

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4076724号
(P4076724)

(45) 発行日 平成20年4月16日 (2008. 4. 16)

(24) 登録日 平成20年2月8日 (2008. 2. 8)

(51) Int. Cl.

F I

H 0 4 L 12/42 (2006. 01)

H 0 4 L 12/42

D

請求項の数 8 (全 27 頁)

(21) 出願番号 特願2000-532850 (P2000-532850)
 (86) (22) 出願日 平成11年2月24日 (1999. 2. 24)
 (65) 公表番号 特表2002-504765 (P2002-504765A)
 (43) 公表日 平成14年2月12日 (2002. 2. 12)
 (86) 国際出願番号 PCT/US1999/003955
 (87) 国際公開番号 W01999/043002
 (87) 国際公開日 平成11年8月26日 (1999. 8. 26)
 審査請求日 平成17年3月15日 (2005. 3. 15)
 (31) 優先権主張番号 60/075, 797
 (32) 優先日 平成10年2月24日 (1998. 2. 24)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 500373758
 シーゲイト テクノロジー エルエルシー
 アメリカ合衆国, カリフォルニア, スコッ
 ツ バレイ, ビー. オー. ボックス 66
 360, ディスク ドライブ 920
 (74) 代理人 100066692
 弁理士 浅村 皓
 (74) 代理人 100072040
 弁理士 浅村 肇
 (74) 代理人 100107467
 弁理士 員見 正文
 (74) 代理人 100094673
 弁理士 林 拓三

最終頁に続く

(54) 【発明の名称】 ダイナミック半二重によるループ・フェアネスの保持

(57) 【特許請求の範囲】

【請求項 1】

ループ・フェアネスを保持する通信チャネル・システムであって、

ファイバ・チャネル・アービトレーテッド・ループ・シリアル通信チャネルに取り付けられた第1のポートを有する第1のチャネル・ノードであって、前記第1のポートが、該ポートが取り付けられた通信チャネルの制御のためにアービトレーションする、第1のチャネル・ノードと、

前記第1のポートに作動的に結合されたフェアネス保持装置と、
 を具備し、

前記第1のポートが前記通信チャネルのループの制御のためにアービトレーションし、
 制御が達成された後に、前記フェアネス保持装置が、所定量の使用が第1のポートと前記通信チャネルとの間で発生したか否かに少なくとも一部基づいて前記通信チャネルの制御を解放させ、前記所定量の使用が送信される第1の所定量のデータを有し、

送信されているデータの量を監視する第1のカウントと、

前記第1のカウントと作動的に結合し、前記第1のカウントによってモニターされるデータの量が前記第1の所定量に達したか否かに少なくとも一部基づいて通信チャネルを解放するための制御を生ずる第1のコンパレータ回路と、

送信されるべく残っているデータの量を監視する第2のカウントと、

前記第2のカウントによって監視されるデータの量が第2の所定量のデータより少ないか否かに少なくとも一部基づいて、前記通信チャネルの制御を解放することを禁止する前

10

20

記第 2 のカウンタに作動的に結合する第 2 のコンパレータ回路と、
を備える前記通信チャネル・システム。

【請求項 2】

前記第 1 のチャネル・ノードが、ダイナミック半二重をサポートし、
 前記第 1 のカウンタおよび前記第 1 のコンパレータが、ダイナミック半二重コマンドの
 第 1 の部分による受信時に初期化される、
請求項 1 記載のシステム。

【請求項 3】

前記第 1 の所定量のデータおよび前記第 2 の所定量のデータがプログラム可能な量であ
 る、請求項 2 記載のシステム。

【請求項 4】

ファイバ・チャネル・アービトレーテッド・ループ・シリアル通信チャネルと、
 前記第 1 のチャネル・ノードに作動的に結合された磁気ディスク記憶ドライブと、
 第 2 のチャネル・ノードを有するコンピュータ・システムと、
 をさらに具備し、
 前記第 2 のチャネル・ノードが、前記ファイバ・チャネル・アービトレーテッド・ル
 ープ・シリアル通信チャネルを介して前記第 1 および第 2 のチャネル・ノード間でデータ
 を転送するために、ファイバ・チャネル・ループ内の前記第 1 のチャネル・ノードに作動的
 に結合されている、
請求項 1 記載のシステム。

【請求項 5】

(a) ファイバ・チャネル・アービトレーテッド・ループ・シリアル通信チャネルのル
 ープの制御のためにアービトレーションするステップと、
 (b) 所定量の使用が第 1 のポートと前記通信チャネルとの間で発生したか否かに少な
 くとも一部基づいて前記通信チャネルの制御を解放させるステップであって、前記解放さ
 せるステップ (b) がさらに、
(b) (i) 第 1 の所定量のデータが送信されているか否かを判断するステップと、
(b) (i i) 前記判断するステップ (b) (i) に基づいてループの制御を解放するステ
 ップと、
(b) (i i i) 第 2 の所定量のデータが転送されるべく残っているか否かを判断するス
 テップと、
(b) (i v) 前記判断するステップ (b) (i i i) に基づいてループの制御を解放す
 るステップと、
 を含む、通信方法。

【請求項 6】

前記判断するステップ (b) (i i i) が、
 (b) (i i i) (A) 第 1 の値を得るために転送されたフレームの数を監視するステ
 ップと、
 (b) (i i i) (B) 前記第 1 の値を前記第 1 の所定量のデータと比較するステッ
 プと、
 をさらに含む、請求項 5 記載の方法。

【請求項 7】

(c) ダイナミック半二重コマンドを受信するステップと、
 (d) 該ダイナミック半二重コマンドの受信時に前記判断するステップ (b) (i i i)
) を初期化するステップと、
 をさらに含む、請求項 5 記載の方法。

【請求項 8】

前記初期化するステップ (d) が、
 (d) (i) 前記第 1 の所定量のデータおよび前記第 2 の所定量のデータをプログラム
 可能にセットするステップをさらに含む、

請求項 5 記載の方法。

【発明の詳細な説明】

【 0 0 0 1 】

本出願は、1998年2月24日出願の米国仮出願第60/075,797号の35U.S.C.119(e)に基づく恩典を請求する。

(発明の分野)

本発明は、大容量記憶装置の分野に関する。特に、本発明は、データ伝送のダイナミック半二重モードによりループ・フェアネスを保持する装置および方法に対する改良ファイバ・チャンネル・アービトレテッド・ループ(FC-AL)装置および方法に関する。

【 0 0 0 2 】

10

(発明の背景)

あらゆるコンピュータ・システムの重要な一構成要素は、データを記憶する装置である。コンピュータ・システムは、データを記憶できる多くの異なる装置を有する。コンピュータ・システムにおける膨大な量のデータを記憶する一つの共通の場所は、ディスク・ドライブに関する。ディスク・ドライブの最も基本的な部品は、回転するディスクと、トランスデューサをディスク上の種々の位置に移動させるアクチュエータと、ディスクにデータを書き込んだりディスクからデータを読み出すために使用される回路とである。ディスク・ドライブは、首尾よくディスク面から読み出したりディスク面に書き込んだりできるようにデータを符号化する回路も備えている。マイクロプロセッサは、ディスク・ドライブの大部分の動作を制御するとともに、要求コンピュータにデータを返送したり、要求コンピュータからデータを取り込んでディスクに記憶する。

20

【 0 0 0 3 】

ディスク・ドライブとコンピュータ・システムの他の部分との間でデータを転送するためのインターフェースは、典型的には、スモール・コンピュータ・システム・インターフェース(SCSI)またはファイバ・チャンネルのようなバスまたはチャンネルである。このようなインターフェースのいくつかの機能は、異なる複数の製造者からの種々の装置が相互に交換できるように、かつ、全てが共通インターフェースに接続できるように、しばしば標準化される。このような標準は、典型的には、米国規格協会(ANSI)のような機構のいくつかの標準協会によって規定されている。

【 0 0 0 4 】

30

種々の記憶装置と種々のコンピュータとの間でデータを交換するための1つの標準インターフェースは、ファイバ・チャンネルである。いくつかの実施の形態では、ファイバ・チャンネル標準は、複数のアービトレテッド・ループ(以下で更に説明する)を含む。いくつかの実施の形態では、ファイバ・チャンネル標準は、データ転送を制御するSCSI類似プロトコルをサポートしている。

【 0 0 0 5 】

ファイバ・チャンネルは、スモール・コンピュータ・システム・インターフェース(SCSI)設計よりも顕著な利点を有する。ファイバ・チャンネルは、従来のSCSI設計の2~20メガバイト/秒に比較して、現在の約106メガバイト/秒に達するかなり高い帯域幅を提供する。ファイバ・チャンネルは、典型的なSCSI環境における最大7個または15個の装置に比較して、126個に現在達する装置(ホストを含む)を接続することができる点で、更に大きな接続機能を提供する。ファイバ・チャンネルは、単一コネクタで取り付けることができ、かつ、スイッチを必要としない。同軸導体を使用するファイバ・チャンネルは、SCSI環境における25メートルに達する最大総延長に比較して、装置間が30メートルに達する距離で動作し、また、全チャンネルにわたり光ファイバを使用すると10キロメートルに達する。

40

【 0 0 0 6 】

SCSI環境では、データ伝送での誤りはパリティを使用して検出されるのに対し、ファイバ・チャンネルでは、誤りはランニング・ディスパリティおよびサイクリック・リダンダンシー・コード・チェック(CRCチェック)情報によって認識される。更なる情報は、本

50

発明者のウェストビーおよび本出願人のシーゲート・テクノロジー（株）による「マルチポート設計においてCRC発生器を使用したCRCチェック」と題する米国特許第5,802,080号および「16B/20Bエンコーダ」と題する米国特許第5,663,724号に見い出すことができる。

【0007】

ファイバ・チャンネル・アービトレーテッド・ループ（FC - AL）は、バイト・オリエンテッドDC平衡（0, 4）ラン・レングス制限8B/10Bのパーティションによるブロック転送コード方式を採用した工業標準システムである。FC - ALは、106.25MHzのクロック周波数で動作する。8B/10Bエンコーダ/デコーダの一形式は、フラナジェックらによる「バイト・オリエンテッドDC平衡（0, 4）8B/10Bパーティ

10

【0008】

ファイバ・チャンネル・アービトレーテッド・ループ（FC - AL）は、それぞれ「ノード」と呼ばれる多数の装置を相互に接続可能にする。ノードは、ファイバ・チャンネル「トポロジ」（この後で定義されている）に接続可能にしたインターフェースを有するコンピュータ・システムのあらゆる装置（コンピュータ、ワークステーション、プリンタ、ディスク・ドライブ、スキャナなど）である。各ノードは、他のノードに対するアクセスを得るためにNLポート（ノード・ループ・ポート）と呼ばれる少なくとも1つのポートを有する。2以上のポートを相互に接続する構成要素は、ひとまとめに、「トポロジ」または

20

【0009】

これらのポートはファイバ・チャンネル・ノードにおける接続であり、データはファイバ・チャンネルを介して他のノード（外側世界）のポートに転送できる。代表的なファイバ・チャンネル・ドライブは、そのドライブのノード内にパッケージされた2つのポートを有する。各ポートは、情報をポートに搬送するものとポートから情報を搬送するものの一对の「ファイバ」を含む。各「ファイバ」はシリアル・データ接続であり、また、一実施の形態では、各ファイバは、実際には、同軸ケーブル（例えば、ノードが互いに接近しているときに使用される同軸銅導体）であり、他の実施の形態では、ファイバは、（例えば、

30

【0010】

各ファイバは1方向のみにデータを搬送するので、複数のノードはループに沿って相互に接続され、そこでは、これらのノードは、転送すべきデータを有するときは、ループを制御するためのアービトレーションをしなければならない。「アービトレーション」は、複数のノードを協調させて、どのノードがループを制御するのかを決定する処理である。ファイバ・チャンネル・アービトレーテッド・ループは、ハブまたはスイッチなしでループに多数のノードを取り付ける。これらのノード・ポートは、ポイント・ツー・ポイント・データ転送回路を確立するためにアービトレーション動作を使用する。FC - ALは、各ポートが回路を確立するために少なくとも最小必要機能を含む分散トポロジである。アービトレーション・ループ・トポロジは、2および126ノード・ポート間で、任意の数のノードを接続するために使用される。

40

【0011】

いくつかの実施の形態では、各ノードは、冗長性をもたす二重ポート（それぞれ別個のポ

50

ートに接続される)を含み、その結果、一方のループが故障すると、他方のループがループ役割を遂行することができる。二重ポートはまた、2つのホスト(例えば、2つのホスト・コンピュータ)が1台のドライブを共有することを可能にする。

【0012】

(フェアネス背景)

本発明の文脈における「フェアネス」は、ファイバ・チャネル・ループのような共有リソースが複数のユーザの各々に対して、例えば複数のディスク・ドライブ100のそれぞれに対して、タイムリ・ベースで、すなわち、不都合な遅延なしに、かつ、各ディスク・ドライブ100が転送するために必要とするデータ量に比例した時間量に対して、利用可能にされる動作モードである。「ブレイング・アンフェア」は、1ユーザ、例えば1台のディスク・ドライブが、過度に頻繁にまたは過度に長期間にわたって、共有リソース、例えばファイバ・チャネル・ループの制御を行う動作モードであり、したがって、ループを使用する必要がある他のディスク・ドライブがそれを使用することを妨げる。

10

【0013】

「フェアネス」は、各ポートが(優先度に基づいて)ループに対するアクセスを獲得可能にさせるとともに、(時間制限なしに)それが欲するだけの情報を送信可能にさせる。そのために、他の全てのポートがアクセスを使用してしまうまで、アービトレーション獲得を待機する必要がある。

【0014】

ループ用のプロトコルは、各L__ポート(ループ・ポートとも呼ばれ、L__ポートは、ノーマル・ループ・ポート「NL__ポート」、または、2つのループを相互接続するために一般に使用されるファブリック・ループ・ポート「FL__ポート」である。)がループをアクセスするために連続的にアービトレーションをすることを可能にさせる。優先度は、アービトレーション・ループ物理アドレス(AL__PA)に基づいて各参加L__ポートに割り付けられる。他の優先度によるプロトコルのように、これは、低い優先度のL__ポートがループに対するアクセスを獲得できないという状況に至る恐れがある。アクセス・フェアネス・アルゴリズムは、全てのL__ポートがアービトレーションしかつループに対するアクセスに勝つ機会を与えられるアクセス・ウィンドウを設定する。全てのL__ポートがループを1回アクセスする機会を有するときは、新しいアクセス・ウィンドウが開始される。L__ポートは、再びアービトレーションすることができ、究極的に新しいアクセス・ウィンドウにおいてループに対するアクセスに勝つ。あらゆるL__ポートが任意の1アクセス・ウィンドウにおいてループをアクセスすることを必要とすることはない。

20

30

【0015】

アクセス・フェアネス・アルゴリズムを使用するL__ポートがアービトレーションをしかつループに対するアクセスに勝ったときは、少なくとも1アイドルがL__ポートによって送信されてしまうまで、L__ポートは再びアービトレーションすべきでない。第1のL__ポートがアービトレーションに勝ちアイドルを送信する間の時間は、アクセス・ウィンドウである。アクセス・ウィンドウの初期リセットを阻止するために、特殊なアービトレーション・プリミティブ信号(すなわち、ARB(F0))を使用する。アクセス・フェアネス・アルゴリズムの詳細は、ファイバ・チャネルFC-AL仕様(ANSI標準:ファイバ・チャネルFC-AL X3T11/プロジェクト960D/Rev. 4.5およびファイバ・チャネルFC-AL2 X3T11/プロジェクト1133D/Rev. 6.3)に記載されている。

40

【0016】

アクセス・フェアネス・アルゴリズムは、ANSI X3.230, FC-PHがクラス1接続に関する時間を制限しないように、L__ポートがアービトレーションに勝てば、L__ポートがループを制御する時間を制限することはしない。しかしながら、ENDTOVよりも長くアクセスが拒否されると、アクセス・ウィンドウがリセットされて、L__ポートがアービトレーションを開始することができる。

【0017】

50

全ての L__ポートがフェアネス・アルゴリズムを実施するけれども、F L__ポートまたは N L__ポートは常にフェアネス・アルゴリズムを使用することが要求される。例えば、1 つの L__ポートが他の L__ポートよりも多くのループ・アクセスを要求するときは、L__ポートはアンフェアとなることを選択することができる。

【0018】

全ての L__ポート用のループに同等のアクセスを提供するために、各 L__ポートはアクセス・フェアネス・アルゴリズムを使用することが推奨される。L__ポートがアクセス・フェアネス・アルゴリズムを使用しているときは、それは「フェア」L__ポートと呼ばれる。

【0019】

フェア L__ポートがループに対してアービトレーションしかつループへのアクセスに勝ちかつ他の L__ポートがアービトレーションしていることを検出しないときは、その L__ポートは、既存の回路を無期限に開き続けるかこの回路を閉じ、かつ、ループの所有権を保持して（すなわち、再アービトレーションなしに）ループ上の他の L__ポートを開くことができる。

【0020】

フェア L__ポートがループに対するアクセスを有しかつ他の L__ポートがアービトレーションしていることを検出すると、L__ポートは最小可能時間でループを閉じることができる。L__ポートは、ループを閉じて、異なる L__ポートを開く前に次のアクセス・ウィンドウにおいて再びアービトレーションする。

【0021】

いくつかのループの構成は、若干の L__ポートが 1 回 / アクセス・ウィンドウよりもループに対して多くのアクセスを有することを必要とすることがある。これらの L__ポートの例は、限定的ではないが、サブシステム・コントローラまたはファイル・サーバを含む。

【0022】

L__ポートは、アクセス・フェアネス・アルゴリズムを使用しないように初期化することができる（または、一時的に選択することができる）。L__ポートがフェアネス・アルゴリズムを使用していないときは、それは「アンフェア」L__ポートと呼ばれる。アクセス・フェアネスに参加するか否かの判断は、実施する必要性に対して残される。

【0023】

アンフェア L__ポートがループに対してアービトレーションしかつループへのアクセスに勝ちかつ他の L__ポートがアービトレーションしていることを検出しないときは、その L__ポートは、既存の回路を無期限に開き続けるかこの回路を閉じ、かつ、ループの所有権を保持して（すなわち、再アービトレーションなしに）ループ上の他の L__ポートを開いてもよい。

【0024】

アンフェア L__ポートがループを制御しかつ他の L__ポートがアービトレーションしていることを検出したときは、アンフェア L__ポートは最小可能時間でループを閉じてもよい。アンフェア L__ポートは、ループの所有権を保持するとともに（すなわち、再アービトレーションなしに）、ループ上の他の L__ポートを開いてもよい。

【0025】

参加している F L__ポートは、常に、その A L__P A に基づきループ上の最優先 L__ポートである。F L__ポートは、そのトラヒックの大部分はファブリックの残りをもつので、アクセス・フェアネス・アルゴリズムを使用することを免除される。

【0026】

F L__ポートがループを制御しかつ他の N L__ポートがアービトレーションしていることを検出すると、F L__ポートは最小可能時間でループを閉じることができる。F L__ポートは最高優先度を有しかつフェアネスから免除されているので、それは、常に、アービトレーションに勝つことになる。したがって、他の N L__ポートとの通信が必要である場合には、F L__ポートは、ループに対するそのアクセスを保持するとともに（すなわち、再

10

20

30

40

50

アービトレーションなしに)、ループ上の他のNL__ポートを開いてもよい。

【0027】

したがって、特にダイナミック半二重機能を含むファイバ・チャンネル実施では、ループ・フェアネスを保持する構成に対する必要性が存在する。

【0028】

(発明の概要)

ループ・フェアネスを保持する方法および装置を説明する。いくつかの実施の形態は、ダイナミック半二重特徴を含む。本発明の第1の態様は、1以上のポートを有する第1のチャンネル・ノードを含むループ・フェアネスを保持する通信チャンネル・システムを含み、各ポートは、ファイバ・チャンネル・アービトレーテッド・ループ・シリアル通信チャンネルをサポートしかつこれに取り付けられている。これらのポートの1つは、そのポートの取り付けられたチャンネルを制御するためにアービトレーションし、チャンネル・ループの制御において、アービトレーションに勝つと、フェアネスを保持している装置は、所定量の使用が第1のポートと通信チャンネルとの間に発生したか否かに少なくとも一部基づいて通信チャンネルの制御を解放させる。

【0029】

本発明の他の態様は、ファイバ・チャンネル・アービトレーテッド・ループ・シリアル通信チャンネルのループを制御するためにアービトレーションすることと、第1のポートと通信チャンネルとの間で所定量の使用が発生したか否かに少なくとも一部基づいて通信チャンネルの制御を解放させることとを含む通信方法を提供する。

【0030】

上記システムおよび方法のいくつかの実施の形態では、所定量の使用は、第1の所定量のデータの転送を含む。いくつかのそのような実施の形態では、第2の所定量よりも少ないデータが転送されるべく残っているときは、チャンネルの制御の解放が禁止される。

【0031】

(詳細な説明)

好ましい実施の形態の以下の詳細な発明では、その一部を形成するとともに、説明の手段として本発明を実施可能な特定の実施の形態を示す添付図面を参照する。他の実施の形態を利用することもでき、かつ、本発明の範囲から逸脱することなく構造的な変更を行うことができることを理解すべきである。

【0032】

この出願で説明した本発明は、ハード・ディスク・ドライブ、ZIPドライブ、フロッピー・ディスク・ドライブ、光ディスク・ドライブ、CDROM(「コンパクト・ディスク読出し専用メモリ」)ドライブおよび他の任意形式のドライブ、複数のドライブ・システム(例えば、「インエクスペンシブ/インディペンデント・ディスク・ドライブのリダント・アレー」、すなわち、RAID構成)、または、複数のドライブと他の複数のドライブまたは複数の情報処理システムとの間でデータを通信する他の複数の装置を含むあらゆる形式のディスク・ドライブに有用である。いくつかの実施の形態では、本発明は、(例えば、複数のファイバ・チャンネル・ループを互いに接続するために使用される)ハブおよびスイッチ、ワークステーション、プリンタ、ファイバ・チャンネル・アービトレーテッド・ループ上に接続されている他の装置または情報処理システムのような非ディスク装置用のノード・インターフェースにおいて有用である。

【0033】

図1は、ファイバ・チャンネル・ノード・インターフェースを有するディスク・ドライブ100のブロック図を示す。更なる情報は、「ループ初期化および応答用の方法および専用フレーム・バッファ」と題する米国特許出願第09/193,387号に見い出すことができる。

【0034】

図1と図2とを参照すると、ファイバ・チャンネル・ループ・インターフェース回路1220(ノード・インターフェース1220とも呼ばれる)は、ループ初期化および応答用の

専用送信フレーム・バッファ 73 を含む。一実施の形態では、各ノード・インターフェース回路 1220 は、2つのポート 116 (ポート A およびポート B を表す) を含む。(「ループ初期化」は、1 以上の特殊非データ・フレームのシーケンスを送信し(かつ、これらのフレームに対する応答を監視し)てファイバ・チャンネル・ループを初期化する。「応答」は、他のノードからの指令または質問に回答して送出された非データ・フレームである。)ファイバ・チャンネル・アービトレーテッド・ループ通信チャンネル 1250 (ループ 1250 またはファイバ・チャンネル・ループ 1250 と呼ばれる)は、ディスク記録ドライブ 100 とコンピュータ(または、他の情報処理システム) 1202 または他の情報処理装置との間でデータを通信するために使用され得る。一実施の形態では、ファイバ・チャンネル・ループ 1250 はシリアル通信チャンネルであり、他の実施の形態では、ファイバ・チャンネル・ループ 1250 を実施するために 2 以上の並列ライン(または、「複数のファイバ」)が使用される。このような専用送信フレーム・バッファ 73 を備えると、他方のポートがデータ・フレームを送受信している間に、二重ポート・ノード 1220 の一方のポート 116 が初期化または応答フレームを送信することを可能にする。ポート 116 は、一方が入力データ用のライン 117 であり他方が出力データ用のライン 118 であるシリアル・ラインであり、両方のライン 117, 118 は、通信チャンネル・ループ 1250 に接続するとともにその一部を形成している。また、2ポート・ノードの各ポート 116 に対して、専用受信バッファ(53, 53', 55)が設けられている。(プライム符号による参照番号を有する各ブロック(例えば、53')はプライムのない対応するブロック(例えば、53)と同一機能を提供することに注意すべきである。)ファイバ・チャンネル 1250 からフレームとともに受信するサイクリック・リダンダンシー・コード情報は、1 以上のフレーム・バッファ(53, 53' または 55)の 1 つに記憶されたのち、チェックされてフレーム・バッファ(53, 53' または 55)にある間にデータの完全さを保証する。ループ 1250 の制御は、ループ 1250 の制御に対してアービトレーションするのに費やされる総時間量を減少させるために、プログラム可能な量のデータが送信に利用できる限り、保持される。

【0035】

いくつかの実施の形態では、ディスク・ドライブ 100 は、1 以上のディスク・プラッタ 134 とディスク・プラッタ当り 1 以上の磁気読出し/書込みトランスデューサ 150 とアーム・アクチュエータ・アセンブリ 126 とを有する磁気記憶ヘッド・ディスク・アセンブリ(HDA) 114 を含む。トランスデューサ(または、「ヘッド」)と HDA インターフェース 113 との間の信号は、データをディスク・プラッタ 134 におよびそれから転送する。したがって、いくつかの実施の形態の「ディスク・ドライブ」(例えば、図 1 のディスク・ドライブ 1256)は、HDA 114 と HDA インターフェース 113 (例えば、通常の SCSI ドライブ)とを含み、1 以上のそのような従来のディスク・ドライブ 1256 は、図 1 に示すように、ループまたはファイバ・チャンネル・トポロジに接続するために、外部ノード・インターフェース 1220 に接続される。他の実施の形態では、「ディスク・ドライブ」は、図 2 のディスク・ドライブ 100 によって代表され、ディスク・ドライブ 100 全体内のディスク・ドライブ 1256 と統合されるノード・インターフェース 1220 を含む。一実施の形態では、データは、続いて、オフ・チップ・バッファ 111 からおよびそれに転送される。本発明は、専用オン・チップ・バッファ 119 を備えており、これは、図示した実施の形態では、各ポート(すなわち、バッファ 53, 53')用の受信非データ・フレーム・バッファ 53 (代わって、「入力非データ・バッファ 53」と呼ばれる。)と、両ポート(他の実施の形態では、一度に 1 ポートによってのみ単一バッファが使用される。)により同時に使用可能とされる送信フレーム・バッファ 73 と、CRC チェッカ 596 とともに共有データ・フレーム・バッファ 55 とを含む。

【0036】

一実施の形態では、ループ・ポート・トランシーバ・ブロック 115 (すなわち、115 および 115')は、ポート A およびポート B を介してこれに接続されているファイバ・

10

20

30

40

50

チャンネル・ループ 1 2 5 0 (図 2 を参照) に転送するデータを直並列変換および並直列変換する複数のポート・トランシーバを含む。いくつかの実施の形態では、トランシーバ 1 1 5 は、外部トランシーバとして実施される。他の実施の形態では、これらのトランシーバは、ブロック 1 1 0 にオン・チップで配置されている。いくつかの実施の形態では、右側 (すなわち、図 1 のトランシーバ 1 1 5 または 1 1 5 ' に対して右側) インターフェースは、1 0 ビット幅である並列入出力信号である。他の実施の形態では、それらは 2 0 ビット幅である。ブロック 1 1 0 , 1 1 1 , 1 1 2 とポート A トランシーバ 1 1 5 とポート B トランシーバ 1 1 5 ' とは、総合して、ファイバ・チャンネル・ノード・インターフェース 1 2 2 0 を形成している。いくつかの実施の形態では、トランシーバ 1 1 5 , 1 1 5 ' は単一チップ 1 1 0 に集積される。他の実施の形態では、それらの直並列 / 並直列変換機能を含むトランシーバ 1 1 5 , 1 1 5 ' は、チップ 1 1 0 から分離した回路上に実施される。

10

【 0 0 3 7 】

他の実施の形態では、トランシーバ 1 1 5 は、シリアル・ループ 1 2 5 0 とチップ 1 1 0 との間の単なるインターフェースであり、1 0 ビット幅または 2 0 ビット幅のデータに対する直並列変換 / 並直列変換がオン・チップで実行される。

【 0 0 3 8 】

図 2 は、コンピュータ・システム 1 2 0 0 の概要図である。都合のよいことに、本発明はコンピュータ・システム 1 2 0 0 に使用するのによく適している。コンピュータ・システム 1 2 0 0 はまた、電子システムまたは情報処理システムと呼ばれてもよく、中央処理装置 (C P U) とメモリとシステム・バスとを含む。コンピュータ・システム 1 2 0 0 は、中央処理装置 1 2 0 4 とランダム・アクセス・メモリ (R A M) 1 2 3 2 と中央処理装置 1 2 0 4 およびランダム・アクセス・メモリ 1 2 3 2 を通信可能に接続するシステム・バス 1 2 3 0 とを有する C P U 情報処理システム 1 2 0 2 を含む。C P U 情報処理システム 1 2 0 2 は、ファイバ・チャンネル・ノード・インターフェース 1 2 2 0 を含む。1 以上のディスク記憶情報処理システム 1 0 0 ~ 1 0 0 ' のそれぞれは、1 以上のディスク・ドライブ装置 1 2 5 6 とファイバ・チャンネル・ノード・インターフェース 1 2 2 0 とを含む。

20

【 0 0 3 9 】

いくつかの実施の形態では、多数のディスク・ドライブ 1 2 5 6 は、装置 1 0 0 ' がディスク・ドライブの R A I D アレーとなるように、単一ノード・インターフェース 1 2 2 0 、例えば R A I D (インエクスペンシブ / インディペンデント・ディスク・ドライブのリダundant・アレー) 構成に接続されている。C P U 情報処理システム 1 2 0 2 は、内部入出力バス 1 2 1 0 と入出力バス 1 2 1 0 に取り付けられたいくつかの周辺装置、例えば 1 2 1 2 , 1 2 1 4 および 1 2 1 6 とを駆動する入出力インターフェース回路 1 2 0 9 を含む。これらの周辺装置は、ハード・ディスク・ドライブ、磁気光学ドライブ、フロッピ・ディスク・ドライブ、モニタ、キーボードおよび他のそのような周辺装置を含むことができる。任意形式のディスク・ドライブまたは他の周辺装置は、ここで説明するファイバ・チャンネル方法および装置 (特に、例えばファイバ・チャンネル・ノード・インターフェース 1 2 2 0 における改良) を使用することができる。各装置では、任意の所与のループ 1 2 5 0 に接続するために A ポートまたは B ポートを使用することができる。

30

40

【 0 0 4 0 】

システム 1 2 0 0 の一実施の形態は、オプションとして、第 2 の C P U 情報処理システム 1 2 0 2 ' (システム 1 2 0 2 と同一または類似している) を含み、C P U 情報処理システム 1 2 0 2 ' は、中央処理装置 1 2 0 4 ' (中央処理装置 1 2 0 4 と同一である) と、ランダム・アクセス・メモリ (R A M) 1 2 3 2 ' (R A M 1 2 3 2 と同一である) と、中央処理装置 1 2 0 4 ' と R A M 2 3 2 ' とを通信可能に結合するシステム・バス 1 2 3 0 ' とを有する。C P U 情報処理システム 1 2 0 2 ' は、それ自身のファイバ・チャンネル・ノード・インターフェース 1 2 2 0 ' (ノード・インターフェース 1 2 2 0 と同一である) を含むが、(ループ 1 2 5 0 から分離され独立した) 第 2 のファイバ・チャンネル・ループ 1 2 5 0 ' を介して 1 以上のディスク・システム 1 0 0 に接続されている (この説明

50

例では、それは、ディスク・システム 100' に単に接続されているが、他の実施の形態では、全ての装置またはディスク・システム 100 ~ 100' に接続されている)。この構成は、2つのCPU情報処理システム 1202, 1202' が各CPU情報処理システム 1202用の個別的なファイバ・チャネル・ループを使用して1以上のディスク・システム 100を共有可能にさせる。更に他の実施の形態では、全ての装置 100 ~ 100' と全てのCPU情報処理システム 1202 ~ 1202' とが両方のループ 1250, 1250に接続されている。

【0041】

本発明の一実施の形態を構築するのに使用されるファイバ・チャネル仕様は、下記のANSI標準を含む。

【0042】

【表1】

ファイバ・チャネル	FC-PH	X3T11/プロジェクト 755D/Rev. 4.3	
		物理的及びシグナリング・インターフェース	
ファイバ・チャネル	FC-AL	X3T11/プロジェクト 960D/Rev. 4.5	
		アービトレーテッド・ループ	
ファイバ・チャネル	FC-AL2	X3T11/プロジェクト 1133D/Rev. 6.3	
		アービトレーテッド・ループ	
ファイバ・チャネル	FCP	X3T10/	Rev. 012
SCSI用プロトコル		X3.2	
		69-199X	

【0043】

I. ダイナミック半二重によるループ・フェアネスの保持

本発明の一実施の形態は、ダイナミック半二重(DHD)コマンドまたは命令によりループ・フェアネスを保持する(または、等価的に提供する)フェアネス保持装置 175および方法と、DHDコマンドを使用したフェアネスを提供する方法とを提供する。

【0044】

ファイバ・チャネルは本質的に全二重リンクである(ファイバ対において複数のフレームが同時に両方向に移動する)が、いくつかのループ・ポート(L__ポート)の実施は1方向データ転送のみをサポートすることができる。あるポートが全二重動作を可能とするとときでも、ループ・アービトレーションを損ない得る状況が存在する。

【0045】

ダイナミック半二重(DHD)は、OPENステートにあるポートによってループ上に送信されるループ・プリミティブであり、OPENEDステートにあるポートに対してそれ以上送信すべきフレームを有していないことを表示する。DHDは、複数のポートが確立された回路を下記により更に効率的に使用することを可能にする。

1. 半二重データ転送のみを可能とするポートが再アービトレーションすることなく逆方向に複数のフレームを転送することを可能にする。

2. OPENステートにあるポートがそのデータ転送を完了したとしても、OPENEDステートにあるポートが全てのフレームを転送することを可能にする。

【0046】

OPENステートにあるポートは、通常、最初のCLRを送信してループを閉じる。全二重回路が存在するときは、OPENEDステートにあるポートは、CLSを受信して、バ

10

20

30

40

50

ッファ・ツー・バッファ・クレジット (B B _ _ C r e d i t - R _ _ R d y) がもうないというまでフレームを送信し続けることができる。O P E N E D ポートがもはやフレームを送信できなくなると、それはO P E N ポートにC L S を返送しなければならない。

【 0 0 4 7 】

C L S を送信するよりもクローズ・イニシアティブを転送するのに有用と思われる少なくとも2つの場合が存在する。

1 . いくつかの実施はノードで同時送受信転送を処理することができない。これらのノードは、しばしば、O P E N ポートに対してペンディングしているフレームを有するが、これらは半二重設計のため、存在する双方向回路の利点を利用できない。

2 . 全二重転送が可能であっても、O P E N ポートがC L S を送信すれば、O P E N E D ポートは既存のクレジットに基づいてフレームを送信することのみが可能である。(例：ポートが新しいコマンドを受信するようにO P E N E D であり、2つのR _ _ R d yを受信し、O P E N ポート用のデータを読み出している。O P E N E D ポートは、それが閉じなければならない前に、2つのデータ・フレームを送信することのみが可能である。)

【 0 0 4 8 】

この追加的な再アービトレーション・サイクルを避けるために、D H D プリミティブ信号が供給される。D H D は、C L S を送信する代わりにO P E N _ _ L ポートによって送信される。D H D を送信すると、O P E N ポートはR _ _ R d y およびリンク制御フレーム(ただし、データ・フレームなし)を送信し続けることができる。O P E N E D ポートは、それがD H Dを受信しクローズ・イニシアティブを保持していることを記憶しており、もはやO P E N ポートに送信すべきフレームがないときは最初のC L S を送信することが期待されている。

【 0 0 4 9 】

この記述における「フェアネス」は、ディスク・ドライブ100のような複数のユーザのそれぞれに対して適時ベースにより、すなわち好ましくない遅延なしに、かつ、各ディスク・ドライブ100が転送する必要があるデータ量に比例した時間量に対して、ファイバ・チャンネル・ループ1250のような共有リソースを利用できるようにした動作モードである。「プレイング・アンフェア」は、1ユーザ、例えばC P U 情報処理システム1202(イニシエータとして作動する)の1ポートが、全ての他のユーザが自分の順番が来るのを待機することなく、共有リソースの制御を獲得する動作モードである。プレイング・アンフェアは、例えばC P U 情報処理システム1202が多数の低速度ディスク・ドライブ100上で動作を開始可能にさせるので、しばしば、制限されたある種の(例えば、高速度装置が低速度装置より多くの順番が巡って来るのを可能にする)環境においては、好ましいものとなる。1ディスク・ドライブ100がアンフェアを実行するのを可能にすると、その装置からのデータ転送がそうでないときよりも速く完了し得るが、他のディスク・ドライブ100がサービスを受けられず、他の動作を実行するように開放させないため、また、C P U 情報処理システム1202が他の複数のディスク・ドライブ100からデータを得るまでに待機するのが長くなり過ぎてしまうために、一般的には、システム・パフォーマンスを低下させる結果となる(図2を参照)。

【 0 0 5 0 】

本発明で目的とするファイバ・チャンネル・アービトレーテッド・ループ設計では、ループ1250をアクセスするために、ループ・ポート1220はアービトレーションしなければならない。どのポート1220がループ1250の制御を獲得するのかを判断するために優先度システムが使用され、また、「フェアネス」方式は複数のポートが飢えていないことを保証するために使用される。目標装置のように、ディスク・ドライブ100は、通常、フェアに行動し、低い優先度が与えられ、これが各ドライブ100がループ1250に対して等しくアクセスをするのを保証する。イニシエータ(例えば、C P U 情報処理システム1202)は、ドライブ・キューをフルに保持するために、アンフェアに行動することができる(ドライブ・キューは、各ディスク・ドライブ100がそれに向けられたコマンドを保持するペンディングおよびカレント・オペレーションのリストである)。「イ

10

20

30

40

50

ニシエータ」は、他の装置（「目標」と呼ばれる）によって実行されるべき入出力（I/O）処理を要求するファイバ・チャネル・ループ上の装置である。

【0051】

「全二重」は、データがポート1220に受信されているのとちょうど同時にデータがそのポートから送信されることができる動作モードである。「半二重」は、データが交互的にのみであって同時ではなくポートから送信されたりそのポートに受信されることができる動作モードである。

【0052】

本発明の一態様は、ファイバ・チャネル・アービトレーテッド・ループ・アーキテクチャ（FC-AL-2）に付加される「ダイナミック半二重」と呼ばれる特徴である。ポートが（下記章III, IIII, IVおよびVで説明されている設計のような）全二重動作を可能とするときであっても、DHD特徴を設けることによってループ・アービトレーション・サイクルの数を減少させることができる状態がある。例えば、OPENポート（OPEN状態にあるポート116）がCLSプリミティブ信号を送信するとき、OPENEDポートは既存のクレジットに基づいてフレームを送信できるだけである。（例えば、ノード・インターフェース1220のAポート116はOPENEDであって新しいコマンドを受信し、それは2つのR__RDYを受信するとともにOPENポート用のデータを読み込んでいた。OPENEDポートは、それが閉じられなければならない前に、2つのデータ・フレームを送信できるだけである。）ダイナミック半二重により、OPENポートは、CLSの代わりにプリミティブDHDを送信するとともに、R__RDYを送信し続ける。これは、データを読み出していたディスク・ドライブ100が、ループに対してアービトレーションすることを必要とせず新しいコマンドを受信するときに転送を完了することを可能にする。このDHD特徴は、アービトレーション・サイクルを減少させることができるが、ループ・フェアネスを損なう恐れがある。

【0053】

例えば、本発明によるシステム1200（図2を参照）では、情報処理システム1202は、（複数のデータ・フレームを返送する）ディスク・ドライブ100に読み出し動作コマンドを送出するとともに、（ディスク・ドライブ100が、制御を放棄することなくかつループ1250の制御のために再アービトレーションする必要なく多数のデータ・フレームの転送を完了するためにループ1250上の制御を延長された期間維持することを可能にする）DHDプリミティブ信号を送出する。したがって、DHDプリミティブ信号を受信したディスク・ドライブ100は、それが既にアービトレーションに勝ちかつ次のフェアネス・ウィンドウを待機していたとしても、再アービトレーション・ウィンドウを再び使用可能にされるからである。ループのフェアネス（全ての装置がループ・リソースに対して同等のアクセスを有する能力）は損なわれる。

【0054】

最高の優先度アドレスを有しかつアンフェアを行う（すなわち、最高の優先度アドレスを使用して、それらが長時間その順番が来るのを待機していたとしても最低の優先度アドレスを有する他の装置の犠牲によりループの制御を行うことによって）イニシエータ情報処理システム1202（図2を参照）は、アービトレーションに勝って、ドライブ100に新しいコマンド（例えば、DHDコマンド）を送出する。ドライブ100は、イニシエータ1202用のデータを読み出していたので、ループ1250に対してアービトレーションを実行していた。ドライブ100は、DHDコマンドを受信すると、そのデータ・フレームの全てと、動作を完了する応答とを送出することが可能とされる。新しいコマンドにより開放にならなかった他のドライブ（例えば、ドライブ100'）は、アービトレーションで「フェア」に勝つまで待機する必要がある。一方、コマンドを完了したドライブ100は、次のコマンドに対するデータを読み出すことが可能であった。前のコマンドは完了したので、イニシエータ1202は、そのドライブ100に新しいコマンドを送出して、そのキューをフルに保持することができる。新しいコマンドは、ドライブ100が他の読み出し転送を完了することを可能にさせ得た。いくつかのドライブは、時間上で、ループ

10

20

30

40

50

1250に対するアクセスを「スターブド(starved)」にするとともに、DHDを受信するのに使用するものよりも少ないコマンドで終了することができる。

【0055】

本発明の改良されたフェアネス保持特徴は、1つのノード116がループ1250の制御を維持することを可能にさせられる時間を制限しており、したがって、他のノード116が合理的な時間内に順番が来ることを保証している(かくして、「スターベーション」(starvation))を回避する。一実施の形態では、一つのノード116がループに対する制御を維持することが可能にさせられる時間は、任意の1つの動作で転送されるフレームの最大数を調節することによって調節される。一実施の形態では、ノード116がそのデータ転送の完了に十分近いときには、そうでないときに許容されるフレームの最大数を超えるものであっても、その転送を完了することが可能にさせられる。

10

【0056】

本発明によれば、ダイナミック半二重を使用するときにループ・フェアネスを保持するために、数値によるフレーム限定特徴がノード・インターフェース1220に付加される。一実施の形態では、転送が完了に近くない限り、装置(例えば、ディスク・ドライブ100)がDHDコマンド(DHD「プリミティブ」とも呼ばれる。)を受信した後にいくつかの読出しデータ・フレーム(すなわち、送信すべきデータ・フレーム)が送信されることを許容されるかについて、その上限値が設定される。いくつかの実施の形態では、数値限定は、プログラム可能であり、(一実施の形態では)販売者固有ロゲイン・パラメータまたは(他の実施の形態では)モード・ページ設定に設定されるか、(他の実施の形態では)マイクロプロセッサを介してデフォルト値に設定される。これは、イニシエータ1220が、ドライブ100が使用されているアプリケーション(すなわち、CPU情報処理システム1202で実行されているソフトウェア)に対する数値限定をプログラム可能に調整することを可能にするか、そのドライブを合理的な値に設定することを可能にする。

20

【0057】

一実施の形態では、本発明は、図1に示すようなフェアネス保持装置175を提供する。一実施の形態では、16ビット・カウンタ60(ここでは、dhd__cnt__out(15:0))とも呼ばれる。図1を参照)は、DHDが受信されるときにint__zero信号159によってゼロにされるとともに、送信された各フレームに対して(インクリメント161を使用して)インクリメントされる。コンパレータ162は、DHDカウンタ160が限界レジスタ163(dhd__max__frms(15:0))とも呼ばれる。)に記憶された数値限定値に達した時をチェックして、データ転送を停止する(ループ1250を閉じる)時を判定する。コンパレータ162は、suspend__xfer信号164を出力する。一実施の形態では、suspend__xfer信号164は、データ・フレーム送信動作を停止させるとともにループを一時的に閉じるようにして、他の装置がループ1250を使用することを可能にする(すなわち、この実施の形態は、以下で説明する信号169および信号179を無視して、suspend__xfer信号172をsuspend__xfer信号164と同等にすることによって限界に達すると、転送を停止する)。他の実施の形態では、このフェアネス保持装置175用のenable__DHD__suspend信号170は、ANDゲート171への入力としても供給される(すなわち、この実施の形態は、以下で説明する信号169を無視して、イネーブル信号170が「オン」であれば限界に達して転送を停止し、イネーブル信号170が「オフ」であれば本発明のフェアネス保持特徴175をディセーブルする)。

30

40

【0058】

他の実施の形態では、(DHDコマンドが受信されたときにのみというよりも)他のデータ転送に対してカウンタ160が活性化される。いくつかの実施の形態では、カウンタ160は、以上で説明した限界レジスタ163にロードされる値をロードすることによって初期化される減算カウンタにより置換され、また、この減算カウンタは、フレームが転送される毎に1つデクリメントされ、ゼロになると信号164をアクティブにし、したがって以上で説明したブロック160, 161, 162, 163と等価な機能を提供する。

50

【 0 0 5 9 】

いくつかの実施の形態では、カウンタ 1 6 0 は、上述したように、転送されたフレームの数をモニタする。他の実施の形態では、カウンタ 1 6 0 は、転送されたデータ量のバイト若しくはワードの数または他の測定量をモニタする。

【 0 0 6 0 】

更に他の実施の形態では、カウンタ 1 6 0 は、転送中に経過した時間をモニタするとともに所定の時間が過ぎると信号 1 6 4 を活性化するタイマによって置換される。これは、以上で説明したものと同様の機能を提供するが、転送されたデータ量によるというよりも時間に基づくフェアネスを提供する。いくつかのそのような実施の形態では、所定量よりも少ないデータが転送されるべく残っているとき、および/または、enable__DHD__suspend 信号 1 7 0 がディセーブルされたとき、AND ゲート 1 7 1 はそのまま使用されて suspend__xfer 信号を禁止する。

10

【 0 0 6 1 】

いくつかの実施の形態では、カウンタ 1 6 5 は、上述したように、転送されたフレームの数を監視する。他の実施の形態では、カウンタ 1 6 5 は、バイト若しくはワードの数または転送されたデータの量の他の測定値を監視する。

【 0 0 6 2 】

更に他の実施の形態では、カウンタ 1 6 5 は、転送中に経過した時間を監視するとともに所定の時間が過ぎると信号 1 6 9 を活性化するタイマで置換される。これは、上述したものと同様の機能を提供するが、転送されたデータの量というよりも時間に基づいたフェアネスを提供する。いくつかのそのような実施の形態では、所定量よりも少ないデータが転送されるべく残っているとき、および/または、enable__DHD__suspend 信号 1 7 0 がディセーブルされたときは、AND ゲート 1 7 1 がそのまま使用されて suspend__xfer 信号を禁止する。

20

【 0 0 6 3 】

一実施の形態では、転送長カウンタ 1 6 5 ロジックは、フレームの代わりにワードの数を使用する。

【 0 0 6 4 】

上述の説明はダイナミック半二重特徴を含む実施の形態について特になされたが、本発明のフェアネス保持特徴 1 7 5 が DHD 特徴から独立して提供される本発明の他の実施の形態が存在する。例えば、いくつかの実施の形態では、フェアネス保持特徴 1 7 5 が提供されるが、DHD 特徴は提供されない。他の実施の形態では、両者は提供されるが、これらの状況に対するフェアネスを強固にするために、DHD コマンドを受信したときおよび他の一定の状況においてフェアネス保持特徴が使用される。例えば、一実施の形態では、情報処理システム 1 2 0 2 からドライブ 1 0 0 への書込みデータの転送である。

30

【 0 0 6 5 】

一実施の形態では、転送長カウンタ 1 6 5 (x f r __ c n t __ o u t (2 6 : 0) と呼ばれる。) が提供される。転送長カウンタ 1 6 5 は、ロジック 1 6 6 によってデクリメントされて、送信すべきワードの残りの数を示す。転送長カウンタ 1 6 5 およびデクリメント 1 6 6 とはデータ転送フレーム長カウンタ 8 2 を形成する。この実施の形態は、転送長カウンタ 1 6 5 における残りのワードのカウントを最小長レジスタ 1 6 8 (d h d __ c m p l __ l e n (2 6 : 0) と呼ばれる。) に保持されたプログラマブル値と比較して、転送が完了に近いときには転送の停止を阻止する (すなわち、転送されるべく残っているワードが最小長レジスタ 1 6 8 に保持された値よりも小さいときには、その転送は「完了に近い」と定義される。) 。一実施の形態では、AND ゲート 1 7 1 は、信号 1 6 4 と信号 1 6 9 との論理積を形成して suspend__xfer 信号 1 7 2 を導き出し、続いてこれがデータ・フレーム送信動作を停止させて、ループを一時的に閉じて他の装置がループ 1 2 5 0 を使用することができるようにする (すなわち、この一実施の形態は、以下で説明する信号 1 7 0 を無視し、残りのワードがレジスタ 1 6 8 における値よりも小さくなっていない限り、限界に達すれば、転送を停止する。) 。他の実施の形態では、フェア

40

50

ネス保持特徴 175 に対する `enable_DHD_suspend` 信号 170 は、AND ゲート 171 への入力としても供給される。転送を停止させる論理式（ゲート 171 の出力）は、以下のようになる。

【0066】

【数 1】

```
suspend_xfer signal 172 = (enable_DHD_suspend = '1'
    AND (dhd_cnt_out(15:0) = dhd_max_frms(15:0))
    AND (xfr_cnt_out(26:0) ≥ dhd_cmpl_len(26:0)) )
```

10

【0067】

したがって、ループ・アービトレーション・フェアネスを保持するために、DHD プリミティブが受信され、かつ、読出しデータが「開」ポート（このポートは「開放」である。）に利用可能であれば、このポートは、`dhd_max_frm` カウントに達するまでデータを送信したのち、ループを閉じる。転送が終りに近い（すなわち、`dhd_cmpl_len` における値よりも小さい）ときは、このポートは、直ちに閉じるよりも転送を完了させる。転送はそのように行われるので、再びアービトレーションを行う必要はない。`dhd_max_frm` および `dhd_cmpl_len` における値は、モード・ページ初期化によるかログ・イン値によりセットされて、この特徴をもっと柔軟性のあるものにする。

20

【0068】

データ・フレームが送信されるべきであるとき、DHD プリミティブが受信されると、カウンタ 60 がリセットされ（すなわち、ゼロの値にセットされ）、次いで、カウンタ 160 は、送信された各データ・フレームに対して“1”だけインクリメントされる。カウンタ 160 が（転送が終りに近くなっていない限り）限界レジスタ 163 に保持されている最大許容値に達すると、ループ 1250 に対する次の時間アービトレーションに勝つまで、それ以上のデータ・フレームは送信されない。

【0069】

「`suspend_xfer`」出力信号 172 は、ポート A, B 開放制御回路 42 への入力としてループ制御回路 40 によって使用されて（図 5 および以下の説明を参照）、DHD が受信されると、ループ上に閉を送信する。そうでなければ DHD モードは非常に長い転送を許容してループのフェアネスを損なうので、`suspend_xfer` 信号 172 は、ループ・フェアネスを回復すなわち保持する。

30

【0070】

一実施の形態では、`suspend_xfer` 信号 172 はまた、入力ポート開制御ステート・マシン 42 / 42' として供給されて、ポート開制御ステート・マシン 42 / 42' にループを「閉」にさせ（すなわち、他の装置のポートがそのループに対するアービトレーションを可能にするためにループ 1250 の制御を解放するようにさせ）、これによって、ポートがループ 1250 の制御をアンフェアに維持することを防止する。いくつかの実施の形態では、所定量のデータ（例えば、ワード数）が転送されると、`suspend_xfer` 信号 172 が活性化される。他の実施の形態では、所定の時間が経過すると、`suspend_xfer` 信号 172 が活性化される。これらの形式の実施の形態のうちのいくつかでは、転送されるべき所定量よりも少ないデータ（例えば、フレーム数）が残っているときは、`suspend_xfer` 信号 172 の活性化が禁止される。上記形式の実施の形態のうちのいくつかでは、`enable_DHD_suspend` 信号 170 がディセーブル（不活性）されると、`suspend_xfer` 信号 172 の活性化が禁止される。

40

【0071】

II. ループ初期化および応答用の専用フレーム・バッファ

50

本発明の一実施の形態に関して、フレーム・バッファは、第3世代特定用途向け集積回路（ASIC）に付加されて、両方のポートが同時に活性であることを可能にさせる。非データ・フレーム（図1の「受信非データ・フレーム・バッファ」53, 53'とも呼ばれる。）を受信する2つのバッファは、ノードの両方のポートで同時に複数のコマンドおよび複数のFCPフレーム（*fiber-channel-protocol*フレーム）を受信できるように（および、全二重動作ができるように、すなわち、あるポートの一方のファイバを介して受信すると同時に、同じポートの他方のファイバを介して送信するように）設けられている。これは、ディスク・ドライブ100が転送の中断または終了まで待機するよりも同じポートを介したおよび/または他のポートを介したデータ転送中に一方のポートを介して新しいコマンド（または、他の非データ・フレーム）を受信できるようにさせる。従来のアプローチにおけるよりも早くコマンドを有することによって、本発明は、データ転送が進行している間にコマンドを分類して最適化することを可能にさせ、これによって、システム1200の性能を改善する。

【0072】

図3は、ファイバ・チャンネル・ノード・インターフェース・チップ110のブロック図である。本発明におけるファイバ・チャンネル・ノード・インターフェース・ロジック110は、アービトラレーテッド・ループ・ロジックおよびフレーミング・ロジックを含むファイバ・チャンネル・プロトコルに対して責任を持つ。一実施の形態は、ファイバ・チャンネル・プロトコル（FCP）標準によって定義されたSCSI上位プロトコルのみを使用してクラス-3SCSI実施（上述したFCAL仕様）に対して最適化される。ファイバ・チャンネル・ノード・インターフェース・ロジック110は、種々のバッファ帯域幅をサポートするとともに、二重ポートおよび全二重動作を支援するように4つのオン・チップ・フレーム・バッファ（53, 53' 55および73）を含む。ファイバ・チャンネル・ノード・インターフェース・ロジック110はまた、マイクロプロセッサ112に対してインターフェースを行い、マイクロプロセッサ112がファイバ・チャンネル・ノード・インターフェース・ロジック110を構築するとともにファイバ・チャンネル・ノード・インターフェース・ロジック110の現在状態についてのステータス情報を読み出せるようにする。

【0073】

ファイバ・チャンネル・ノード・インターフェース・ロジック110は、2つのループ・ポート回路20（一方はポートA用で、他方はポートB用であり、各ポートはループ通信をサポートするためにデータ入力インターフェースおよびデータ出力インターフェースを有する。）と、ループ制御回路40（フレーム送信回路40とも呼ばれる。）と、受信パス・ロジック50と、転送制御回路60と、単一フレーム送信回路70と、送信パス・マルチプレクサ（*mux*）79と、データ・フレーム送信パス・ロジック80と、マイクロプロセッサ・インターフェース90とを有する。これらのブロックは、受信フレーム処理、送信データ・フレーム発生、単一送信フレーム発生、転送制御およびプロセッサ・インターフェース処理のような機能をサポートする。

【0074】

マイクロプロセッサ・インターフェース回路90は、マイクロプロセッサ112がファイバ・チャンネル・ノード・インターフェース・ロジック110のレジスタおよびカウンタにアクセスできるようにする。（「マイクロプロセッサ」を説明するとき、そのような用語は任意の適当なプログラマブル・ロジック装置を含む。）インターフェース・レジスタは、ファイバ・チャンネル・インターフェースの応答よりも前に外部マイクロプロセッサ112によって初期化される。出力転送はこのインターフェースを介して初期化され、また、受信された転送のステータスはこのインターフェースを介して入手可能である。

【0075】

図3用の入力信号は、ファイバ・チャンネル16からのデータ入力をポートA用のループ・ポート回路20に搬送するA__IN3021と、ファイバ・チャンネル16からのデータ入力をポートB用のループ・ポート回路20に搬送するB__IN3022とを含む。DATA FROM OFF CHIP BUFFER3051は、オフ・チップ・バッファ1

10

20

30

40

50

11からのデータを受信パス50に搬送する。TOFFCHIPBUFFER3052は、受信パス50からオフ・チップ・バッファ111にデータを搬送する。BUFFERSTATUS3061は、転送制御60にステータスを提供する。MPUIンターフェース90へのMPUADDRESS3091およびMPUDATA3095はマイクロプロセッサ112からのアドレスおよびデータをそれぞれ供給する。MPUIンターフェース90へのREAD__ENABLE3092およびWRITE__ENABLE3093はマイクロプロセッサ112からのイネーブル信号を供給する。信号MPU3076は、マイクロプロセッサ112が送信フレームバッファ73をアクセスすることを可能にさせる。A__OUT3023は、ポートAに対してループ・ポート回路20からファイバ・チャンネル16にデータを搬送し、また、B__OUT3024は、ポートBに対してループ・ポート回路20からファイバ・チャンネル16にデータを搬送する。

10

【0076】

図4は、ファイバ・チャンネル・ループ・ポート回路20のブロック図である。本発明の一実施の形態のファイバ・チャンネル設計は、周辺装置の直接アタッチメント用の二重ポート・ファイバ・チャンネル・インターフェースをサポートするために2つの同じループ・ポート回路20を含む。一実施の形態では、ファイバ・チャンネル・ループ・ポート回路20は、受信レジスタ21と、8B/10Bデコーダ・ロジック22と、ワード同期ステート・マシン23と、受信クロック喪失検出器24と、同期喪失タイマ25と、アービトレーテッド・ループ・ロジック26と、8B/10Bエンコーダ27とを含む。

【0077】

20

一実施の形態では、各ループ・ポート回路20は、10ビット・データ・インターフェースを使用して外部トランシーバ115（図1を参照）にインターフェースする。そのような実施の形態では、トランシーバ115は、並列インターフェース（例えば、10ビット幅または20ビット幅インターフェース）におよびそれからシリアル・データを並直列変換しかつ直並列変換する。他の実施の形態では、これらのトランシーバ115はチップ110に集積される。（ファイバ・チャンネルから入力される）並列データは、各トランシーバ115の受信機部からの受信クロックを使用して捕捉されて、並列8B/10Bデコーダを使用してデコードする前に20ビット幅フォーマットに変換される。次いで、（特別な順序集合（ordered set）を表すために使用された）16ビット・データ+2kキャラクタは、アービトレーテッド・ループ・ロジック26に配置される前にワードの有効性についてチェックされる。アービトレーテッド・ループ・ロジック26の出力は、送信機クロックに再同期されて、受信フレーミング・ロジックに渡されてもよいし、エンコーダ27を介してループ1250上に再送信されてもよい。一実施の形態では、エンコーダ27は、各動作中に1つの8ビット・キャラクタを1つの10ビット・キャラクタに変換する。他の実施の形態では、2つ以上の8ビット・キャラクタが、対応した数の10ビット・キャラクタに各動作において変換される。（「16B/20Bエンコーダ」と題する米国特許第5,663,724号を参照。）アービトレーテッド・ループ・ロジック26は、ループ・ステート・マシンと順序集合デコーダと弾性挿入および削除機能とを含む。ループ・ポート回路20は、ファイバ・チャンネル・アービトレーテッド・ループANSI標準（すなわち、上述したFC-ALおよび/またはFC-AL2）に定義されているアービトレーション・ループ・プロトコルを実施する。

30

40

【0078】

一実施の形態では、ファイバ・チャンネル・データは、シリアルに送信されて、トランシーバ115によって10ビット並列データに変換される。受信レジスタ21は、トランシーバ115の受信機部によって発生されたクロックを使用してトランシーバ115からの10ビット・データ（A__IN3021またはB__IN3022）を捕捉する。データは、8B/10Bデコーダ22を通過する前に20ビット幅（すなわち、10ビット・キャラクタ幅）に直ちに变換される。一実施の形態では、デコーダ22は「8B/10Bデコーダ」と呼ばれているが、各動作中に1つの10ビット・キャラクタを1つの8ビット・キャラクタに変換する。他の実施の形態では、2つ以上の10ビット・キャラクタは、対応

50

する数の 8 ビット・キャラクタに各動作において変換される。

【 0 0 7 9 】

8 B / 1 0 B デコーダ 2 2 は、受信レジスタ 2 1 によって捕捉された符号化データを入力する。2 つの 1 0 ビット・キャラクタは、並列にデコードされて、2 つの 8 ビット・キャラクタを出力する。入力キャラクタのランニング・ディスパリティがチェックされ、また、エラー・ステータスがワード同期ステート・マシン 2 3 とアービトレテッド・ループ・ロジック 2 6 とに渡される。負のランニング・ディスパリティは、ランニング・ディスパリティ・エラーに従って次の順序集合上に強制設定される。符号化規則の違反もチェックされ、コード違反ステータスがワード同期ステート・マシン 2 3 に渡される。

【 0 0 8 0 】

受信クロック喪失検出器 2 4 は、トランシーバ 1 1 5 からの受信クロックが停止した時を検出する。受信クロック喪失が検出されると、ワード同期ステート・マシン 2 3 がリセットされて、データはアービトレテッド・ループ・ロジック 2 6 の F I F O に行くのが阻止される (F I F O は、ファースト・イン・ファースト・アウト・メモリであり、典型的には、異なる速度を有するバスまたは処理間をインターフェースするのに使用される)。カレント・フィル・ワード (C F W、以下で更に説明する) は、ワード同期が再獲得されるまで送信される。

【 0 0 8 1 】

ワード同期ステート・マシン 2 3 は、ワード同期用の入力ストリームを監視する。3 つの有効な順序集合が正しいバイト / 制御キャラクタ・アライメントで検出されたときに、ワード同期が検出され、かつ、干渉のない無効キャラクタが検出される。「ワード同期喪失」が F C - P H (すなわち、F C - P H 物理およびシグナリング・インターフェース X 3 T 1 1 / プロジェクト 7 5 5 D / R e v . 4 . 3) 標準によって定義されている。ワード同期が達成されると、データがアービトレテッド・ループ・ロジック 2 6 の F I F O に入力される。

【 0 0 8 2 】

同期喪失タイマ 2 5 は、(3 つの有効な順序集合を検出するのに 1 フレーム時間かかるかもしれないので) ワード同期喪失条件が 1 フレーム時間以上の間存在した時を判断するために使用される。このタイマが時間切れとなったときは、マイクロプロセッサ 1 1 2 は、L O S S - O F - S Y N C 割込み信号 4 0 2 5 で割り込みされる結果、処理を開始することができる。

【 0 0 8 3 】

アービトレテッド・ループ・ロジック 2 6 は、ループ弾性 F I F O と、ループ F I F O 制御ロジックと、順序集合デコード・ロジックと、ループ・ステート・マシン・ロジックと、カレント・フィル・ワード選択ロジックと、ループ出力マルチプレクサ・ロジックと、種々の機能とを含む。ループ弾性 F I F O は、(受信クロックによってクロッキングされた) 入力データを送信クロックと再同期するために必要なバッファリングを提供する。ループ F I F O 制御ロジックは、アービトレテッド・ループ・ロジック 2 6 のステートを監視して、挿入または削除動作を必要であるか否かを判断する。順序集合は、順序集合認識ロジックによってデコードされる。これらの順序集合は、フレーム・デリミタおよびアービトレテッド・ループ順序集合を含む F C - P H 定義順序集合 (すなわち、F C - P H 物理およびシグナリング・インターフェース X 3 T 1 1 / プロジェクト 7 5 5 D / R e v . 4 . 3) を含む。カレント・フィル・ワード選択ロジックは、ループ・ステータスおよびデコードされた順序集合を監視してカレント・フィル・ワード (C F W) を決定する。アービトレテッド・ループがイネーブルされると、ハードウェア・ステート・マシンは、順序集合デコーダを使用して、F C - A L 標準 (すなわち、ファイバ・チャンネル F C - A L 1 アービトレテッド・ループ標準 X 3 T 1 1 / プロジェクト 9 6 0 D / R e v . 4 . 5 またはファイバ・チャンネル F C - A L 2 アービトレテッド・ループ標準 X 3 T 1 1 / プロジェクト 1 1 3 3 D / R e v . 6 . 3) に説明されている複数のループ機能を実行する。入力 L O O P A T R A N S M I T C O N T R O L O U T P U T 6 4 2

10

20

30

40

50

5 および LOOP B TRANSMIT CONTROL OUTPUT 6 4 2 7 は、図 5 においてロジックからアービトレテッド・ループ・ロジック 2 6 への入力を供給する。出力 LOOP A STATES AND CONTROL 6 4 2 2, LOOP B STATES AND CONTROL 6 4 3 2 は、各ループの出力を制御して、ループ制御ロジックにステータスを提供し、続いて、これがループ制御 4 0 に要求を発生する（図 5 を参照）。出力 LOOP A DATA 4 0 2 6, LOOP B DATA 4 0 2 7 は、データを各ローカル・ポートに供給する。

【 0 0 8 4 】

一実施の形態では、8 B / 1 0 B エンコーダ・ロジック 2 7 は、アービトレテッド・ループ・ロジック 2 6 から 1 6 ビット・データおよび 2 k キャラクタ（下位 k が常に 0 である。）を受け取る。一実施の形態では、入力は、2 つの 1 0 ビット・キャラクタに符号化され、これらは分離されてトランシーバ 1 1 5（図 1 参照）に 1 度に 1 つ出力され、トランシーバ 1 1 5 はデータをシリアル・ストリームに変換する。他の実施の形態では、両方とも 1 0 ビット・キャラクタ（すなわち、2 0 ビット）が並列にトランシーバ 1 1 5 に送出され、トランシーバ 1 1 5 はデータをシリアル・ストリームに変換する。送信マルチプレクサ 7 9（図 3 を参照）はまた、エンド・オブ・フレーム（E O F）・デリミタが転送されている時を表すためにステータスを提供して、エンコーダ 2 7 が現在ランニング・ディスペリティに基づいて正しい形式（すなわち、フレイバ（f l a v o r））の E O F を選択できるようにする。また、（O p e n d ステートにある）ポートが送信しているとき、または、アービトレテッド・ループ・ロジック 2 6 がプリミティブを送信しているときは、各非 E F O プリミティブの開始時にランニング・ディスペリティが強制的に負にされる。出力信号 A _ _ O U T 3 0 2 3 および B _ _ O U T 3 0 2 4 は各トランシーバ 1 1 5, 1 1 5 ' にデータを送信する。

【 0 0 8 5 】

図 5 は、ループ制御回路 4 0（フレーム送信（X M I T）回路 4 0 とも呼ばれる。）のブロック図である。ループ制御回路 4 0（図 3 および図 5 を参照）は、複数フレームすなわち R _ _ R D Y を送信し始めるために（ポート A およびポート B のアービトレテッド・ループ・ロジック 2 6 における）適当なアービトレテッド・ループ・ステート・マシンに対する要求を発生するとともに送信フレーミング・ステート・マシン 7 2, 8 1 に対する要求を発生する制御ロジックを含む。

【 0 0 8 6 】

送信データ・シーケンサ・ロジック 4 1 は、転送がマイクロプロセッサ 1 1 2 によって要求されるときに活性化されるロジックを含む。送信データ・シーケンサ・ロジック 4 1 は、入力信号 TRANSMIT STATUS INPUTS 6 4 1 1 を使用して転送を監視するとともに、転送の各段階に関して「イネーブル」（すなわち、イネーブル信号 TRANSMIT CONTROL OUTPUTS 6 4 1 3）を発生する。これは、転送レディーおよび F C P 応答がマイクロプロセッサ 1 1 2 の介入なしに発生されることを可能にさせる。

【 0 0 8 7 】

ループ・ポート A / B 開放制御ステート・マシン 4 2（ポート A）, 4 2'（ポート B）は、ポートが他の L _ _ ポートによって開放される場合、または、ループ 1 2 5 0 がフレームを送信するために開放される場合を取り扱う。このロジックは、アービトレーションする要求とループ 1 2 5 0 を閉じる要求と R _ _ R D Y および種々のフレームを送信する要求とを発生し、半二重または全二重動作に対して構成することができる。

【 0 0 8 8 】

次の条件は、アービトレーションする要求を開始するために満足される必要がある。

- マイクロプロセッサ 1 1 2 から送信ポート・イネーブル付きフレームを送信する要求
- 送信ポートが監視ステートにある
- 転送長カウンタがゼロでない
- 転送を中断するマイクロプロセッサ 1 1 2 からの要求がない

10

20

30

40

50

- (非データ転送、または、満足されたデータしきい値でまだ送信されていない転送レディ付きデータ書込み転送、または、満足されたデータしきい値および満足されたデータ・フレーム・バッファしきい値でのデータ読出し転送)。

【0089】

ポートが半二重モード用に構成されているときは、R__RDYはOpened状態においてのみ送信することができる。ポートが全二重モード用に構成されているときは、Opened状態またはOpen状態においてR__RDYを送信することができる。R__RDYを送信させる条件は、“利用可能なBuffer-to-Buffer Credit(BB__Credit)および最大BB__Creditよりも小さい未解決R__RDY”を含む。(Buffer-to-Buffer制御ロジック603は、接続された

10

【0090】

ポートが半二重モード用に構成されているときは、Open状態においてのみフレームを送信することができる。ポートが全二重モード用に構成されているときは、Open状態において、または、ポートがフレーム受信者によって全二重モードで開放されればOpened状態において、フレームを送信することができる。

【0091】

フレームを送信する要求は、次の条件が全て満足されたときに発生される。

- データ・フレーム・バッファ55が、利用可能なデータを有する。
- Buffer-to-Buffer Creditが利用可能である(受信されたR__RDY)。
- (ブロック609において)非データ転送またはデータ読出し転送および転送長カウンタがゼロでない。

20

【0092】

ループ1250を閉にする(ポート116によって解放されるべき通信チャネルの制御)条件は、次のものを含む。

- Opened状態に入ったときにBuffer-to-Buffer Creditが利用可能でない。
- フェアネスのために、すなわち、(図1の制限レジスタ163で指定された)転送可能とされるフレームの数についての所定の制限がDHDカウンタ160によって到達し、かつ、(オプションとして)(最小長レジスタ163によって指定された)最小長よりも多いフレームの数が転送されるべく残っており、かつ、(オプションとして)enable__DHD__suspend信号170が活性である。
- 未解決R__RDYがなく、かつ、Opened状態のときにそれ以上のBB__Creditが利用できない。
- ポートがOpened状態にあるときにプロセッサ・ビジー要求が活性である。
- 転送が完了した。
- データ読出し転送動作およびデータが利用可能でない。
- CLSプリミティブが受信され、かつ、それ以上のBB__Creditが利用可能でない。
- マイクロプロセッサ中断要求が待機中であり、かつ、ロジックがフレーム間にある。

30

40

【0093】

図5において、ループ・ポートA/B開放割込み制御状態・マシン46(ポートA), 46'(ポートB)は、ループ1250がOpen割込み状態にあるときの場合を取り扱う。このロジック46, 46'は、フレームを送信する要求を発生する。各ポート(それぞれ46, 46')に対して1つの状態・マシンがある。これらの状態・マシンは、マイクロプロセッサ112がフレームを要求するときに、フレームを送信する要求を発生し、かつ、EOFの送信を監視する。送信が完了すると、送信完了がマイクロプロセッサ112に対して発生される。

50

【0094】

ブロック40への入力は、PORT_BB_CREDIT_AVAILABLE_TO_TRANSMIT_RDY6017と、PORT_CREDIT_AVAILABLE_TO_TRANSMIT_A_FRAME6020と、LOOP_A_STATES_AND_CONTROL6422と、LOOP_B_STATES_AND_CONTROL6432（図4を参照）と、DATA_AVAILABLE6019とを含む。ブロック40からの出力は、TRANSMIT_CONTROL_OUTPUTS6413と、LOOP_A_TRANSMIT_CONTROL_OUTPUTS6425と、LOOP_B_TRANSMIT_CONTROL_OUTPUTS6427とを含む。

【0095】

III. フレーム受信用の専用フレーム・バッファ

二重ポート・ファイバ・チャネル・アービトレーテッド・ループ設計1200では、オン・チップ・フレーム・バッファ119のバッファを使用して、入フレームおよび出フレームを管理することができる。受信されたフレームおよび送信されたフレームは、通常、より低い転送速度で大きなオフ・チップ領域（例えば、オフ・チップ・バッファ111）に記憶される。オフ・チップ・バッファ111が単一ポートに対して全転送速度が可能なときであっても、二重ポート設計に対しては、必要とされる帯域幅は遙かに大きく、追加コストが加わる。FC-AL ASIC110（図1を参照）におけるオン・チップ・フレーム・バッファ119は、パフォーマンス、シリコン面積およびコスト間のバランスを考慮して種々の方法で構成される。この仕様は、各ポートに同時に非データ形式フレームを受信するとともに、大きな専用データ・フレーム・バッファ55（これもオン・チップ・フレーム・バッファ119全体の構成要素）を提供するように、専用フレーム・バッファ53, 53'（オン・チップ・フレーム・バッファ119全体の構成要素）の使用を詳述している。本発明による二重ポート設計では、複数のフレームを同時に両方のポート116で受信することができる。これらのフレームは、通常、受信された後に大きなオフ・チップ・メモリ111に移動されて記憶される。更なる情報は、「ループ初期化および応答の方法および専用フレーム・バッファ」と題する米国特許出願第09/193,387号に見い出すことができる。

【0096】

IV. オン・チップ・メモリでデータ保全用にファイバ・チャネルCRCを使用すること
本発明の一態様によれば、ファイバ・チャネル・フレームを一時的に記憶するフレーム・バッファは、最大ファイバ・チャネル・インターフェース・データ転送速度でフレームを受信することを可能とする。その後、より低くより管理し易い速度でオフ・チップ・メモリにフレームを転送することができる。パリティ、CRCまたは他の冗長機能のような種々のメカニズムは、オプションとして、フレーム・バッファにデータを記憶している間にデータを保護するために使用される。

【0097】

一実施の形態では、受信したファイバ・チャネル・サイクリック・リダンダンシー・コード（CRC）をデータとともにフレーム・バッファ通過させることによってデータ保全チェックが強化され（すなわち、CRCは、フレームとともにフレーム・バッファに記憶されて、後の時点でフレームとともに読み出される）、RAMをより幅広にする付加的なパリティ・ビットを省略することができる。

【0098】

V. アービトレーテッド・ループ・オーバヘッドを軽減する方法および装置

ファイバ・チャネル・アービトレーテッド・ループ設計1200では、ループ・ポート116のノード・インターフェース1220は、ループ1250にアクセスするためにアービトレーションする必要がある。どのポートがループ1250の制御を獲得するのかを決定するために優先度システムが使用され、また、ポートが足りなくないことを保証するために「フェアネス」方式が使用される。目標装置として、ドライブ100は、通常、CPU情報処理システム1202よりも低い優先度が与えられ、その結果、ドライブ100は

10

20

30

40

50

、複数のより高い優先度の装置がそれらのアクセスを完了するまで、アービトレーションに勝つのを待たねばならないかもしれない。ループ・ポート 1 1 6 のノード・インターフェース 1 2 2 0 がループ 1 2 5 0 の制御を獲得すると、それは、不必要なアービトレーション・サイクルを避けるために、ループ 1 2 5 0 を閉じる前に、可能な限り多くのフレームを送出する。しかし、データがもはや得られないときは、ループ・ポート 1 1 6 のノード・インターフェース 1 2 2 0 は、ループ 1 2 5 0 を閉じて、他のポートがループ 1 2 5 0 にアクセスすることを許容する。これは、ある他のコントローラ・アーキテクチャに使用される方法である。本発明は、ポートのデータ利用可能性に基づいて、ループ 1 2 5 0 を閉じるのか否かの判断に関するルールを変更することによってループ・パフォーマンスを高め、したがって総合ループ・オーバーヘッドを軽減する機構を提供する。

10

【 0 0 9 9 】

ある他のコントローラ・アーキテクチャでは、エンド・オブ・フレーム・デリミタが送信されているとき、ポートは、他のフレームが利用可能か否かを判断する。データがもはや得られないときは（例えば、全フレームが送信用に利用できないときは）、ループ 1 2 5 0 が閉じられる。その直ぐ後にデータが再び利用可能となるので、ポートは、後で再びアービトレーションし、転送を継続する前にアービトレーションに勝つ必要がある。転送の最終フレームが利用可能となるときにこれが発生すれば、転送の完了は遅延され、次のコマンドが可能となる前に遅延をもたらすかもしれない。本発明は、データがポートに対して直ぐに利用可能となるのであれば、あるポートによってループ 1 2 5 0 を開放のままにすることができるコントローラ・アーキテクチャ設計用の機構を提供する。一実施の形態では、次の条件の両方が満足されるときは、更なるデータがポートに利用可能となるのが十分に予測される際は（ポートがそのループの制御を保持するのを正当化するために）、ループ 1 2 5 0 を開放したままとする。

20

- 少なくとも X フレームが利用可能なオフ・チップであり、かつ

- データの少なくとも Y ワードがデータ・フレーム・バッファ 5 5 で利用可能である。

一実施の形態では、所定の量のデータが利用可能（オン・チップでは少なくとも 1 / 2 フレームが利用可能で、オフ・チップでは少なくとも 1 フレームが利用可能）であれば、ループ 1 2 5 0 は開に保持されるが、フレームの転送は、1 つのフレーム全部がオン・チップで利用可能になるまで、開始しないであろう。

【 0 1 0 0 】

30

本発明の 1 つの目的は、データがポート 1 1 6 に対して直ぐに利用可能であるときは、ループ 1 2 5 を開に保持して、付加アービトレーション・サイクルを回避することである。延長された期間（例えば、ヘッド切換えを実行するのに必要な時間）の間待機となるのであれば、データが利用可能となるまで待機してループ 1 2 5 0 を開のままにしておくはならない。なぜならば、これは、ループ 1 2 5 0 上で他のポートが転送を実行するのを阻止するからである。

【 0 1 0 1 】

（結論）

以上、1 つ以上のポート（1 1 6）を有する第 1 のチャネル・ノード（1 2 2 0）を含むループ・フェアネスを保持する改良された通信チャネル・システム（1 2 0 0）について説明した。各ポート（1 1 6）は、ファイバ・チャネル・アービトレーテッド・ループ・シリアル・通信チャネル（1 2 5 0）をサポートし、かつ、これに取り付けられる。これらのポートの 1 つは、そのポートが取り付けられたチャネル（1 2 5 0）の制御のためにアービトレーションするのである。ここで、チャネル・ループ（1 2 5 0）の制御は、アービトレーションに勝つと、第 1 のポートと通信チャネル（1 2 5 0）との間で所定量の使用が発生したか否かに少なくとも一部基づいてフェアネス保持装置（1 7 5）は通信チャネルの制御を解放させる。

40

【 0 1 0 2 】

一実施の形態では、所定量の使用は、通信チャネルの制御を所定の時間保持することを含む。

50

【 0 1 0 3 】

他の実施の形態では、所定量の使用は、第 1 の所定量のデータを転送することを含む。そのような一実施の形態では、システム (1 2 0 0) は、転送されたデータの量を監視する第 1 のカウンタ (1 6 2) をさらに含む。第 1 のカウンタ (1 6 0) に作動的に結合された第 1 のコンパレータ回路 (1 6 2) は、第 1 のカウンタ (1 6 0) によって監視されたデータの量が第 1 の所定量のデータに達したか否かに少なくとも一部基づいて通信チャネルの制御を解放させる。一実施の形態では、第 1 の所定量のデータはレジスタ (1 6 3) に保持される。

【 0 1 0 4 】

いくつかの実施の形態では、システム (1 2 0 0) は、転送されるべく残っているデータの量を監視する第 2 のカウンタ (1 6 5) をさらに含む。第 2 のコンパレータ回路 (1 6 7) は、第 2 のカウンタ (1 6 5) に作動的に結合され、第 2 のカウンタ (1 6 5) によって監視されたデータの量が第 2 の所定量のデータよりも小さいか否かに少なくとも一部基づいて通信チャネルの制御の解放を禁止する。一実施の形態では、第 2 の所定量のデータはレジスタ (1 6 8) に保持される。

10

【 0 1 0 5 】

いくつかの実施の形態では、システム (1 2 0 0) はダイナミック半二重をサポートし、ここで、第 1 のカウンタ (1 6 0) および第 1 のコンパレータ (1 6 2) はダイナミック半二重コマンドの第 1 の部分による受信時に初期化される。

【 0 1 0 6 】

いくつかの実施の形態では、データの第 1 の所定量およびデータの第 2 の所定量は、プログラム可能な量である。

20

【 0 1 0 7 】

システム (1 2 0 0) のいくつかの実施の形態は、第 1 のチャネル・ノード (1 2 2 0) に作動的に結合された磁気ディスク記憶装置 (1 1 4) をさらに含む。コンピュータ・システム (1 2 0 2) は第 2 のチャネル・ノード (1 2 2 0) を含み、ここで、第 2 のチャネル・ノード (1 2 2 0) は、ファイバ・チャネル・アービトレーテッド・ループ・シリアル通信チャネルを介して第 1 および第 2 のチャネル・ノード間でデータを転送するために、ファイバ・チャネル・ループ (1 2 2 0) で第 1 のチャネル・ノード (1 2 2 0) に作動的に結合されている。

30

【 0 1 0 8 】

本発明の他の態様は、(a) ファイバ・チャネル・アービトレーテッド・ループ・シリアル通信チャネルのループの制御のためにアービトレーションすることと、(b) 第 1 のポートと通信チャネルとの間で所定量の使用が発生したか否かに少なくとも一部基づいて通信チャネルの制御を解放させることとを含む。

【 0 1 0 9 】

方法の一実施の形態では、解放するステップ (b) は、(b) (i) 通信チャネルの制御が所定の時間維持されたか否かを判断するステップと、(b) (i i) 判断するステップ (b) (i) に基づいてループの制御を解放するステップとをさらに含む。

【 0 1 1 0 】

方法の他の実施の形態では、解放するステップ (b) は、(b) (i i i) 第 1 の所定量のデータが転送されたか否かを判断するステップと、(b) (i v) 判断するステップ (b) (i i i) に基づいてループの制御を解放するステップとをさらに含む。そのような一実施の形態では、判断するステップ (b) (i i i) は、(b) (i i i) (A) 第 1 の値を得るために転送されたフレームの数を監視するステップと、(b) (i i i) (B) 第 1 の値を第 1 の所定量のデータと比較するステップとをさらに含む。

40

【 0 1 1 1 】

方法の他の実施の形態では、解放するステップ (b) は、(b) (v) 第 2 の所定量のデータが転送されたか否かを判断するステップと、(b) (v i) 判断するステップ (b) (v) に基づいてループの制御の解放を禁止するステップとをさらに含む。

50

【 0 1 1 2 】

上述した方法のいくつかの実施の形態は、(c) ダイナミック半二重コマンドを受信するステップと、(d) 判断するステップ(b) (i i i) をダイナミック半二重コマンドの受信時に初期化するステップとをさらに含む。これらの実施の形態のいくつかでは、初期化するステップ(d) は、(d) (i) 第 1 の所定量のデータおよび第 2 の所定量のデータをプログラム可能にセットするステップをさらに含む。

【 0 1 1 3 】

本発明の別の他の態様は、ループ・フェアネスを保持するファイバ・チャネル・ノード・コントローラ・システムを提供する。このシステムは、ファイバ・チャネル・アービトレートッド・ループ・シリアル通信チャンネル(1 2 5 0) と、ファイバ・チャネル・アービトレートッド・ループ・シリアル通信チャンネル(1 2 5 0) に取り付けられた第 1 のポート(1 1 6) を有する第 1 のチャネル・ノード(1 2 2 0) であって、第 1 のポートが、そのポートが取り付けられた通信チャンネルの制御のためにアービトレーションする、第 1 のチャネル・ノード(1 2 2 0) と、所定量の使用が第 1 のポートと通信チャンネルとの間で発生したか否かに少なくとも一部基づいて通信チャンネルの制御を解放させる、ここで説明したようなフェアネス保持手段とを含む。

10

【 0 1 1 4 】

上述の説明は例示を意図しかつ限定的でないことを理解すべきである。本発明の種々の実施の形態の多数の特徴および効果が種々の実施の形態の構造および機能の詳細とともに以上の説明で明らかにされたが、詳細に対する多くの他の実施の形態および変更は、以上の説明を再検討すると当該技術分野に習熟する者に明らかである。したがって、本発明の範囲は、付記する請求の範囲をこのような請求の範囲が主張する等価物の完全な範囲とともに参照して、判断すべきである。

20

【 図面の簡単な説明 】

【 図 1 】 ファイバ・チャネル・ノード・インターフェースを有するディスク・ドライブ 1 0 0 のブロック図である。

【 図 2 】 本発明を組み入れている情報処理システム 1 2 0 0 のブロック図である。

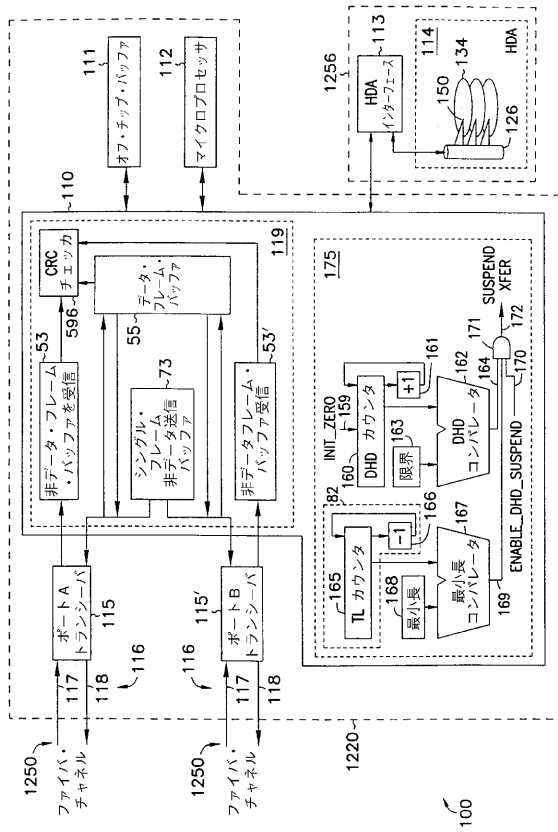
【 図 3 】 ファイバ・チャネル・ノード・インターフェース・チップ 1 1 0 のブロック図である。

【 図 4 】 ファイバ・チャネル・ループ・ポート回路 2 0 のブロック図である。

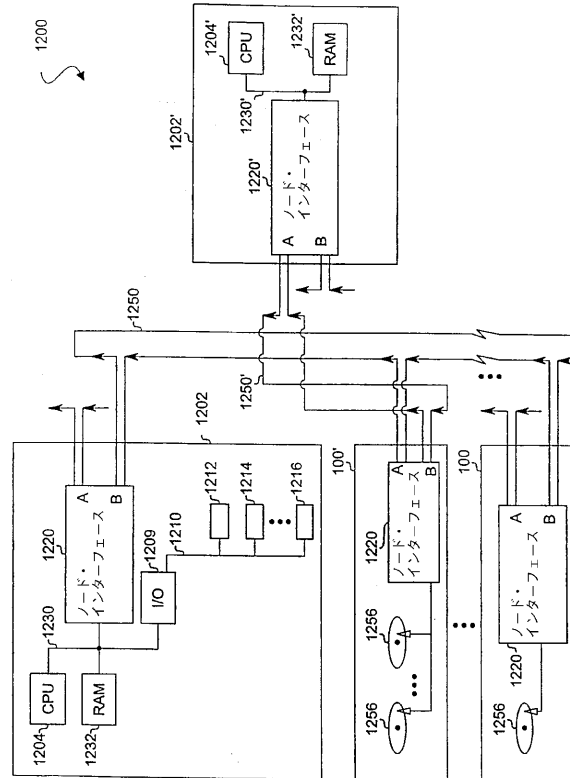
30

【 図 5 】 ファイバ・チャネル・ループ制御回路 4 0 のブロック図である。

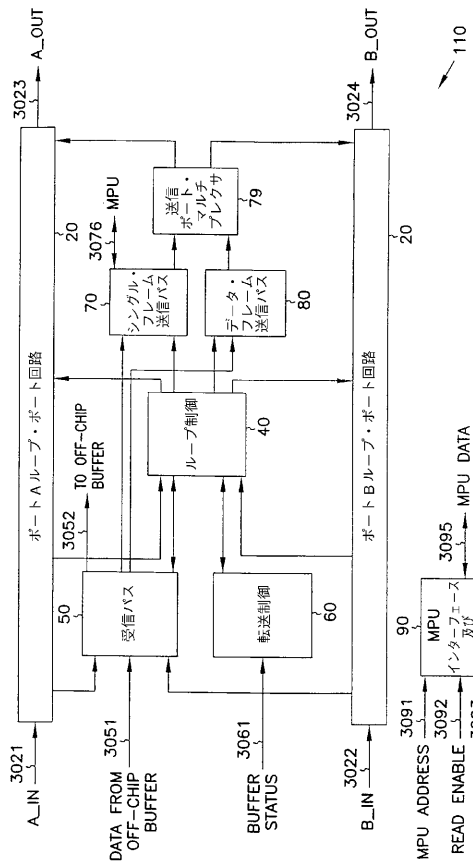
【 図 1 】



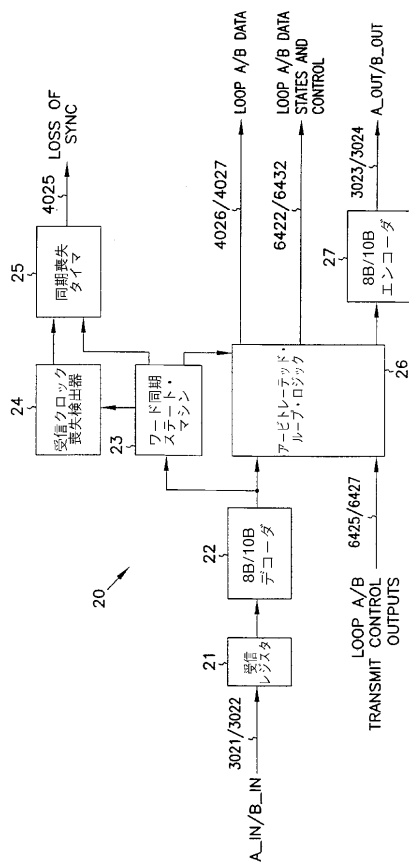
【 図 2 】



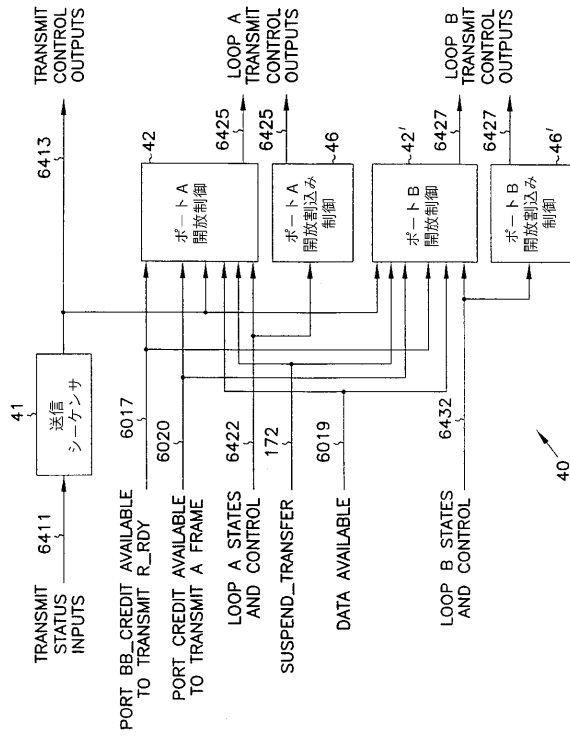
【 図 3 】



【 図 4 】



【図 5】



フロントページの続き

(72)発明者 ミラー、マイケル、エイチ

アメリカ合衆国 ミネソタ、イーデン プレイリー、 パーク ビュー レーン 6 8 5 0

(72)発明者 ウエストビー、ジュディ、リン

アメリカ合衆国 ミネソタ、ブルーミントン、ダブリュ、 ワンハンドレッド アンド イレブン
ス ストリート 7 9 0 6

審査官 矢頭 尚之

(58)調査した分野(Int.Cl. , D B名)

H04L 12/42