与えるアドレスすなわちLBAが含まれている。SPN 2 0 1 は、DE 2 1 1 内の各サブチャUNKについて同時に一つずつタイムスタンプ付き読み出し要求（TSRR: time stamped read request）を開始する。ICE 1 0 0 では、要求はすぐに直接なされる。換言すれば、SPN 2 0 1 はサブチャUNKの要求を行うことを、そのデータを保存している特定のユーザー-SPN 1 0 3 に対してすぐに直接に開始する。図示した構成では、たとえローカルに保存されていても同じようにして要求がなされる。つまり、要求されたサブチャUNKがSPN 2 0 1 のローカルディスクドライブ上にあったとしても、リモートにあるかのようにスイッチ 2 0 1 を通じて要求を送出するのである。ネットワークは、あるSPNからその同じSPNに要求が送られようとしていることを認識するよう構成されていてもよい位置である。すべての場合を同じように扱うほうが簡単である。特に、要求が実際にローカルであることが比較的ありそうもないような大きな施設ではそうである。

#### 【 0 0 2 5 】

要求はすぐに直接に送出されるが、サブチャUNKはそれぞれ完全に管理された仕方です返される。各TSRRは当該SPNIDを使っている特定のユーザー-SPNに対するものであり、対象となるユーザー-SPNがそのデータを取得して返すためのDD#およびLBAを含んでいる。TSRRはさらに、要求されたサブチャUNKが適正に適切な要求者に返されることを保証し、要求者がそのサブチャUNKを識別することを可能にするのに十分な他のいかなる識別情報をも含みうる（たとえば宛先SPN上で実行されている複数UPの間の区別をするUP識別子、各データチャUNKについてのサブチャUNKの間の区別をするサブチャUNK識別子など）。各TSRRはもともとの要求がなされた特定の時点と同定するタイムスタンプ（TS: timestamp）を

も含む。TSは同期伝送のために要求の優先度を同定する。ここで、優先度は時間に基づいており、早い要求ほど高い優先度を得る。要求されたタイトルの返されたサブチャックは、受け取られると、さらなる処理およびそのタイトルを要求したユーザーへの配送のためにローカルなタイトルメモリ 2 1 3 に保存される。

#### 【 0 0 2 6 】

ユーザーSPN 2 0 3 は、TSRRを受け取るため、そして要求されたサブチャックを返すために各ユーザーSPN (たとえば 2 0 1、2 0 3) 上で実行されている、転送プロセス (TP: transfer process) 2 1 5 および支援機能の動作を図解している。TP 2 1 5 はストレージ・プロセス (図示せず) を含むか、あるいは他の仕方ですトレージ・プロセスとのインターフェースをもつ。ストレージ・プロセスは、保存されているサブチャックの要求およびアクセスのためのSPN 2 0 3 上のローカルなディスクドライブ 1 1 1 のインターフェースとなるものである。ストレージ・プロセスは状態機械などのような、いかなる所望の仕方です実装されてもよく、TP 2 1 5 とローカルディスクドライブ 1 1 1 との間でインターフェースされる別個のプロセスであってもよい。これは当業者には既知のとおりである。図示したように、TP 2 1 5 は他のユーザーSPN 1 0 3 上で実行されている一つまたは複数のTSRRを一つまたは複数のUPから受け取り、各要求をローカルメモリ 1 0 9 内の読み取り要求待ち行列 (RRQ: read request queue) 2 1 7 に保存する。RRQ 2 1 7 はサブチャックSCA、SCBなどのための要求のリストを保存する。要求されたサブチャックを保存しているディスクドライブは、対応する要求をRRQ 2 1 7 から除去し、物理的な順番にソートし、それからそれぞれの読み出しをソートされた順番で実行する。各ディスク上のサブチャックへのアクセスはグループで管理される。各グループは「エレベーター・シーク」動作に従って物理的順序にソートされている (低位から高位へ一回の掃引、次に高位から低位への掃引などといった具合にディスクヘッドがシークス上次のサブチャックを読むために一時停止しながらディスク表面を行き来する)。うまくいった読み込みの要求は成功読み込み待ち行列 (SRQ: successful read queue) 2 1 8 にTS順にソートされて保存される。失敗した読み込みの要求 (もしあれば) は失敗読み込み待ち行列 (FRQ: failed read queue) 2 2 0 に保存され、失敗した情報がネットワーク管理システム (図示せず) に転送され、該ネットワーク管理システムがエラーおよび適切な是正動作を決定する。図示した構成では待ち行列 2 1 7、2 1 8 および 2 2 0 が実際のサブチャックではなく要求情報を保存していることを注意しておく。

#### 【 0 0 2 7 】

読み込みに成功した各サブチャックは最近要求されたサブチャックのLRUキャッシュのためにリザーブされているメモリに入れられる。取得された各サブチャックについて、TP 2 1 5 は対応するメッセージ (MSG: message) を生成するが、これはそのサブチャックのTS、そのサブチャックのソース (SRC: source) (たとえばそのサブチャックが送信されてくるもとのSPNIDおよびその物理メモリ位置ならびにその他任意の識別情報) ならびにそのサブチャックを送信すべき先の宛先 (DST: destination) SPN (たとえばSPN 2 0 1) を含む。図示したように、SRQ 2 1 8 はサブチャックSCA、SCBなどに対してそれぞれメッセージMSGA、MSGBなどを含んでいる。要求されたサブチャックが読み出されてキャッシュされたのち、TP 2 1 5 は対応するMSGを、MGMT SPN 2 0 5 で実行されている同期スイッチマネージャ (SSM: synchronized switch manager) 2 1 9 に送る。

#### 【 0 0 2 8 】

SSM 2 1 9 は諸ユーザーSPNのTPから複数MSGを受け取って優先順位付けし、最終的には、SRQ 2 1 8 内のMSGのうちの一つをメッセージ識別子 (MSGID: message identifier) などを利用して同定する送信要求 (TXR: transmit request) をTP 2 1 5 に送る。SSM 2 1 9 がSRQ 2 1 8 内のサブチャックを同定するMSGIDをもつTXRをTP 2 1 5 に送ると、その要求リスト項目はSRQ 2 1 8 からネットワーク転送プロセス (NTP: network transfer process) 2 2 1 に移され、そこで当該サブチャックを宛先のユーザーSPNに転送するのに使われるパケットが構築される (ここで、「移す」とはその要求をSRQ 2 1 8 からは除去することを表している)。サブチャック要求リスト項目がSRQ 2 1 8 から除去される順序

は、リストがタイムスタンプ順になっているにもかかわらず、必ずしもその順ではない。適正な順序を決定するのはSSM 2 1 9 だけだからである。SSM 2 1 9 は、送るべきサブチャックを少なくとも一つもつ他のすべてのSPN 1 0 3 に対して一つのTXRを送る。ただし、サブチャックを送るべき先がすでに同等以上の優先度のサブチャックを受け取るようスケジュールされているSPN 1 0 3 上のUPである場合は別であり、これについてはのちにさらに説明する。SSM 2 1 9 は次いで全ユーザーSPN 1 0 3 に対して単一の送信コマンド (TX CMD: transmit command) をブロードキャストする。SSM 2 1 9 によってブロードキャストされたTX CMDコマンドに反応して、TP 2 1 5 はNTP 2 2 1 に、そのサブチャックをユーザーSPN 1 0 3 の要求UPに送信するよう命令する。このようにして、SSM 2 1 9 からTXRを受け取った各SPN 1 0 3 は同時に別の要求元ユーザーSPN 1 0 3 に送信を行う。

10

## 【 0 0 2 9 】

MGMT SPN 2 0 5 上のVFS 2 0 9 はタイトルのリストおよびICE 1 0 0 内でのその位置を管理する。典型的なコンピュータシステムでは、ディレクトリ (データの情報) は通例当該データが存在するのと同じディスク上にある。ところがICE 1 0 0 では、VFS 2 0 9 は中央に位置して分散データを管理する。各タイトルのデータがディスクアレイ中の複数のディスクにまたがって分散しており、前記複数のディスク自身も複数のユーザーSPN 1 0 3 にまたがって分散しているからである。先に述べたように、ユーザーSPN 1 0 3 上のディスクドライブ 1 1 1 は主としてタイトルのサブチャックを保存する。VFS 2 0 9 は、先に述べたようにSPNID、DD#およびLBAを通じて各サブチャックの位置を決めるための識別子を含んでいる。VFS 2 0 9 はまた、光記憶装置のようなICE 1 0 0 の外部にある他の部分のための識別子をも含んでいる。ユーザーがあるタイトルを要求すると、そのユーザー要求を受け付けたユーザーSPN 1 0 3 上で実行されているUPにはディレクトリ情報 (ID / アドレス) の完全なセットが利用可能となる。そこからは、すべきことはサブチャックをディスクドライブからメモリ (バッファ) へと転送し、該サブチャックをスイッチ 1 0 1 を通じて要求元ユーザーSPN 1 0 3 に移すことである。要求元ユーザーSPN 1 0 3 はバッファ中で完全なチャックを組み立て、それをユーザーに配送し、終了するまで繰り返す。

20

## 【 0 0 3 0 】

SSM 2 1 9 は、「準備完了」メッセージ (RDY MSG: ready message) リスト 2 2 3 中にタイムスタンプ順の準備完了メッセージのリストを作成する。ユーザーSPN 1 0 3 上でTPからメッセージが受け取られる順番は必ずしもタイムスタンプ順ではないが、RDY MSG リスト 2 2 3 中ではTS順に整理されている。次の一組の転送の直前に、SSM 2 1 9 は最も早いタイムスタンプから始まってRDY MSG リスト 2 2 3 をスキャンする。SSM 2 1 9 はまずRDY MSG リスト 2 2 3 内で最も早いTSを同定し、対応するTXRメッセージを生成して、対応するサブチャックを保存しているユーザーSPN 1 0 3 のTP 2 1 5 に送り、それによりそのサブチャックのペンディングの転送を開始させる。SSM 2 1 9 は後続の各サブチャックについてTS順にリスト 2 2 3 のスキャンを続け、すでにペンディングのサブチャック転送に関わっているのではないソースと宛先をもつ各サブチャックについてTXRメッセージを生成する。全ユーザーSPN 1 0 3 にブロードキャストされる一つ一つのTX CMDに対しては、各ユーザーSPN 1 0 3 が送信するサブチャックは同時に一つのみであり、受信するサブチャックは同時に一つのみである。ただし、この両者は同時にできる。たとえば、SPN#2 へのペンディングのサブチャック転送をスケジュールするためにTXRメッセージがSPN#10 のTPに送られた場合、SPN#10は同時に別のサブチャックを送ることはできない。しかし、SPN#10は同時に別のSPNからサブチャックを受信することはできる。さらに、SPN#2はSPN#10から前記サブチャックを受信しながら同時に別のサブチャックを受信することはできないが、SPN#2は同時に別のSPNに送信をすることはできる。これはスイッチ 1 0 1 の各ポートの全二重性のためである。

30

40

## 【 0 0 3 1 】

SSM 2 1 9 はユーザーSPN 1 0 3 が全部ふさがってしまうまで、あるいはRDY MSG リスト 2 2 3 の終わりに達するかするまでRDY MSG リスト 2 2 3 のスキャンを続ける。RDY MSG リスト 2 2 3 中のTXRメッセージに対応する各項目は最終的にはRDY MSG リスト 2 2 3 か

50

ら除去される（TXRメッセージが送られたとき、あるいは転送が完了したあとに）。前の周期の最後の転送が終わったとき、SSM 2 1 9 は全ユーザーSPN 1 0 3 に次のラウンドの送信を開始するよう合図するTX CMDパケットをブロードキャストする。各転送は、図示した特定の構成については約4ないし5msの周期内で同期的に起こる。各転送ラウンドの間、追加的なMSGがSSM 2 1 9 に送られ、次の送信ラウンドをスケジュールするために新しいTX Rメッセージが諸ユーザーSPN 1 0 3 に送出され、そのプロセスは繰り返される。TX CMDから次のTX CMDまでの間の周期は近似的には、一つのサブチャUNKのバイトすべてを送信するために必要な、パケットオーバーヘッドおよびパケット間遅延を含めた期間に、そのサブチャUNKの送信中にスイッチ内に発生したかもしれないすべてのキャッシュをクリアするための期間（典型的には60マイクロ秒（ $\mu s$ ））と、個々のSPNによるTX CMDの認識における遅延によって引き起こされるあらゆるジッターのための期間（典型的には100  $\mu s$ 未満）とを加えたものに等しい。

10

**【 0 0 3 2 】**

ある実施形態では、重複した、すなわちミラーされたMGMT SPN（図示せず）が主MGMT SPN 2 0 5 をミラーしており、SSM 2 1 9、VFS 2 0 9 およびディスプレイパッチャーがそれぞれ一対の冗長な専用MGMT SPNの上で重複される。ある実施形態では、同期TX CMDブロードキャストがMGMT SPN 2 0 5 の健全性を示す鼓動としてはたらく。鼓動とは、副次的なMGMT SPNに対する万事良好という信号である。鼓動がないと、たとえば5msなど所定の時間期間内のいっさいの管理機能を副次的MGMT SPNが引き継ぐのである。

**【 0 0 3 3 】**

20

本発明はある種の好ましいバージョンを参照しつつかなりの詳細において記載されてきたが、他のバージョンおよび変形も可能であり、考えられている。当業者は開示された概念および具体的な実施例を、請求項によって定義される本発明の精神および範囲から外れることなく、本発明と同じ目的に資するための他の構造を考案したり修正したりする基礎として容易に使うことができる。

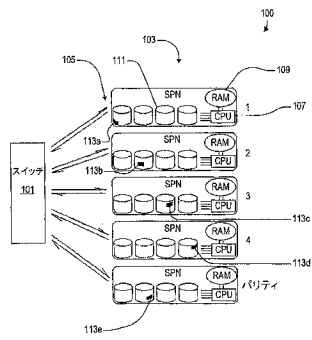
**【 図面の簡単な説明 】****【 0 0 3 4 】**

【 図 1 】 本発明の例示的な実施形態に基づいて実装された対話的コンテンツエンジン（ICE: Interactive Content Engine）の一部分の単純化したブロック図である。

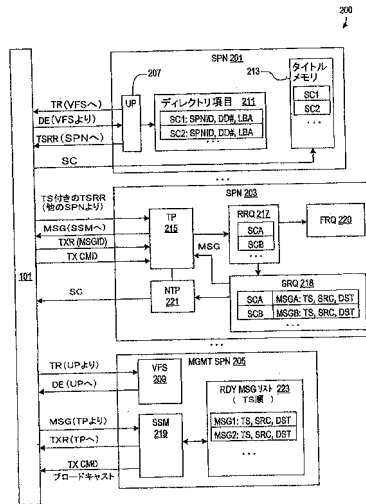
30

【 図 2 】 本発明のある実施形態に基づいて実装された同期データ転送システムを図解する、図 1 のICEの一部分の論理ブロック図である。

【図1】



【図2】



---

フロントページの続き

審査官 古川 哲也

- (56)参考文献 特開平07-141232(JP,A)  
特開平07-236132(JP,A)  
特開平08-107422(JP,A)  
特開平08-107542(JP,A)  
特開平09-097136(JP,A)  
特開平09-138735(JP,A)  
米国特許第06134596(US,A)  
国際公開第01/10125(WO,A1)

(58)調査した分野(Int.Cl., DB名)

H04N 7/16 - 7/173  
G06F 3/06  
H04L 12/40