

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4211282号
(P4211282)

(45) 発行日 平成21年1月21日(2009.1.21)

(24) 登録日 平成20年11月7日(2008.11.7)

(51) Int. Cl.	F I
G06F 13/00 (2006.01)	G06F 13/00 520B
G06F 12/00 (2006.01)	G06F 13/00 520C
	G06F 13/00 520R
	G06F 12/00 545M

請求項の数 22 (全 17 頁)

(21) 出願番号	特願2002-139232 (P2002-139232)	(73) 特許権者	000002185
(22) 出願日	平成14年5月14日(2002.5.14)		ソニー株式会社
(65) 公開番号	特開2003-330619 (P2003-330619A)		東京都港区港南1丁目7番1号
(43) 公開日	平成15年11月21日(2003.11.21)	(74) 代理人	100067736
審査請求日	平成17年4月13日(2005.4.13)		弁理士 小池 晃
前置審査		(74) 代理人	100096677
			弁理士 伊賀 誠司
		(74) 代理人	100106781
			弁理士 藤井 稔也
		(74) 代理人	100113424
			弁理士 野口 信博
		(74) 代理人	100150898
			弁理士 祐成 篤哉

最終頁に続く

(54) 【発明の名称】 データ蓄積方法及びデータ蓄積システム、並びに、データ記録制御装置、データ記録指令装置、データ受信装置及び情報処理端末

(57) 【特許請求の範囲】

【請求項1】

ネットワークを介して接続された複数（n個）の情報処理端末にデータを分散させて記録するデータ蓄積方法において、

上記複数の情報処理端末にデータを分散して記録する記録指令装置が、入力された第1のデータをp個のブロックに分割し符号化率q/pでFEC（Forward Error Collection）符号化してq個の符号化ブロックを有する第2のデータに変換する符号化工程と、

上記記録指令装置が、上記第2のデータにおけるq個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化工程と、

上記記録指令装置が、上記パケット化工程において生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成工程と、

上記記録指令装置が、上記パケットを上記複数の情報処理端末に送信する第1の送信工程と、

上記情報処理端末が、上記記録確率で該パケットを記録する記録工程とを有し、

上記記録確率で上記複数の情報処理端末の何れかに記録された上記パケットを送信指令装置が読み出す処理として、

上記送信指令装置が、上記情報処理端末が返信確率で上記パケットを返信することを要求する返信要求パケットを生成する返信要求パケット生成工程と、

上記送信指令装置が、上記返信要求パケットを上記複数の情報処理端末に送信する第2の送信工程と、

10

20

上記複数の情報処理端末が、上記返信確率 で該パケットを返信する返信工程とを有し

、
 上記符号化率 q/p 、記録確率 、返信確率 、および上記情報処理端末の個数 n を、
 $p \times q \times \dots \times n \times$ を満たすように設定する
 データ蓄積方法。

【請求項 2】

上記記録確率 は、上記パケットの上記制御情報に含まれ、該パケットにより上記複数の
 情報処理端末に送信され、

上記返信確率 は、上記返信要求パケットの制御情報に含まれ、該返信要求パケットに
 より上記複数の情報処理端末に送信される請求項 1 記載のデータ蓄積方法。

10

【請求項 3】

上記記録確率 は、上記パケットとは別のパケットで送信され、

上記返信確率 は、上記返信要求パケットとは別のパケットで送信される請求項 1 記載
 のデータ蓄積方法。

【請求項 4】

上記返信工程において返信されるパケット数が上記第 1 のデータのブロック数 p と同じ
 かそれ以上である請求項 1 記載のデータ蓄積方法。

【請求項 5】

上記符号化率 q/p 、上記記録確率 、上記返信確率 を変更し上記パケットの分散を
 制御するパケット分散制御工程を有する請求項 4 記載のデータ蓄積方法。

20

【請求項 6】

上記パケット分散制御工程では、上記記録確率 と、上記返信確率 と、上記情報処理
 端末の個数 n と、第 2 のデータのブロック数 q との積が上記第 1 のデータのブロック数 p
 と同じかそれ以上になるように各値を変更する請求項 4 記載のデータ蓄積方法。

【請求項 7】

上記パケット分散制御工程では、上記ネットワークにおけるデータ損失確率 a と、上記
 記録確率 と、上記返信確率 と、上記情報処理端末の個数 n と、第 2 のデータのブロッ
 ク数 q との積が上記第 1 のデータのブロック数 p と同じかそれ以上になるように各値を変
 更する請求項 6 記載のデータ蓄積方法。

【請求項 8】

30

上記パケット分散制御工程では、上記複数の情報端末が応答しない確率 b と、上記記録
 確率 と、上記返信確率 と、上記情報処理端末の個数 n と、第 2 のデータのブロック数
 q との積が上記第 1 のデータのブロック数 p と同じかそれ以上になるように各値を変更す
 る請求項 4 記載のデータ蓄積方法。

【請求項 9】

上記記録確率生成工程および上記返信要求パケット生成工程は、

上記記録確率 、上記返信確率 、または上記符号化率 q/p を大きくすることで、上
 記返信工程において返信されるパケットが重複する確率を低くする請求項 1 記載のデー
 タ蓄積方法。

【請求項 10】

40

ネットワークを介して接続された複数 (n 個) の情報処理端末にデータを分散させて記
 録するデータ蓄積システムにおいて、

入力された第 1 のデータを p 個のブロックに分割し符号化率 q/p で F E C (Forward
 Error Collection) 符号化して q 個の符号化ブロックを有する第 2 のデータに変換する符
 号化手段と、

上記第 2 のデータにおける q 個の符号化ブロックの各々に制御情報を付加してパケット
 化するパケット化手段と、

上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する
 確率を表した記録確率 を生成する記録確率生成手段と、

上記パケットを上記複数の情報処理端末に対して送出する送信手段と

50

を有するデータ記録指令装置と、
 上記データ記録指令装置との間で上記パケットを送受する送受信手段と、
 上記パケットを記録する記録手段と、
 上記受信したパケットを上記記録確率 で上記記録手段に記録し、記録したパケットを
 返信確率 で返信するように制御する制御手段と
 を有する複数の情報処理端末と、
 上記情報処理端末との間でパケットを送受する送受信手段と、
 上記情報処理端末の上記記録手段に記録されたパケットの返信を要求する返信要求パケ
 ットを生成する返信要求パケット生成手段と
 を有する送信指令装置と
 を備え、
 これらのデータ記録指令装置、情報処理端末、及び送信指令装置がネットワークにより
 互いに接続されているデータ蓄積システム。

10

【請求項 1 1】

上記記録確率 は、上記パケットの上記制御情報に含まれ、該パケットにより上記複数の
 情報処理端末に送信され、

上記返信確率 は、上記返信要求パケットの制御情報に含まれ、該返信要求パケットに
 より上記複数の情報処理端末に送信される請求項 1 0 記載のデータ蓄積システム。

【請求項 1 2】

上記記録確率 は、上記パケットとは別のパケットで送信され、

上記返信確率 は、上記返信要求パケットとは別のパケットで送信される請求項 1 0 記
 載のデータ蓄積システム。

20

【請求項 1 3】

上記データ記録制御装置における受信手段にて受信されるパケット数が上記第 1 のデー
 タのブロック数 p と同じかそれ以上である請求項 1 0 記載のデータ蓄積システム。

【請求項 1 4】

上記符号化率 q / p 、上記記録確率 、上記返信確率 を変更し上記パケットの分散を
 制御するパケット分散制御手段を有する請求項 1 3 記載のデータ蓄積システム。

【請求項 1 5】

上記パケット分散制御手段は、上記記録確率 と、上記返信確率 と、上記情報処理端
 末の個数 n と、第 2 のデータのブロック数 q との積が上記第 1 のデータのブロック数 p と
 同じかそれ以上になるように各値を変更する請求項 1 3 記載のデータ蓄積システム。

30

【請求項 1 6】

上記パケット分散制御手段は、上記ネットワークにおけるデータ損失確率 a と、上記記
 録確率 と、上記返信確率 と、上記情報処理端末の個数 n と、第 2 のデータのブロック
 数 q との積が上記第 1 のデータのブロック数 p と同じかそれ以上になるように各値を変更
 する請求項 1 5 記載のデータ蓄積システム。

【請求項 1 7】

上記パケット分散制御手段は、上記複数の情報端末が応答しない確率 b と、上記記録確
 率 と、上記返信確率 と、上記情報処理端末の個数 n と、第 2 のデータのブロック数 q
 との積が上記第 1 のデータのブロック数 p と同じかそれ以上になるように各値を変更する
 請求項 1 3 記載のデータ蓄積システム。

40

【請求項 1 8】

上記記録確率生成手段および上記返信要求パケット生成手段は、
 上記記録確率 、上記返信確率 、または上記符号化率 q / p を大きくすることで、上
 記複数の情報処理端末が返信するパケットが重複する確率を低くする請求項 1 0 記載のデ
 ータ蓄積システム。

【請求項 1 9】

ネットワークを介して接続された複数 (n 個) の情報処理端末にデータを分散して記録
 するデータ記録制御装置において、

50

入力された第1のデータをp個のブロックに分割し符号化率 q/p でFEC (Forward Error Collection) 符号化してq個の符号化ブロックを有する第2のデータに変換する符号化手段と、

上記第2のデータにおけるq個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化手段と、

上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成手段と、

上記情報処理端末に記録されたパケットを上記情報処理端末が返信する確率を示す返信確率を含む制御情報が付加された返信要求パケットを生成する返信要求パケット生成手段と、

10

上記パケット及び上記返信要求パケットを上記複数の情報処理端末に対して送出する送信手段と、

上記複数の情報処理端末から送られるパケットを受信する受信手段と
を備えるデータ記録制御装置。

【請求項20】

ネットワークを介して接続された複数(n個)の情報処理端末にデータを分散して記録するデータ記録指令装置において、

入力された第1のデータをp個のブロックに分割し符号化率 q/p でFEC (Forward Error Collection) 符号化してq個の符号化ブロックを有する第2のデータに変換する符号化手段と、

20

上記第2のデータにおけるq個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化手段と、

上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成手段と、

上記パケットを上記複数の情報処理端末に対して送出する送信手段と
を備えるデータ記録指令装置。

【請求項21】

ネットワークを介して接続された複数(n個)の情報処理端末からデータを受信するデータ受信装置において、

上記情報処理端末に記録されたパケットを上記情報処理端末が返信する確率を示す返信確率を含む制御情報が付加された返信要求パケットを生成する返信要求パケット生成手段と、

30

上記返信要求パケットを上記複数の情報処理端末に対して送出する送信手段と、
上記複数の情報処理端末から送られるパケットを受信する受信手段と
を備えるデータ受信装置。

【請求項22】

ネットワークで接続された外部装置との間でパケット化されたデータを送受する送受信手段と、

上記パケットを記録する記録手段と、

上記パケットを該パケットの付加情報に記述された記録確率で記録するか否か、および記録したパケットを返信確率で返信するか否かを制御する制御手段と

40

を備える情報処理端末。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、データ蓄積方法及びデータ蓄積システム、並びにデータ記録制御装置、データ受信装置及び情報処理端末に関し、ネットワークを介して接続されたノード間でデータを分散して蓄積するデータ蓄積方法及びデータ蓄積システムに関する。また、このデータ蓄積システムに用いて好適なデータ記録制御装置、データ記録指令装置、データ受信装置及び情報処理端末に関する。

50

【 0 0 0 2 】

【 従来 の 技 術 】

近年、あるデータを、ネットワークを介して互いに接続された多数の情報処理端末に分散して記録する大規模ストレージシステムが注目されている。このような分散型ストレージシステムにおいて、データを記録管理するサーバは、マルチキャスト等によって情報処理端末やその他のサーバにデータを送信して、データを情報処理端末や他のサーバに備えられたローカルの記録媒体に記録している。

【 0 0 0 3 】

この場合、オンデマンドでデータを取り出せるようにするためには、記録媒体に多量のデータを記録しなければならない。例えば、1本当たり約2ギガバイトのデータ容量になる映画の場合、このような映像データを500本分記録するとすれば、1テラバイト以上の容量が必要となる。

10

【 0 0 0 4 】

また、ストリーミングによってデータを提供する場合の例として、サーバがデータを要求しているクライアントに対してユニキャストでデータを提供する際には、エラーのない伝送を行うために、例えば、TCP/IPの到着済み信号(ACK)のようにデータの再送を要求するプロトコルが用いられる。

【 0 0 0 5 】

ところが、この手法は、サーバ側に多大な負担がかかるため、高性能なサーバ1台を用いたとしても、現状では、数百台のクライアントにしかサービスを提供することができない。また、UDP/IPのようなACKを用いないプロトコルを使用したとしても、サービス可能なクライアントの数は、数千台程度である。このように、ストリーミングによってデータを提供しようとする、サーバ側のコストが増大し、クライアントの数が制限されてしまう。

20

【 0 0 0 6 】

そこで近年では、マルチキャスト技術にFEC(Forward Error Correction)を用いて、データの再送を要求することなく複数のクライアントにデータを送信する方式が提案されている。これは、サーバがマルチキャストでストリームを繰り返し送信し、クライアントは、このストリームから必要な信号を拾い上げ、拾い上げたデータを復号して再生する方式である。

30

【 0 0 0 7 】

【 発 明 が 解 決 し よ う と す る 課 題 】

この方式を利用して、1本2ギガバイトになる映画の映像データ500本分を10分以内に送信する場合には、約14.7ギガビット/秒の伝送帯域が必要である。さらに、同量の映像データを1分以内に送信する場合には、約147ギガビット/秒の伝送帯域が必要になる。これは、理論値であるが、このような容量及び伝送方式に耐えうるサーバは、非常にコストがかかり、実現したとしても実用的でない。また、複数のホストにデータを分散して記録するという方式もあるが、このシステムを実現しようすると、巨大なデータを複数のサーバで管理しなければならないため、データ管理やデータ通信のための処理が増大してしまう。

40

【 0 0 0 8 】

そこで本発明は、上述したような従来の実情に鑑みて提案されたものであり、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して管理できるデータ蓄積方法及びデータ蓄積システム、並びに、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録するとともに記録場所からデータを取り出すよう制御するデータ記録制御装置、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録するデータ記録指令装置、個々の端末に分散して記録されたデータを受信するデータ受信装置及び多量のデータを分散して記録できる情報処理端末を提供することを目的とする。

50

【 0 0 0 9 】

【課題を解決するための手段】

上述した目的を達成するために、本発明に係るデータ蓄積方法は、ネットワークを介して接続された複数（ n 個）の情報処理端末にデータを分散させて記録するデータ蓄積方法において、上記複数の情報処理端末にデータを分散して記録する記録指令装置が、入力された第1のデータを p 個のブロックに分割し符号化率 q/p でFEC（Forward Error Collection）符号化して q 個の符号化ブロックを有する第2のデータに変換する符号化工程と、上記記録指令装置が、上記第2のデータにおける q 個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化工程と、上記記録指令装置が、上記パケット化工程において生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成工程と、上記記録指令装置が、上記パケットを上記複数の情報処理端末に送信する第1の送信工程と、上記情報処理端末が、上記記録確率で該パケットを記録する記録工程とを有し、上記記録確率で上記複数の情報処理端末の何れかに記録された上記パケットを送信指令装置が読み出す処理として、上記送信指令装置が、上記情報処理端末が返信確率で上記パケットを返信することを要求する返信要求パケットを生成する返信要求パケット生成工程と、上記送信指令装置が、上記返信要求パケットを上記複数の情報処理端末に送信する第2の送信工程と、上記複数の情報処理端末が、上記返信確率で該パケットを返信する返信工程とを有し、上記符号化率 q/p 、記録確率、返信確率、および上記情報処理端末の個数 n を、 $p \times q \times n$ を満たすように設定する。

10

20

【 0 0 1 0 】

また、本発明に係るデータ蓄積システムは、ネットワークを介して接続された複数（ n 個）の情報処理端末にデータを分散させて記録するデータ蓄積システムにおいて、入力された第1のデータを p 個のブロックに分割し符号化率 q/p でFEC（Forward Error Collection）符号化して q 個の符号化ブロックを有する第2のデータに変換する符号化手段と、上記第2のデータにおける q 個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化手段と、上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成手段と、上記パケットを上記複数の情報処理端末に対して送出する送信手段とを有するデータ記録指令装置と、上記データ記録指令装置との間で上記パケットを送受する送受信手段と、上記パケットを記録する記録手段と、上記受信したパケットを上記記録確率で上記記録手段に記録し、記録したパケットを返信確率で返信するように制御する制御手段とを有する複数の情報処理端末と、上記情報処理端末との間でパケットを送受する送受信手段と、上記情報処理端末の上記記録手段に記録されたパケットの返信を要求する返信要求パケットを生成する返信要求パケット生成手段とを有する送信指令装置とを備え、これらのデータ記録指令装置、情報処理端末、及び送信指令装置がネットワークにより互いに接続されている。

30

【 0 0 1 1 】

また、上述した目的を達成するために、本発明に係るデータ記録制御装置は、ネットワークを介して接続された複数（ n 個）の情報処理端末にデータを分散して記録するデータ記録制御装置において、入力された第1のデータを p 個のブロックに分割し符号化率 q/p でFEC（Forward Error Collection）符号化して q 個の符号化ブロックを有する第2のデータに変換する符号化手段と、上記第2のデータにおける q 個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化手段と、上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成手段と、上記情報処理端末に記録されたパケットを上記情報処理端末が返信する確率を示す返信確率を含む制御情報が付加された返信要求パケットを生成する返信要求パケット生成手段と、上記パケット及び上記返信要求パケットを上記複数の情報処理端末に対して送出する送信手段と、上記複数の情報処理端末から送られるパケットを受信する受信手段とを備える。

40

50

【0012】

また、上述した目的を達成するために、本発明に係るデータ記録指令装置は、ネットワークを介して接続された複数（ n 個）の情報処理端末にデータを分散して記録するデータ記録指令装置において、入力された第1のデータを p 個のブロックに分割し符号化率 q/p でFEC（Forward Error Collection）符号化して q 個の符号化ブロックを有する第2のデータに変換する符号化手段と、上記第2のデータにおける q 個の符号化ブロックの各々に制御情報を付加してパケット化するパケット化手段と、上記パケット化手段によって生成されたパケットを上記複数の情報処理端末で記録する確率を表した記録確率を生成する記録確率生成手段と、上記パケットを上記複数の情報処理端末に対して送出する送信手段とを備える。

10

【0013】

また、本発明に係るデータ受信装置は、ネットワークを介して接続された複数（ n 個）の情報処理端末からデータを受信するデータ受信装置において、複数の情報処理端末に記録されたパケットを返信する確率を示す返信確率を含む制御情報が付加された返信要求パケットを生成する返信要求パケット生成手段と、返信要求パケットを複数の情報処理端末に対して送出する送信手段と、複数の情報処理端末から送られるパケットを受信する受信手段とを備える。

【0014】

また、本発明に係る情報処理端末は、ネットワークで接続された外部装置との間でパケット化されたデータを送受する送受信手段と、パケットを記録する記録手段と、パケットを該パケットの付加情報に記述された記録確率で記録するか否か、および記録したパケットを返信確率で返信するか否かを制御する制御手段とを備える。

20

【0015】

【発明の実施の形態】

以下、本発明を適応した具体例について図を参照して詳細に説明する。図1は、ネットワークを構成する端末にデータを分散して記録する分散型ストレージシステム的具体例を示す図である。分散型ストレージシステムは、網の目のように接続された n 個のノード 20_1 、 20_2 、...、 20_{n-1} 、 20_n と、各ノードへのデータの記録を指示制御する記録指令装置10と、各ノードに記録されたデータを読み出す送信指令装置30とを有する。本具体例では、記録指令装置10と、ノード20と、送信指令装置30とを別の装置として説明するが、これら両装置の機能を有する記録送信指令装置のような装置であってもよい。また、ノードの各々に記録指令装置10や送信指令装置30としての機能を装備することもできる。この場合、ネットワークを構成する各装置を区別な使用できる。

30

【0016】

また、図1では、ネットワークを構成するルータなどの伝送制御装置を省略してあるが、実際には、ノード20を通過するパケットの経路を選択するルータなどの伝送制御装置が設けられている。伝送制御装置は、ノード20は別に設けられていてもよいし、ノード20が伝送制御装置のような機能を有していてもよい。

【0017】

次に、記録指令装置10について説明する。図2は、記録指令装置10の内部構成を示す図である。記録指令装置10は、FEC（Forward Error Collection）符号化を施すFECエンコーダ11、符号化したデータをインターリーブするインターリーバ12、インターリーバ12から出力されたデータをパケットに変換するパケット生成部13、及びネットワークとの接続を行うネットワークインターフェース14を有しており、データ入力部15を介してデータが入力される。

40

【0018】

FECエンコーダ11は、FEC符号化を用いて入力されたデータを符号化する。ここで、FEC符号化とは、トルネード符号化方式、リードトルネード符号化方式、ターボ符号方式などの受信側で誤り訂正を行う符号化方式の総称であり、FECエンコーダ11は、データ入力部15から入力されたデータを p 個のブロックに分割し、この p 個のブロック

50

に F E C 符号化を施して q 個のブロックに変換している。この p 個のブロックから q 個のブロックに符号化することを符号化率 q / p の符号化といい、この符号化率 q / p を変更することによって、この分散型ストレージシステムの記録効率や伝送効率を変更することができる。

【 0 0 1 9 】

インターリーバ 1 2 は、符号化されたデータの順番を並び換える。インターリーバ 1 2 は、このようにインターリーブすることによってデータを分散させ、パケットの消失によって発生するバーストエラーがランダムエラーになるようにすり替える。

【 0 0 2 0 】

パケット生成部 1 3 は、インターリーバ 1 2 から送信されたデータを所定のサイズに分割し、この分割したデータに、制御情報を付加してパケット化する。図 3 は、パケット生成部 1 3 によって生成されるパケット 4 0 の構造を示す図である。パケット 4 0 は、ヘッダ 4 1、記録確率記述部 4 2、ペイロード 4 3、フッタ 4 4 から構成される。ペイロード 4 3 には、F E C エンコーダ 1 1 によって変換された送信情報が記述されている。ヘッダ 4 1 とフッタ 4 4 には、データの種類を示すデータ I D、C R C (Cyclic Redundancy Check) のチェックサム、ノード 2 0 の G U I D (Global Unique ID)、ネットワークアドレスなどの制御情報が記述されている。

10

【 0 0 2 1 】

記録確率記述部 4 2 には、後述する各ノードがこのパケットを記録する確率が記述されている。ノードの制御部 2 2 は、この記録確率に基づいてパケットを記録する。

20

【 0 0 2 2 】

分散型ストレージシステムを構成する全てのノードは、この記録確率に基づいてパケットを記録するか否かを決定する。これにより、分散型ストレージシステムを構成するノードに確率でデータが記録されることになる。この分散型ストレージシステムでは、ノードの数 n が十分に大きく、且つ復号されたブロック数 q が十分に大きい場合に、各ノードに均等な確率でデータを分散することができる。

【 0 0 2 3 】

ネットワークインターフェース 1 4 は、パケット化されたデータを入力し、ユニキャスト又はマルチキャストで各ノード 2 0 にデータを送信する。

【 0 0 2 4 】

次に、ノード 2 0 の構成について説明する。図 4 は、ノード 2 0 の構成を示す図である。ノード 2 0 は、ネットワークインターフェース 2 1、制御部 2 2、記録部 2 3 を有し、ネットワークインターフェース 2 1 を介してネットワークと接続している。

30

【 0 0 2 5 】

制御部 2 2 は、入力したパケット 4 0 を記録部 2 3 に記録する記録処理と、他のノード 2 0 からの要求に応じてデータを返信する返信処理を実行する。

【 0 0 2 6 】

記録処理は、入力したパケット 4 0 の記録確率をもとに、記録部 2 3 にパケット 4 0 を記録するか否かを判定する処理であり、上述したように、分散型ストレージシステムを構成するノードの数 n が十分に大きく、且つ復号化されたブロックの個数 q が十分に大きい場合、全てのノードにパケットが均等に記録され、分散型ストレージシステム全体としての確率でデータが記録される。

40

【 0 0 2 7 】

返信処理は、送信指令装置から送られたデータの送信要求をするパケットを受信すると開始する処理である。このデータ送信要求パケットのデータ構造については後述するが、このパケットには、データを返信する確率が記述されており、制御部 2 2 は、記録されたデータを返信確率で返信する。ここで、データを返信すると判定したとき、制御部 2 2 は、目的のデータを検索して返信パケットを生成する。

【 0 0 2 8 】

次に、送信指令装置 3 0 について説明する。図 5 は、送信指令装置 3 0 の構成を示す図で

50

ある。送信指令装置 30 は、ネットワークを介しての外部とのデータの送受信を行うネットワークインターフェース 31、ノードにデータの送信を要求するデータ要求部 32、パケットに分割されたデータを結合するパケット結合部 33、デインターリーブ 34、FEC デコーダ 35 を有し、復号したデータは応用処理部 36 に入力され、モニタやスピーカ（図示省略）などの外部機器に出力したり、図示しない記録装置に格納されたりする。

【0029】

データ要求部 32 は、分散型ストレージシステムを構成する各ノードにデータを要求するパケットを送信する。図 6 は、データを要求するパケット 50 の構成を示す図である。パケット 50 は、ヘッダ 51、返信確率記述部 52、リクエスト記述部 53、フッタ 54 から構成される。リクエスト記述部 53 には、要求するデータを識別するためのデータ ID が記録される。ヘッダ 51 とフッタ 54 には、CRC のチェックサム、ノードのネットワークアドレスや GUID、データの順序を示すシーケンス番号などの制御情報が記録される。

10

【0030】

返信確率記述部 52 には、返信確率 が記述されている。返信確率 は、パケット 50 を受信したノードがデータを返信するか否かの判定を行うための変数である。この変数をもとに、データを返信すると判定するノードもあれば、データを返信しないと判定するノードもあるが、返信確率 は、分散型ストレージシステム全体をマクロ的にみたときの値であり、分散型ストレージシステム全体では、各ノードのデータを返信する確率が になる。そのため、分散型ストレージシステムに n 個のノード 20 が存在する場合に、返信されるパケットの割合は、ノード 20 の個数 n と返信確率 を掛け合わせた値 $n \times$ となる。

20

【0031】

パケット結合部 33 は、各ノード 20 から返信されたパケットを結合する。図 7 は、ノードから返信されたパケット 60 のデータ構造を示す図である。図 7 に示すようにパケット 60 は、ヘッダ 61、ペイロード 62、フッタ 63 から構成され、ペイロード 62 には送信するデータが格納され、ヘッダ 61 とフッタ 63 には、CRC のチェックサムや受信側のノードのネットワークアドレス、パケットの順序を示すシーケンス番号など、制御情報が格納される。

【0032】

図 7 のパケット 60 を受信すると、パケット結合部 33 は、シーケンス番号を読み取り、受信したパケット 60 の順序を入れ替え、ヘッダ 61 やフッタ 63 などの制御情報を除去して、シーケンス番号の順にパケットを結合する。

30

【0033】

デインターリーブ 34 は、結合されたパケットにデインターリーブをかけ、データの並びを整列させる。FEC デコーダ 35 は、デインターリーブされたデータに FEC 復号を施し、元のデータを復元する。

【0034】

FEC デコーダ 35 によって復号されたデータは、応用処理部 36 に出力される。応用処理部 36 は、復号されたデータを、図示しない記録部に保存したり、モニタやスピーカなどの入出力インターフェースに出力する。

40

【0035】

このように、本具体例における分散型ストレージシステムは、記録確率 で各ノードにデータを記録し、各ノードに記録したデータを返信確率 で返信させるシステムであり、記録指令装置 10 から出力される元データは、 $\times n \times$ の割合で返信される。例えば、 p 個のブロックを q 個のブロックに符号化すると、 $q \times \times n \times$ 個のブロックが返信される。後述する論文 R I Z Z 0 9 7 に記載のように、この返信されたブロックの個数が復号前のブロックの個数 p よりも多い場合、データは復号可能である。そのため、返信されるブロックの個数が p 個より多くなるように、 \times 、 \times 、 q / p の値を決定しておくこと、目的のデータを復号することができる。

【0036】

50

次に、分散型ストレージシステムの動作について、図 8 ~ 図 12 にしたがって説明する。まず、データの記録動作について説明する。データを記録する際、記録指令装置 10 は、データ入力部 15 から目的のデータを入力し、入力したデータを F E C エンコーダ 11 に出力する。

【 0 0 3 7 】

F E C エンコーダ 11 は、図 8 (a) に示すように、入力したデータを p 個のブロックに分割する。そして、符号化率 q / p の F E C 符号化を施し、 p 個のブロックに分割したデータを q 個の符号化ブロックに変換する。

【 0 0 3 8 】

上述したように、F E C 符号化とは、トルネード符号化方式、リードトルネード符号化方式、ターボ符号化方式などの受信側で誤り訂正を行う符号化方式の総称であり、F E C 符号化を用いて、あるデータを符号化率 q / p で符号化した場合、論文 R I Z Z 0 9 7 (<http://www.iet.unipi.it/~luigi/fec.html#fec.ps>) に発表されているように、 p 個以上の符号化されたブロックが残存すれば、幾つかのブロックが消失しても、元のメッセージが復元できるようになっている。

【 0 0 3 9 】

F E C 符号化を施されたデータは、インターリーバ 12 に出力され、インターリーバ 12 は、符号化されたデータの順番を並び換え、データを分散させる。インターリーバ 12 を通過したデータは、パケット生成部 13 に出力される。パケット生成部 13 は、分散されたデータを所定の大きさに分割し、分割したデータにヘッダやフッタを付加して、図 8 (b) に示すようなパケット 40 を生成する。

【 0 0 4 0 】

このように、生成されたパケット 40 は、ネットワークインターフェース 14 に出力される。ネットワークインターフェース 14 は、図 9 に示すように、ユニキャスト若しくはマルチキャストを用いて、パケット化したデータを、分散型ストレージシステムを構成する各ノードに送信する。

【 0 0 4 1 】

記録指令装置 10 から送信されたパケット 40 を受信したノードは、記録処理を開始し、パケット 40 の記録確率 を読み出して、パケット 40 を記録する否かの判定を行う。そして、記録すると判定した場合には、データを記録部 22 にデータを記録し、記録しないと判定した場合には、受信したパケット 40 をその他のノードに出力する。記録確率 は、分散型ストレージシステム全体で、どの程度のデータを記録するかを示す値であり、送信データが q 個の符号化ブロックであるとき、分散型ストレージシステム全体には、 $q \times$ 個のデータが記録される。

【 0 0 4 2 】

次に、データを取り出す処理について説明する。データを取り出すとき、送信指令装置 30 は、パケット 50 を生成し、図 10 に示すように、ユニキャスト若しくはマルチキャストを用いて、各ノードにパケット 50 を出力する。

【 0 0 4 3 】

パケット 50 を受信したノードは、返信処理を開始する。パケット 50 には、通常の情報以外の他に、データを識別するデータ I D、パケットを返信する返信確率 などが記録されており、返信処理では、パケット 50 の返信確率 をもとにデータを返信するか否かの判定する。ここで、データを返信すると判定した場合、ノードの制御部 22 は、データを包含したパケット 60 を生成し、送信指令装置 30 に返信する。また、データを返信しないと判定した場合には、ノードの制御部 22 は、何も実行しない。

【 0 0 4 4 】

図 11 は、各ノードから送信指令装置 30 にパケットが送信される様子を示した図である。図 11 に示すように、分散型ストレージシステムには、データを返信するノードとそうでないノードがあるが、全体で見ると、各ノードは返信確率 の割合でデータを返信するようになっており、分散型ストレージシステムに n 個のノードが存在すると、 $n \times$ の割

10

20

30

40

50

合のデータが返信されることになる。上述したように、分散型ストレージシステムには、記録確率 r をもとに $q \times n$ のデータが記録されているので、 $q \times n \times r$ 個のデータが返信される。

【0045】

次に、図12を参照して、データの復号化について説明する。各ノードから送信されるデータを受信した送信指令装置30は、受信したデータをパケット結合部32に入力し、パケットの結合を行う。図12(a)は、ノードから送信指令装置30に返信されるパケット60を示す図である。図12(a)において、斜線で描かれた箇所がデータの損失された箇所を示す。分散型ストレージシステムでは、元のデータを n の確率で返信するため、図12(a)に示すように、全てのパケットが返信される訳ではなく、幾つかのパケットの消失が起きている。

10

【0046】

しかしながら、このデータはインターリーブされているため、デインターリーブをかけると、パースト的なデータの消失は、分散し、ランダムなデータの消失に変換される。

【0047】

上述したように、FEC符号化では、符号化率 q/p でデータを符号化した場合、符号化されたデータが p ブロック以上存在すれば、データを復号することができる。送信した q 個のブロックが n の割合で返信されたときは、返信されるデータの個数は、 $q \times n$ である。この個数 $q \times n$ が p 個以上であるように q/p 、 r の値を設定すれば、元のデータを復元できる。

20

【0048】

また、本具体例における分散型ストレージシステムは、 $p > q \times n$ を満たすように、符号化率 q/p 、記録確率 r 、返信確率 s を設定すればよいので、上述の式を満たす範囲で符号化率 q/p 、記録確率 r 、返信確率 s を変更することにより、データの記録効率や伝送効率を変更させることができる。以下、 q/p 、 r 、 s の各パラメータの設定例について説明する。

【0049】

例えば、非常に多くの返信要求があるデータに対して、記録確率 r の値を大きくし、返信確率 s の値を小さくすると、各ノードから送信されるデータが少なくなり、ノードにおけるデータの検索処理やデータの送信処理が簡略化されるようになる。

30

【0050】

また、記録確率 r の値を大きくするかわりに、符号化率 q/p の値を大きくして、返信確率 s を小さくしても、各ノードにおけるデータの検索処理やデータの送信処理を簡略化することができる。

【0051】

また、符号化率 q/p を小さくし、記録確率 r を大きくすると、送信するパケットの数を抑えることができる。これは、 p が十分大きいときに効果的である。また、記録確率 r を小さくし、符号化率 q/p を小さくすることで、同一のパケットを複数のノードに記録することを避けることができる。これは、 p が十分小さいときに効果がある。

【0052】

また、返信確率 s を大きくし、記録確率 r 又は符号化率 q/p を小さくすることにより、分散型ストレージシステム全体に記録される符号化データの容量を小さくすることができる。或いは、データの記録時、出力時又は送信時などに、パケットの消失する確率を a とすると、 $a \times n \times q$ が p よりも十分大きくなるように r 、 s 、 q/p の値を制御することによって、十分な数のデータが返信される。

40

【0053】

また、複数のノードから返信されるユニークなパケットの個数を数学的に推定し、記録確率 r 、返信確率 s 、または符号化率 q/p を大きくすることで、ユニークなパケットが到着する確率を高くすることができる。

【0054】

50

以上のように、本具体例における分散型ストレージシステムは、システムを構成するノードに記録確率 を記述したパケット40を送信することにより、データを分散して記録する。そして、ノードに記録されたデータを返信確率 で返信させることにより、データの取り出しを行う。このように、データを分散して記録すると、一つのサーバにデータ管理の負荷が集中することなく、データを格納することができる。また、本具体例における分散型ストレージシステムでは、一つのデータを複数のノードで共有するため、システム全体で必要となるデータ容量が少なくて済む。

【0055】

また、複数のノードからデータを受信するようにすると、一つのサーバにトラフィックが集中することがなくなり、安定した通信量でデータの送受信が行える。

10

【0056】

また、 q/p の値を変更することにより、伝送効率や記録するデータ量を変更することができる。さらに、データの記録時、出力時、送信時などにパケットが消失するとき、消失する確率 a を考慮して、 q/p のパラメータの値を変更して、パケットが消失しても復号に十分なデータを取り出せるようにすることもできる。

【0057】

また、この分散型ストレージシステムでは、記録確率と返信確率に基づいた演算によってデータの分散記録・読み出しを実行できるため、データの管理が単純であり、カムコーダや携帯電話などの処理能力の少ない家庭用の機器であってもこのシステムに適応できる。また、カムコーダや携帯電話などの処理能力の低い家庭用の機器をノードとして用いることができるようになるため、数百万台規模の分散型ストレージシステムが容易に構築できる。

20

【0058】

また、上記分散ストレージシステムにおいては、記録指令装置は、記録させるデータと記録確率を同じパケットで送信したが、記録させるデータと記録確率を異なるパケットで送信したり、外部の記録装置に記録し、各情報処理ノードから参照するようにしてもよい。

【0059】

なお、上述した具体例は、発明の一例に過ぎず、本発明の要旨を含む変形例は、本発明に含まれるものとする。例えば、ノードとしては、カムコーダ、パーソナルビデオレコーダーやホームゲートウェイなどが考えられるが、データを記録する記録部と、所定の演算を行う制御部と、データの送受信を行うネットワークインターフェースを有していればその他の構成を備える装置であってもよい。

30

【0060】

また、記録確率、返信確率 をパケット中に記録したが、記録確率、返信確率 を任意の記録装置、若しくはパケットなどに記録し、各ノードがその値を参照するようにしてもよい。特に、FEC符号化としてリードトルネード符号化方式を利用した場合には、インターリーブ処理を省略することもできる。

【0061】

【発明の効果】

以上詳細に説明したように、本発明に係るデータ蓄積方法によれば、第1のデータを符号化して得た第2のデータをパケット化し、個々のパケットを複数の情報処理端末で記録するか否かを表す記録確率を生成し、記録確率とパケットを複数の情報処理端末に送り出し、この記録確率で複数の情報処理端末の何れかに記録されたパケットを読み出す際には、パケットを返信するか否かを示す返信確率を付して返信要求データを生成して返信要求データを複数の情報処理端末に送信し、複数の情報処理端末は、返信確率に基づいて該パケットを返信することにより、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録するとともに記録場所からデータを取り出して再構成することができる。

40

【0062】

また、本発明に係るデータ蓄積システムによれば、第1のデータを符号化して得た第2の

50

データをパケット化し、個々のパケットを複数の情報処理端末で記録するか否かを表す記録確率を生成し、パケットおよび記録確率を送り出し、この記録確率で複数の情報処理端末の何れかに記録されたパケットを読み出す際には、パケットを返信するか否かを示す返信確率を付して返信要求データを生成して返信要求データを複数の情報処理端末に送信し、複数の情報処理端末は、返信確率に基づいて該パケットを返信することにより、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録するとともに記録場所からデータを取り出して再構成することができる。

【0063】

また、本発明に係るデータ記録制御装置によれば、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録するとともに記録場所からデータを取り出すよう制御し、これを再構成できる。

10

【0064】

また、本発明に係るデータ記録指令装置によれば、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に負担をかけることなく多量のデータを分散して記録することができる。

【0065】

さらに、本発明に係るデータ受信装置によれば、複雑な処理や膨大な伝送帯域を必要とせず、個々の端末に分散して記録されたデータを受信して再構成することができる。

【0066】

さらにまた、本発明に係る情報処理端末によれば、複雑な処理や膨大な伝送帯域を必要とせず、負担なく多量のデータを分散して記録できる。

20

【図面の簡単な説明】

【図1】分散型ストレージシステムの構成を示す図である。

【図2】記録指令装置の内部構成を示す図である。

【図3】記録するデータを格納するパケットの構成を示す図である。

【図4】ノードの内部構成を示す図である。

【図5】送信指令装置の内部構成を示す図である。

【図6】データの送信を要求するパケットの構成を示す図である。

【図7】ノードから返信されるパケットの構成を示す図である。

【図8】パケットの生成工程を示す図である。

30

【図9】記録するデータを含むパケットの送信経路を示す図である。

【図10】データの返信を要求するパケットの送信経路を示す図である。

【図11】ノードから送信指令装置に返信されるパケットの送信経路を示す図である。

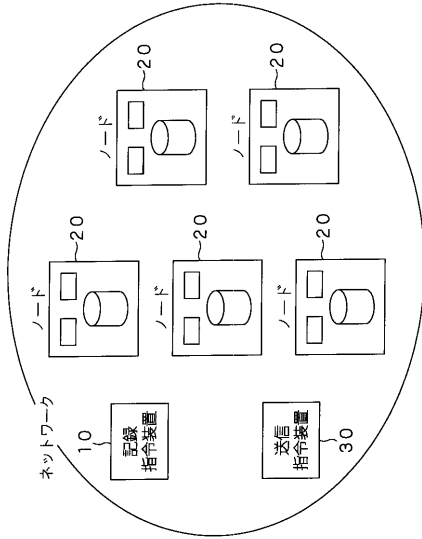
【図12】ノードから送信されたパケットを元データに復元する工程を示す図である。

【符号の説明】

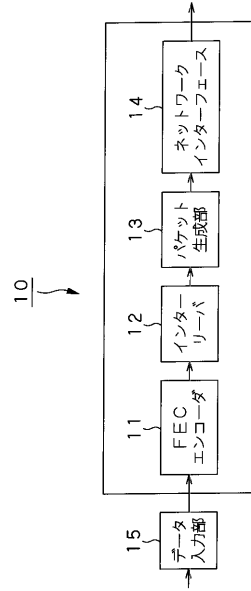
10 記録指令装置、11 FECエンコーダ、12 インターリーバ、13 パケット生成部、20 ノード、22 制御部、23 記録部、30 送信指令装置、32 データ要求部、33 パケット結合部、34 デインターリーバ、35 FECデコーダ、40 パケット、42 記録確率記述部、50 パケット、52 返信確率記述部、53 リクエスト記述部、60 パケット、62
ペイロード

40

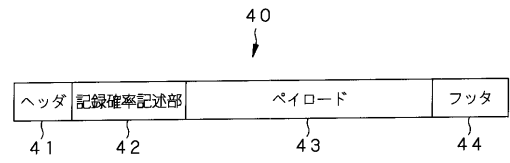
【図1】



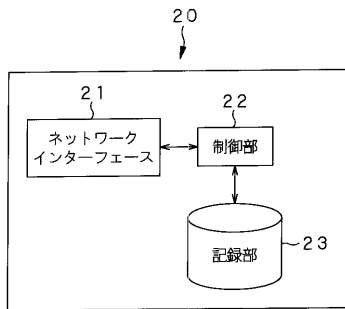
【図2】



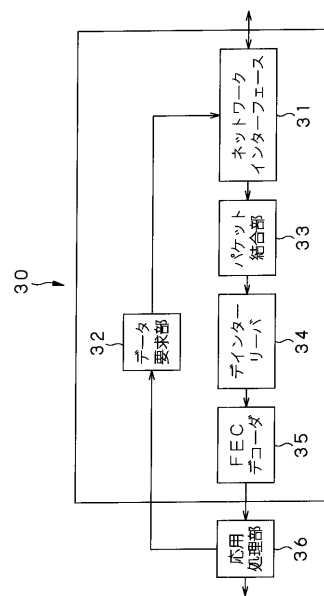
【図3】



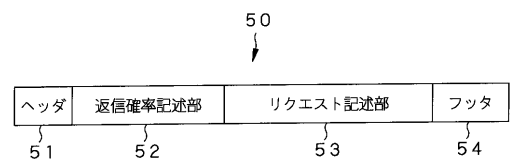
【図4】



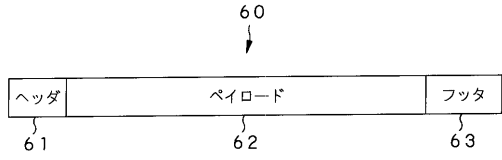
【図5】



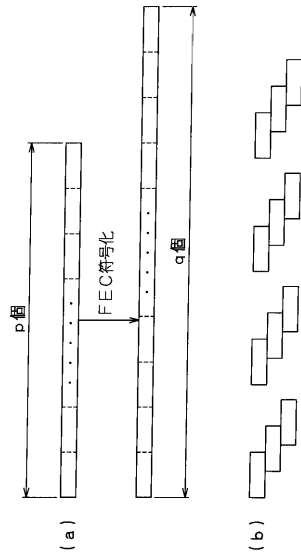
【図6】



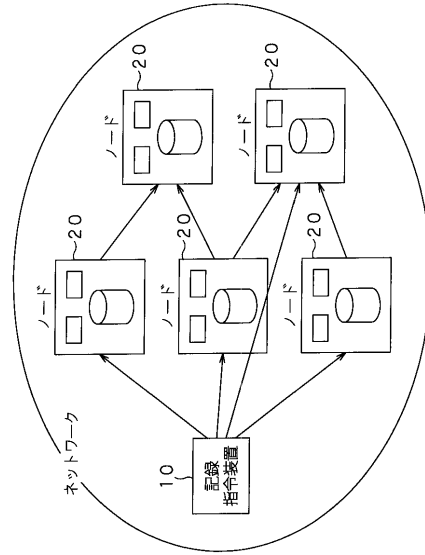
【図7】



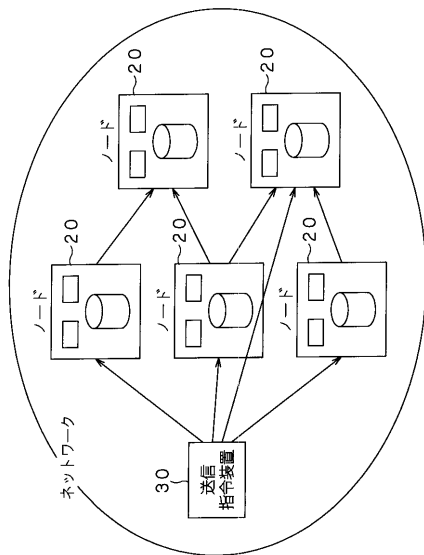
【図8】



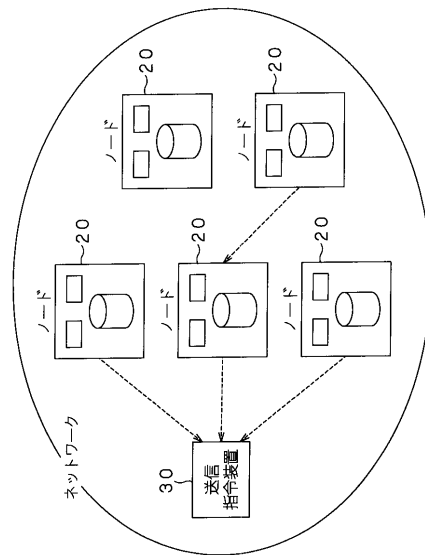
【図9】



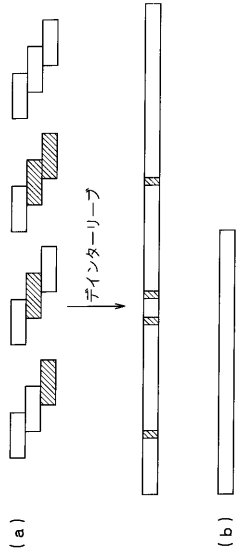
【図10】



【図11】



【 図 12 】



フロントページの続き

(72)発明者 片山 靖
東京都品川区北品川6丁目7番35号 ソニー株式会社内

審査官 須藤 竜也

(56)参考文献 特開2000-155712(JP,A)
国際公開第99/034291(WO,A1)
国際公開第00/031945(WO,A1)
特開平03-064141(JP,A)
特開平07-87092(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 13/00

H04L 12/56