

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5320678号
(P5320678)

(45) 発行日 平成25年10月23日(2013.10.23)

(24) 登録日 平成25年7月26日(2013.7.26)

(51) Int.Cl.		F I			
G06F 12/00	(2006.01)		G06F 12/00	531D	
G06F 3/06	(2006.01)		G06F 12/00	501B	
			G06F 3/06	540	
			G06F 3/06	304P	
			G06F 3/06	304F	

請求項の数 39 (全 46 頁)

(21) 出願番号	特願2007-40093 (P2007-40093)	(73) 特許権者	000004237
(22) 出願日	平成19年2月20日(2007.2.20)		日本電気株式会社
(65) 公開番号	特開2008-204206 (P2008-204206A)		東京都港区芝五丁目7番1号
(43) 公開日	平成20年9月4日(2008.9.4)	(74) 代理人	100103090
審査請求日	平成22年1月19日(2010.1.19)		弁理士 岩壁 冬樹
		(74) 代理人	100124501
			弁理士 塩川 誠人
		(72) 発明者	大和 純一
			東京都港区芝五丁目7番1号 日本電気株式会社社内
		審査官	池田 聡史

最終頁に続く

(54) 【発明の名称】 データ分散格納システム及びデータ分散方法、それに用いる装置並びにそのプログラム

(57) 【特許請求の範囲】

【請求項1】

コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムであって、

分割データの複製数を決定する複製数計画手段と、

コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する分割データ管理手段と、

前記分割データに対応する複製データの格納先を示す複製管理情報を記憶する複製管理情報記憶手段と、

前記複製管理情報に基づいて、分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定するアクセス先決定手段とを備え、

前記複製数計画手段は、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、

前記分割データ管理手段が、記憶装置が追加されたときに、前記記憶装置群における複製データの再配置を行う

ことを特徴とするデータ分散格納システム。

【請求項2】

アクセス先決定手段は、少なくともコンテンツの一部を読み出すホスト端末に対し、前

記少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、通知する

請求項1に記載のデータ分散格納システム。

【請求項3】

アクセス先決定手段は、少なくともコンテンツの一部を読み出すホスト端末で、前記少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定する

請求項1または請求項2に記載のデータ分散格納システム。

【請求項4】

アクセス先決定手段は、分割データに対応する複数の複製データが格納されている複数の記憶装置から、乱数に基づいてランダムに少なくとも1つの記憶装置を決定する

請求項1から請求項3のうちのいずれか1項に記載のデータ分散格納システム。

【請求項5】

アクセス先決定手段は、分割データに対応する複数の複製データが格納されている複数の記憶装置のうち負荷の低い記憶装置をアクセス先として決定する

請求項1から請求項3のうちのいずれか1項に記載のデータ分散格納システム。

【請求項6】

複製管理情報記憶手段は、前記分割データ管理手段によって決定された複製データの格納先を示す複製管理情報を記憶する

請求項1から請求項5のうちのいずれか1項に記載のデータ分散格納システム。

【請求項7】

分割データ管理手段は、分割データに対応する複数の複製データが各々記憶装置に均等化されるように、複製データの格納先を決定する

請求項1から請求項6のうちのいずれか1項に記載のデータ分散格納システム。

【請求項8】

分割データ管理手段は、コンテンツにおける連続する所定数の分割データに対応する複数の複製データが分散配置されるように、複製データの配置先を決定する

請求項1から請求項7のうちのいずれか1項に記載のデータ分散格納システム。

【請求項9】

分割データ管理手段は、2つの記憶装置間で共有する同じ分割データに対応する複製データの数を共有数とした場合に、各々の記憶装置の組み合わせにおける共有数が均等化されるように、複製データの配置先を決定する

請求項1から請求項8のうちのいずれか1項に記載のデータ分散格納システム。

【請求項10】

分割データ管理手段は、複製データの配置先を決定する際に、該複製データと同じ分割データに対応する複製データの配置先が既に決定している場合に、配置先として決定済みの記憶装置を除いた記憶装置の中から、前記決定済み記憶装置との組み合わせにおける共有数が小さい記憶装置を該複製データの配置先として決定する

請求項9に記載のデータ分散格納システム。

【請求項11】

分割データ管理手段は、記憶装置が使用不能となった場合に、複製管理情報記憶手段から、前記使用不能となった記憶装置を格納先として示している情報を削除する

請求項1から請求項10のうちのいずれか1項に記載のデータ分散格納システム。

【請求項12】

分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行う複製処理手段を備え、

分割データ管理手段は、使用不能となった記憶装置が配置先として割り当てられている複製データについて、新たな配置先を決定するとともに、該複製データと同じ分割データに対応する複製データを格納している記憶装置の中から1つの記憶装置をコピー元記憶装

10

20

30

40

50

置に決定して、前記複製処理手段に複製データのコピー処理を行わせる

請求項 1 から請求項 1 1 のうちのいずれか 1 項に記載のデータ分散格納システム。

【請求項 1 3】

分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行う複製処理手段を備えたデータ分散格納システムであって、

分割データ管理手段は、所定のタイミングで、複製管理情報に基づき、複製データについて新たな配置先を決定して、前記複製処理手段に複製データのコピー処理および前の配置先からの削除処理を行わせる

請求項 1 から請求項 1 2 のうちのいずれか 1 項に記載のデータ分散格納システム。

10

【請求項 1 4】

分割データ管理手段は、複製数計画手段によって決定された分割データの複製数に応じて、該分割データに対応する複製データの配置先を決定する

請求項 1 から請求項 1 3 のうちのいずれか 1 項に記載のデータ分散格納システム。

【請求項 1 5】

複製数計画手段は、コンテンツにおける分割データの位置に基づいて、各分割データの複製数を決定する

請求項 1 4 に記載のデータ分散格納システム。

【請求項 1 6】

複製数計画手段は、コンテンツの先頭に位置するデータを含む分割データ、およびコンテンツにチャプターが付与されている場合に、各チャプターの先頭に位置するデータを含む分割データの複製数を、該コンテンツにおける他の分割データの複製数よりも多くする

請求項 1 5 に記載のデータ分散格納システム。

20

【請求項 1 7】

複製数計画手段は、コンテンツのアクセス要求予測量またはアクセス要求量と、各記憶装置の処理能力と、システムで規定した安全係数とに基づいて、分割データの複製数を決定する

請求項 1 から請求項 1 6 のうちのいずれか 1 項に記載のデータ分散格納システム。

【請求項 1 8】

分割データ管理手段からの指示に従い記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行う複製処理手段を備えたデータ分散格納システムであって、

分割データ管理手段は、複製数計画手段によって決定された分割データの複製数と、記憶装置に格納されている該分割データに対応する複製データの数が一致するように、前記分割データに対応する複製データの新たな配置先または削除対象とする配置先を決定し、前記複製処理手段に該複製データのコピー処理または削除処理を行わせる

請求項 1 から請求項 1 7 のうちのいずれか 1 項に記載のデータ分散格納システム。

30

【請求項 1 9】

コンテンツが複数の分割データに分割され、少なくとも 1 つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムに適用され、少なくともコンテンツの一部を読み出すホスト端末に対し、前記コンテンツを構成する分割データへのアクセス先を通知する分割データ管理装置であって、

分割データの複製数を決定する複製数計画手段と、

コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する分割データ管理手段と、

前記ホスト端末のコンテンツ読み出し範囲に含まれる分割データに対応する複製データの配置先を示す複製管理情報に基づいて、前記分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも 1 つの記憶装置を決定するアクセス先決定手段とを備え、

前記複製数計画手段は、記憶装置に対するアクセスが不能になったときの複製データの

40

50

他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、前記分割データ管理手段が、記憶装置が追加されたときに、前記記憶装置群における複製データの再配置を行う

ことを特徴とする分割データ管理装置。

【請求項 20】

コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムに適用されるデータ分散方法であって、

分割データ管理手段が、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定し、

複製管理情報記憶手段が、前記分割データに対応する複製データの格納先を示す複製管理情報を記憶し、

アクセス先決定手段が、前記複製管理情報に基づいて、分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、

複製数計画手段が、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、

前記分割データ管理手段が、記憶装置が追加されたときに、前記記憶装置群における複製データの再配置を行う

ことを特徴とするデータ分散方法。

【請求項 21】

アクセス先決定手段が、少なくともコンテンツの一部を読み出すホスト端末に対し、前記少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、通知する

請求項 20 に記載のデータ分散方法。

【請求項 22】

アクセス先決定手段が、少なくともコンテンツの一部を読み出すホスト端末で、前記少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定する

請求項 20 または請求項 21 に記載のデータ分散方法。

【請求項 23】

アクセス先決定手段が、分割データに対応する複数の複製データが格納されている複数の記憶装置から、乱数に基づいてランダムに少なくとも1つの記憶装置を決定する

請求項 20 から請求項 22 のうちのいずれか1項に記載のデータ分散方法。

【請求項 24】

アクセス先決定手段が、分割データに対応する複数の複製データが格納されている複数の記憶装置のうち負荷の低い記憶装置をアクセス先として決定する

請求項 20 から請求項 23 のうちのいずれか1項に記載のデータ分散方法。

【請求項 25】

複製管理情報記憶手段が、前記分割データ管理手段によって決定された複製データの格納先を示す複製管理情報を記憶する

請求項 20 から請求項 24 のうちのいずれか1項に記載のデータ分散方法。

【請求項 26】

分割データ管理手段が、分割データに対応する複数の複製データが各々記憶装置に均等化されるように、複製データの格納先を決定する

請求項 20 から請求項 25 のうちのいずれか1項に記載のデータ分散方法。

【請求項 27】

分割データ管理手段が、コンテンツにおける連続する所定数の分割データに対応する複数の複製データが分散配置されるように、複製データの配置先を決定する

10

20

30

40

50

請求項 20 から請求項 26 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 28】

分割データ管理手段が、2つの記憶装置間で共有する同じ分割データに対応する複製データの数を共有数とした場合に、各々の記憶装置の組み合わせにおける共有数が均等化されるように、複製データの配置先を決定する

請求項 20 から請求項 27 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 29】

分割データ管理手段が、複製データの配置先を決定する際に、該複製データと同じ分割データに対応する複製データの配置先が既に決定している場合に、配置先として決定済みの記憶装置を除いた記憶装置の中から、前記決定済み記憶装置との組み合わせにおける共有数が小さい記憶装置を該複製データの配置先として決定する

10

請求項 28 に記載のデータ分散方法。

【請求項 30】

分割データ管理手段が、記憶装置が使用不能となった場合に、複製管理情報記憶手段から、前記使用不能となった記憶装置を格納先として示している情報を削除する

請求項 20 から請求項 29 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 31】

分割データ管理手段が、使用不能となった記憶装置が配置先として割り当てられている複製データについて、新たな配置先を決定するとともに、該複製データと同じ分割データに対応する複製データを格納している記憶装置の中から 1 つの記憶装置をコピー元記憶装置に決定し、

20

複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理を行う

請求項 20 から請求項 30 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 32】

分割データ管理手段が、所定のタイミングで、複製管理情報に基づき、複製データについて新たな配置先を決定し、

複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および前の配置先からの削除処理を行う

請求項 20 から請求項 31 のうちのいずれか 1 項に記載のデータ分散方法。

30

【請求項 33】

分割データ管理手段が、複製数計画手段によって決定された分割データの複製数に応じて、該分割データに対応する複製データの配置先を決定する

請求項 20 から請求項 32 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 34】

複製数計画手段が、コンテンツにおける分割データの位置に基づいて、各分割データの複製数を決定する

請求項 33 に記載のデータ分散方法。

【請求項 35】

複製数計画手段が、コンテンツの先頭に位置するデータを含む分割データ、およびコンテンツにチャプターが付与されている場合に、各チャプターの先頭に位置するデータを含む分割データの複製数を、該コンテンツにおける他の分割データの複製数よりも多くする

40

請求項 34 に記載のデータ分散方法。

【請求項 36】

複製数計画手段が、コンテンツのアクセス要求予測量またはアクセス要求量と、各記憶装置の処理能力と、システムで規定した安全係数とに基づいて、分割データの複製数を決定する

請求項 20 から請求項 35 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 37】

分割データ管理手段が、複製数計画手段によって決定された分割データの複製数と、記

50

憶装置に格納されている該分割データに対応する複製データの数が一致するように、前記分割データに対応する複製データの新たな配置先または削除対象とする配置先を決定し

、複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行う

請求項 20 から請求項 36 のうちのいずれか 1 項に記載のデータ分散方法。

【請求項 38】

コンテンツが複数の分割データに分割され、少なくとも 1 つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムにおいて、少なくともコンテンツの一部を読み出すホスト端末に対し、前記コンテンツを構成する分割データへのアクセス先を通知する分割データ管理装置に適用されるデータ分散用プログラムであって、

コンピュータに、

コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する処理、

前記ホスト端末のコンテンツ読み出し範囲に含まれる分割データに対応する複製データの配置先を示す複製管理情報に基づいて、前記分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも 1 つの記憶装置を決定する処理、

記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定する処理、および、

記憶装置が追加されたときに、前記記憶装置群において複製データを再配置する処理を実行させるためのデータ分散用プログラム。

【請求項 39】

コンピュータに、

コンテンツを複数の分割データに分割し、少なくとも 1 つの分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する処理

を実行させる請求項 38 に記載のデータ分散用プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ分散格納システム、データ分散方法、分割データ管理装置、およびデータ分散用プログラムに関し、特に、複数の記憶装置にデータを分散させて格納するデータ分散格納システム、データ分散方法、分割データ管理装置、およびデータ分散用プログラムに関する。

【背景技術】

【0002】

あるデータを複数の記憶装置に分散させるデータ分散方法に関し、ストリーミングを再生するために複数のコンテンツを形成する各データを分割して複数の記憶装置に格納させる方法がある。例えば、特許文献 1 には、一連のデータストリームとして処理されて意味をなすデータを、複数のファイル（対応する磁気ディスク装置）に跨り連続的に分割して書き込むデータ制御方法が記載されている。また、例えば、特許文献 2 には、ファイル上の連続するファイルブロックをそれぞれ別の物理ブロック群に割り付けられるよう配置制御するディスクアレイ装置が記載されている。また、例えば、特許文献 3 には、予め仮想アドレスの指定順序が明らかになっている場合に、それらが各記憶装置に対する均一なアクセスとなり、複数のクライアントによる記憶装置群へのアクセスが同一の記憶装置に重ならないように、分割したデータの配置先を決定するディスクアレイ装置が記載されている。

【0003】

また、非特許文献 1 に記載されているように、データを分散させるだけでなく、データ

の冗長化を行う方法も考えられている。非特許文献 1 には、RAID 1 + 0 や RAID 0 + 1 として、ブロックストレージのレベルでデータの分散および冗長化を行う方法が開示されている。また、特許文献 4 には、複製したデータブロックをアドレスに応じて均等に他の記憶装置に分散させて格納する方法が開示されている。

【0004】

また、特許文献 5 に記載されているように、アクセス負荷が限界値に達する前に、コンテンツの複製の作成やコンテンツとその他のコンテンツとで所在の再配置を行う方法も考えられている。

【0005】

【特許文献 1】特開 2002 - 244893 号公報

【特許文献 2】特開平 09 - 223049 号公報

【特許文献 3】特許第 3052877 号公報

【特許文献 4】特許第 2853624 号公報

【特許文献 5】特開平 11 - 085604 号公報

【非特許文献 1】John L. Hennessy, David A. Paterson, "Computer Architecture: A Quantitative Approach", 3rd Edition, Morgan Kaufmann Pub, 2001, pp.707

【発明の開示】

【発明が解決しようとする課題】

【0006】

しかしながら、特許文献 1 ~ 3 に記載されている方法は、データを複数に分割して複数の記憶装置に分散させて格納することによってスループットを上げることはできるが、障害時の対応については何ら考慮がされていない。なお、特許文献 2 に記載されている方法では、論理ブロックと物理ブロックの対応関係を予め設定しておかなければならず、一度ディスクアレイを構成するとハードディスクの追加が容易ではないという問題もある。

【0007】

また、非特許文献 1 に記載されているようなブロックベースの分散では、冗長化することで障害時の耐久性（信頼性）の向上という点では効果があるが、スループットが向上するとは限らない。その理由は、ブロックベースの分散では、データがコンテンツとしては認識されていないため、同一コンテンツに関するアクセスの並列性を生かそうということが考えられていないためである。

【0008】

例えば、ブロックベースの分散を行うストレージ上にファイルシステムを構築した場合、ファイルシステムは仮想的なブロックデバイスのどこが実際のブロックデバイスに割り当てられているかを認識していない。このため、一つのコンテンツ（ファイル）を複数の仮想的なブロックデバイスに割り当てたとしても、複数の物理的なブロックデバイスに割り当てられる保障はない。

【0009】

また、特許文献 4 に記載されている方法は、障害時の対応についても考慮されているが、通常用の記憶装置と障害用の記憶装置とを分けて扱っているため、記憶装置の容量効率がよくない。また、ある記憶装置に対応するデータ（すなわち、その記憶装置が記憶しているデータと同一内容のデータ）が、他の一つの記憶装置に記憶されているため、障害時やその回復時にアクセス性能が維持できないという問題がある。例えば、ある記憶装置が故障している間、その記憶装置が記憶していたデータのアクセス負荷が、同じデータを記憶している一つの記憶装置にだけかかることになる。また、故障した記憶装置の代わりに新しい記憶装置を追加した際には、コピー処理が完了するまでの間、コピー処理の読み出し元における負荷がその記憶装置にだけかかることになり、スループットを制限しつつコピー処理を行うとしても、コピー処理の時間がかかる分、その記憶装置のアクセス性能は長い時間低下することとなる。

【0010】

さらに、特許文献 4 に記載されている方法では、アドレスを元にデータ（複製されたデ

10

20

30

40

50

ータも含む)の配置を決定しているので、一度ディスクアレイを構成するとハードディスクの追加が容易ではないという問題もある。

【0011】

また、特許文献5に記載されている方法は、デマンドに対応させることで、再配置後にスループットを向上させることはできるが、複製の追加や再配置の処理の際にアクセス性能を維持しようという点は、何ら考慮がされていない。

【0012】

すなわち、従来のデータ分散方法の問題点は、第1に、容量効率を下げることなく信頼性の向上とスループットの向上とを同時に充足することができないことである。第2に、障害時や複製の追加時や再配置時等、通常時以外でのアクセス性能が維持できないことである。第3に、記憶装置の追加等、スケーラビリティの向上が図られていないという点である。

【0013】

そこで、本発明は、以上の問題を解決すべく、データ分散における特性・性能の向上をさらに図ることを目的とする。具体的には、信頼性の向上とスループットの向上とを同時に充足できるデータ分散格納システム及びデータ分散方法、それに用いる装置並びにそのプログラムを提供することを目的とする。また、通常時以外であってもアクセス性能が維持できるデータ分散格納システム及びデータ分散方法、それに用いる装置並びにそのプログラムを提供することを目的とする。また、記憶装置の追加等、スケーラビリティの向上が図れるようなデータ分散格納システム及びデータ分散方法、それに用いる装置並びにそのプログラムを提供することを目的とする。

【課題を解決するための手段】

【0014】

本発明によるデータ分散格納システムは、コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムであって、分割データの複製数を決定する複製数計画手段と、コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する分割データ管理手段と、分割データに対応する複製データの格納先を示す複製管理情報を記憶する複製管理情報記憶手段(例えば、分割データ管理DB7)と、複製管理情報に基づいて、分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定するアクセス先決定手段(例えば、分割データ管理部3やホスト1のアクセス先決定機能)とを備え、複製数計画手段は、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、分割データ管理手段が、記憶装置が追加されたときに、記憶装置群における複製データの再配置を行うことを特徴とする。

【0015】

また、アクセス先決定手段は、少なくともコンテンツの一部を読み出すホスト端末に対し、ホスト端末が読み出す少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、通知してもよい。このアクセス先決定手段は、例えば、分割データ管理部3のアクセス先決定機能によって実現される。

【0016】

また、アクセス先決定手段は、少なくともコンテンツの一部を読み出すホスト端末で、ホスト端末が読み出す少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定してもよい。このアクセス先決定手段は、例えば、ホスト1のアクセス先決定機能によって実現される。

【0017】

また、アクセス先決定手段は、分割データに対応する複数の複製データが格納されてい

10

20

30

40

50

る複数の記憶装置から、乱数に基づいてランダムに少なくとも1つの記憶装置を決定してもよい。

【0018】

また、アクセス先決定手段は、分割データに対応する複数の複製データが格納されている複数の記憶装置のうち負荷の低い記憶装置をアクセス先として決定してもよい。

【0019】

また、複製管理情報記憶手段は、分割データ管理手段によって決定された複製データの格納先を示す複製管理情報を記憶してもよい。

【0020】

また、分割データ管理手段は、分割データに対応する複数の複製データが各々記憶装置に均等化されるように、複製データの格納先を決定してもよい。

10

【0021】

また、分割データ管理手段は、コンテンツにおける連続する所定数の分割データに対応する複数の複製データが分散配置されるように、複製データの配置先を決定してもよい。

【0022】

また、分割データ管理手段は、2つの記憶装置間で共有する同じ分割データに対応する複製データの数を共有数とした場合に、各々の記憶装置の組み合わせにおける共有数が均等化されるように、複製データの配置先を決定してもよい。

【0023】

また、分割データ管理手段は、複製データの配置先を決定する際に、その複製データと同じ分割データに対応する複製データの配置先が既に決定している場合に、配置先として決定済みの記憶装置を除いた記憶装置の中から、決定済み記憶装置との組み合わせにおける共有数が小さい記憶装置を複製データの配置先として決定してもよい。

20

【0024】

分割データ管理手段は、記憶装置が使用不能となった場合に、複製管理情報記憶手段から、使用不能となった記憶装置を格納先として示している情報を削除してもよい。

【0025】

また、データ分割可能システムは、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行う複製処理手段（例えば、複製処理部5）を備え、分割データ管理手段は、使用不能となった記憶装置が配置先として割り当てられている複製データについて、新たな配置先を決定するとともに、該複製データと同じ分割データに対応する複製データを格納している記憶装置の中から1つの記憶装置をコピー元記憶装置に決定して、複製処理手段に複製データのコピー処理を行わせてもよい。

30

【0026】

また、分割データ管理手段は、所定のタイミングで、複製管理情報に基づき、複製データについて新たな配置先を決定して、複製処理手段に複製データのコピー処理および前の配置先からの削除処理を行わせてもよい。

【0027】

また、分割データ管理手段は、複製数計画手段によって決定された分割データの複製数に応じて、該分割データに対応する複製データの配置先を決定してもよい。

40

【0028】

また、複製数計画手段は、コンテンツにおける分割データの位置に基づいて、各分割データの複製数を決定してもよい。

【0029】

また、複製数計画手段は、コンテンツの先頭に位置するデータを含む分割データ、およびコンテンツにチャプターが付与されている場合に、各チャプターの先頭に位置するデータを含む分割データの複製数を、そのコンテンツにおける他の分割データの複製数よりも多くしてもよい。

【0031】

50

また、複製数計画手段は、コンテンツのアクセス要求予測量またはアクセス要求量と、各記憶装置の処理能力と、システムで規定した安全係数とに基づいて、分割データの複製数を決定してもよい。

【0032】

また、分割データ管理手段は、複製数計画手段によって決定された分割データの複製数と、記憶装置に格納されているその分割データに対応する複製データの数が一致するように、その分割データに対応する複製データの新たな配置先または削除対象とする配置先を決定し、複製処理手段に該複製データのコピー処理または削除処理を行わせてもよい。

【0033】

また、本発明による分割データ管理装置は、コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムに適用され、少なくともコンテンツの一部を読み出すホスト端末に対し、そのコンテンツを構成する分割データへのアクセス先を通知する分割データ管理装置であって、分割データの複製数を決定する複製数計画手段と、コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する分割データ管理手段と、ホスト端末のコンテンツ読み出し範囲に含まれる分割データに対応する複製データの配置先を示す複製管理情報に基づいて、その分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定するアクセス先決定手段（例えば、分割データ管理部3のアクセス先決定機能）とを備え、複製数計画手段は、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、分割データ管理手段が、記憶装置が追加されたときに、記憶装置群における複製データの再配置を行うことを特徴とする。

【0035】

また、本発明によるデータ分散方法は、コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群に複数の記憶装置に格納されるデータ分散格納システムに適用されるデータ分散方法であって、分割データ管理手段が、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定し、複製管理情報記憶手段が、分割データに対応する複製データの格納先を示す複製管理情報を記憶し、アクセス先決定手段が、複製管理情報に基づいて、分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、複製数計画手段が、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定し、分割データ管理手段が、記憶装置が追加されたときに、記憶装置群における複製データの再配置を行うことを特徴とする。

【0036】

また、データ分散方法は、アクセス先決定手段が、少なくともコンテンツの一部を読み出すホスト端末に対し、ホスト端末が読み出す少なくともコンテンツの一部を構成する分割データへのアクセス先として該分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定し、通知してもよい。

【0037】

また、データ分散方法は、アクセス先決定手段が、少なくともコンテンツの一部を読み出すホスト端末で、ホスト端末が読み出す少なくともコンテンツの一部を構成する分割データへのアクセス先として分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定してもよい。

【0038】

また、データ分散方法は、アクセス先決定手段が、分割データに対応する複数の複製データが格納されている複数の記憶装置から、乱数に基づいてランダムに少なくとも1つの記憶装置を決定してもよい。

【0039】

10

20

30

40

50

また、データ分散方法は、アクセス先決定手段が、分割データに対応する複数の複製データが格納されている複数の記憶装置のうち負荷の低い記憶装置をアクセス先として決定してもよい。

【0040】

また、データ分散方法は、複製管理情報記憶手段が、分割データ管理手段によって決定された複製データの格納先を示す複製管理情報を記憶してもよい。

【0041】

また、データ分散方法は、分割データ管理手段が、分割データに対応する複数の複製データが各々記憶装置に均等化されるように、複製データの格納先を決定してもよい。

【0042】

また、データ分散方法は、分割データ管理手段が、コンテンツにおける連続する所定数の分割データに対応する複数の複製データが分散配置されるように、複製データの配置先を決定してもよい。

【0043】

また、データ分散方法は、分割データ管理手段が、2つの記憶装置間で共有する同じ分割データに対応する複製データの数を共有数とした場合に、各々の記憶装置の組み合わせにおける共有数が均等化されるように、複製データの配置先を決定してもよい。

【0044】

また、データ分散方法は、分割データ管理手段が、複製データの配置先を決定する際に、その複製データと同じ分割データに対応する複製データの配置先が既に決定している場合に、配置先として決定済みの記憶装置を除いた記憶装置の中から、決定済み記憶装置との組み合わせにおける共有数が小さい記憶装置を該複製データの配置先として決定してもよい。

【0045】

また、データ分散方法は、分割データ管理手段が、記憶装置が使用不能となった場合に、複製管理情報記憶手段から、使用不能となった記憶装置を格納先として示している情報を削除してもよい。

【0046】

また、データ分散方法は、分割データ管理手段が、使用不能となった記憶装置が配置先として割り当てられている複製データについて、新たな配置先を決定するとともに、その複製データと同じ分割データに対応する複製データを格納している記憶装置の中から1つの記憶装置をコピー元記憶装置に決定し、複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理を行ってもよい。

【0047】

また、データ分散方法は、分割データ管理手段が、所定のタイミングで、複製管理情報に基づき、複製データについて新たな配置先を決定し、複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および前の配置先からの削除処理を行ってもよい。

【0048】

また、データ分散方法は、分割データ管理手段が、複製数計画手段によって決定された分割データの複製数に応じて、分割データに対応する複製データの配置先を決定してもよい。

【0049】

また、データ分散方法は、複製数計画手段が、コンテンツにおける分割データの位置に基づいて、各分割データの複製数を決定してもよい。

【0050】

また、データ分散方法は、複製数計画手段が、コンテンツの先頭に位置するデータを含む分割データ、およびコンテンツにチャプターが付与されている場合に、各チャプターの先頭に位置するデータを含む分割データの複製数を、そのコンテンツにおける他の分割データの複製数よりも多くしてもよい。

10

20

30

40

50

【 0 0 5 2 】

また、データ分散方法は、複製数計画手段が、コンテンツのアクセス要求予測量またはアクセス要求量と、各記憶装置の処理能力と、システムで規定した安全係数とに基づいて、分割データの複製数を決定してもよい。

【 0 0 5 3 】

また、データ分散方法は、分割データ管理手段が、複製数計画手段によって決定された分割データの複製数と、記憶装置に格納されているその分割データに対応する複製データの数が一致するように、分割データに対応する複製データの新たな配置先または削除対象とする配置先を決定し、複製処理手段が、分割データ管理手段からの指示に従い、記憶装置間の複製データのコピー処理および記憶装置からの複製データの削除処理を行ってもよい。

10

【 0 0 5 4 】

また、本発明によるデータ分散用プログラムは、コンテンツが複数の分割データに分割され、少なくとも1つの分割データに対応する複数の複製データが記憶装置群における複数の記憶装置に格納されるデータ分散格納システムにおいて、少なくともコンテンツの一部を読み出すホスト端末に対し、そのコンテンツを構成する分割データへのアクセス先を通知する分割データ管理装置に適用されるデータ分散用プログラムであって、コンピュータに、コンテンツを複数の分割データに分割し、分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する処理、ホスト端末のコンテンツ読み出し範囲に含まれる分割データに対応する複製データの配置先を示す複製管理情報に基づいて、分割データへのアクセス先としてその分割データに対応する複製データが格納されている記憶装置のうち少なくとも1つの記憶装置を決定する処理、記憶装置に対するアクセスが不能になったときの複製データの他の記憶装置からのアクセスを保障できる数に応じて、分割データの複製数を決定する処理、および、記憶装置が追加されたときに、記憶装置群において複製データを再配置する処理を実行させることを特徴とする。

20

【 0 0 5 5 】

また、コンピュータに、コンテンツを複数の分割データに分割し、少なくとも1つの分割データに対応する複数の複製データの格納先として、複数の記憶装置を決定する処理を実行させてもよい。

【 発明の効果 】

30

【 0 0 5 7 】

本発明によれば、データ分散における特性・性能の向上を図ることができる。その理由は、少なくとも1つの分割データに対応する複数の複製データが複数の記憶装置に格納された各分割データの配置先を複製管理情報記憶手段が複製管理情報として記憶し、アクセス先決定手段が、その複製管理情報に基づいて、通常用と障害用とを区別することなく複数の記憶装置のうち少なくとも1つの記憶装置をアクセス先として決定するからである。

【 発明を実施するための最良の形態 】

【 0 0 5 8 】

まず、本発明の概要について説明する。本発明は、コンテンツを複数のデータに分割して記憶するデータ分散格納システムにおいて、個々の分割データをそれぞれ2以上に複製した上で、通常用と障害用等を区別することなく、使用可能な全記憶装置を配置先候補にして、同一内容の分割データがそれぞれ異なる記憶装置に格納されるように配置決定を行うことを特徴とする。さらに、同一内容の分割データが異なる記憶装置に格納されるだけでなく、各記憶装置間の相関（ある記憶装置と他のいずれの記憶装置との間における、2記憶装置間で共有する分割データの数）が低くなるように配置決定を行うことを特徴とする。さらに、ストリーミング再生用のコンテンツなど、データの連続性が認められるようなコンテンツを対象とする場合には、1度にアクセスされる可能性の高いデータサイズ（例えば、アクセスに用いるバッファ量相当）分を目安に、連続する分割データが同一の記憶装置に格納されないように配置決定を行うことを特徴とする。

40

【 0 0 5 9 】

50

また、配置決定に基づき各記憶装置に分散配置させた上で、同一内容の分割データに対するアクセス先の決定処理や、障害時の回復処理、および複製数の変更処理を行うことを特徴とする。

【0060】

ここで、本発明を説明するために使用する用語を以下に定義する。本発明において、コンテンツとは、データ分散格納における分散格納の対象とするデータのまとまりをいう。例えば、コンテンツは、1つのファイルとしてまとめられたデータである。また、分割データとは、コンテンツを複数に分割したもの(データ)をいう。なお、各分割データのサイズは固定でなくてもよい。コンテンツ毎に異なるサイズであってもよいし、分割データ毎に異なるサイズであってもよい。例えば、音声や映像のようなエンコーディングされているコンテンツを対象とする場合には、そのエンコーディングの単位やその単位を複数まとめた単位を分割データのサイズとしてもよい。

10

【0061】

以下、本発明の実施の形態について図面を参照して説明する。図1は、本発明によるデータ分散格納システムの構成例を示すブロック図である。図1に示すデータ分散格納システムは、ホスト1と、複数の記憶装置2(記憶装置2-1~n。nは自然数)と、分割データ管理部3と、障害通知部4と、複製処理部5と、複製数計画部6と、分割データ管理データベース(分割データ管理DB)7とを備える。なお、図1では1つのホスト1しか示していないが、ホスト1は、本システムを利用するユーザの数に応じて複数存在する。

【0062】

ホスト1は、複数の記憶装置2に分散格納されたコンテンツにアクセスするためのユーザ端末である。記憶装置2は、コンテンツを分散格納するための記憶装置である。各記憶装置2は、コンテンツを分割してできる各分割データおよびその複製を分散して記憶する。

20

【0063】

分割データ管理部3は、コンテンツを構成する各分割データに対応する複数の複製データの配置先を管理する。なお、分割データ管理部3は、コンテンツ内におけるある1つの分割データに対し、その複製数に応じて複数の複製データの配置先を管理する。複製データとは、コンテンツを分割した分割データと同一内容のデータであって、実際に記憶装置に格納されるデータのことをいう。本発明において、“分割データ”をコンテンツ内における論理的な分割データとして表現するのに対し、“複製データ”を実際に記憶装置に格納する物理的な分割データとして表現するために用いている。なお、分割データ管理部3は、コンテンツを構成する各分割データに対応する複製データを複数の記憶装置に分散配置するために配置先を決定してその配置先を管理するだけでなく、既に分散配置されている状態からその配置先を管理してもよい。配置先の管理として、具体的には、分割データへのアクセス先としての記憶装置を決定したり、複製数の変更や記憶装置の追加、使用不能に伴う複製数の維持処理(複製データの追加、削除、再配置等)を行う。

30

【0064】

なお、分割データへのアクセス先としての記憶装置の決定機能(以下、単にアクセス先決定機能という。)は、ホスト1が有していてもよい。その場合、分割データ管理部3は、ホスト1からの要求に応じて、所望の分割データに対応する複製データの配置先を示す情報をホスト1に通知してもよい。

40

【0065】

障害通知部4は、記憶装置2の障害(具体的には、使用可否)を分割データ管理部3に通知する。複製処理部5は、記憶装置2への分割データの追加処理や削除処理やコピー処理など、記憶装置2に対するアクセス処理を行う。複製数計画部6は、各分割データの複製数を決定する。分割データ管理DB7は、分割データ管理部3が各分割データに対応する複製データの配置先を管理するための情報を記憶するための記憶装置である。分割データ管理DB7は、コンテンツと分割データの関連性に関する情報や、各複製データの配置先に関する情報や、配置先対象となる記憶装置に関する情報を記憶する。

50

【 0 0 6 6 】

以下、本実施の形態では、分割データのことを”チャンク”と言う。なお、”チャンク”という表現には、一般に意味的にまとまりのある分割データを指す場合があるが、本実施の形態においては、意味的にまとまりのある分割データに限定されず、意味的にまとまりのない分割データを含めてチャンクという。

【 0 0 6 7 】

図2は、本実施の形態によるデータ分散格納システムの構成例を示すブロック図である。図2に示すように、本実施の形態によるデータ分散格納システムは、ホスト1と、複数の記憶装置2（記憶装置2-1～n。nは自然数）と、チャンク管理部3と、障害通知部4と、複製処理部5と、複製数計画部6とを備える。なお、図2では1つのホスト1しか示していないが、ホスト1は、本システムを利用するユーザの数に応じて複数存在する。

10

【 0 0 6 8 】

ホスト1は、具体的には、複数の記憶装置2に分散格納されたコンテンツにアクセスするデータ処理装置である。本実施の形態では、ホスト1は、チャンク管理部3から取得した、そのコンテンツを構成するチャンクを記憶している記憶装置2を示すチャンク配置情報に基づいて、必要なチャンクにアクセスする。

【 0 0 6 9 】

記憶装置2-1～nは、各チャンクを格納する記憶装置である。記憶装置2-1～nは、ブロックベースの記憶装置に限らず、例えば、NAS等のファイルベースやOSD(Object based storage device)等のオブジェクトベースの記憶装置システムであってもよい。そのような場合には、チャンクは、ファイルやオブジェクトとして格納されることとなる。なお、記憶装置2-1～nは、記憶装置を制御する制御装置を含む。例えば、記憶装置2-1～nは、磁気記憶装置、光磁気記憶装置、不揮発性の半導体記憶装置、およびそれらのアレイ装置、並びにそれらを制御する制御装置を備えたサーバ装置によって実現される。また、記憶装置2-1～nには、それぞれを識別するための記憶装置IDが割り振られているものとする。

20

【 0 0 7 0 】

なお、本実施の形態において、記憶装置に記憶させるチャンクが複製かそうでないかを区別する必要はなく、単に、あるコンテンツを構成するチャンクと同一内容のチャンクが複数格納されている、として認識すれば足る。以下、”各チャンク”という表現には、コンテンツ内における各々のチャンクを示す場合と、その複製も含め記憶装置内における各々のチャンクを示す場合とを含んでいるものとする。なお、コンテンツ内において識別されるチャンクと、実際に複数の記憶装置に記憶させるチャンクとを区別するために、記憶装置に記憶させるチャンクを”複製チャンク”と表現する場合がある。

30

【 0 0 7 1 】

チャンク管理部3は、コンテンツを構成する各チャンクに対応する複数の複製チャンクの配置先を管理する。本実施の形態では、チャンク管理部3は、コンテンツを複数のチャンクに分割して、各チャンクに対しそれぞれ複数の配置先を決定して、その配置先を管理する。なお、チャンク管理部3は、コンテンツを構成する各チャンクに対しそれぞれ複数の配置先を決定してその配置先を管理するのではなく、既に分散配置されている状態からその配置先を管理してもよい。チャンク管理部3は、各チャンクの配置先を管理するためのチャンク管理データベース7（以下、チャンク管理DB7という。）を有する。なお、チャンク管理DB7は、独立したデータベースシステムとして存在していてもよい。その場合、少なくともチャンク管理部3と複製処理部5とがアクセス可能に接続されていればよい。

40

【 0 0 7 2 】

チャンク管理部3は、配置先の管理として、例えば、チャンクの複製数の変更が行われたときや障害時、他の処理部からの通知等、所定のタイミングでチャンクの削除や追加やコピー（再配置）を行う。なお、本実施の形態では、チャンク管理部3は各チャンクの配置先や参照元を決定するに留まり、実際の記憶装置への反映は複製処理部5が行う。また

50

、例えば、ホスト1からのコンテンツのアクセス要求を受けて、そのコンテンツを構成するチャンクのアクセス先を決定する。なお、このアクセス先決定機能は、ホスト1が有していてもよい。その場合、チャンク管理部3は、コンテンツを構成する各チャンクに対応する複製チャンクが格納されている記憶装置の情報をホスト1に通知してもよい。

【0073】

障害通知部4は、障害が発生した等によって記憶装置2-1~nの使用できなくなった旨をチャンク管理部3に通知する。なお、障害通知部4は、自ら記憶装置の障害を検出する障害検出手段であってもよいし、あるいはシステム管理者の操作に応じて、障害対応動作指示を示す情報を入力する入力手段(例えば、キー入力手段)であってもよい。

【0074】

複製処理部5は、チャンク管理部3からの指示に応じて、記憶装置2-1~nへのチャンクの追加処理、記憶装置2-1~nからのチャンクの削除処理、または記憶装置2-1~2-n間のチャンクのコピー処理を行う。

【0075】

複製数計画部6は、システムの状況に応じて、各チャンクの複製数を決定する。複製数計画部6は、例えば、各コンテンツの要求量、要求予測量、および可用化数(障害等でも同一内容のチャンクが失われなことを目安に定められた数など)に基づいて、各チャンクの複製数を決定し、複製数が変わった場合にチャンク管理部3にチャンクの管理情報の変更を行わせる手段である。

【0076】

チャンク管理DB7は、コンテンツとチャンクの関連性に関する情報、各チャンクの配置先に関する情報、および配置先とする記憶装置に関する情報を記憶する。具体的には、本システムが扱う各コンテンツの分割管理情報、チャンク毎の複製管理情報、および記憶装置情報を記憶する。

【0077】

なお、チャンク管理部3、複製処理部5、複製数計画部6は、具体的には、CPU等のプログラムに従って動作する情報処理装置によって実現される。チャンク管理部3、複製処理部5、複製数計画部6は、それぞれ異なる情報処理装置によって実現されてもよいし、同一の情報処理装置によって実現されてもよい。なお、通常アクセスに与える影響をなるべく無くすという点で、少なくともチャンク管理部3と他の処理部(複製処理部5や、複製数計画部6)とは、異なる情報処理装置によって実現されることが好ましい。

【0078】

ここで、分割管理情報とは、各コンテンツについて、そのコンテンツをどのようなチャンクに分割したかを示す情報であって、例えば、あるコンテンツを構成する各チャンク(ここでは複製を含まず。)についてのチャンク参照情報を含む。チャンク参照情報は、具体的には、コンテンツ内においてそのチャンクを識別するためのチャンクコンテンツIDと、そのチャンクがコンテンツのどこかのデータに該当するかを示す情報(例えば、先頭からのオフセットやチャンクサイズ等)とを含む情報である。なお、コンテンツ内の順序に従ってチャンク参照情報を登録する場合など、チャンク参照情報のエントリの順番によってチャンクコンテンツIDが特定できる場合には、チャンク参照情報はチャンクコンテンツIDを含んでいなくてもよい。また、チャンクサイズが全てのチャンクで固定である場合には、各チャンクのチャンク参照情報としてチャンクサイズを登録するのではなく、分割管理情報として1つの共通するチャンクサイズを含んでいけばよい。

【0079】

また、複製管理情報とは、各コンテンツ内における各チャンクについて、そのチャンクと同一内容のチャンク(すなわち、そのチャンクの複製チャンク)が、少なくともどの記憶装置に記憶されているかを示す情報であって、例えば、コンテンツ内における各チャンクに対して割り当てたチャンクコンテンツIDと、そのチャンクの複製数と、そのチャンクの各複製チャンクについての配置先情報とを対応づけた情報である。チャンク(複製チャンク)の配置先情報は、例えば、そのチャンクを記憶している記憶装置を識別するため

10

20

30

40

50

の記憶装置IDとその記憶装置内でそのチャンクにアクセスするためのアクセス情報（ブロック番号やオフセット、または各記憶装置内においてチャンクを識別するためのチャンクアクセスIDなど）とを含む情報である。

【0080】

また、記憶装置情報とは、本システムにおいて使用可能な記憶装置を示す情報であって、例えば、本システムが備える記憶装置を識別するための記憶装置IDに対応づけて、その記憶装置の稼働状況（稼働中または停止中）と、記憶容量とを記憶する。

【0081】

次に、チャンク管理部3でのコンテンツの管理方法について説明する。チャンク管理部3において、各コンテンツは、システム内においてそのコンテンツを識別するために定められるコンテンツIDにより一意に識別される。コンテンツ内における各チャンクは、そのコンテンツ内においてそのチャンクを識別するために定められるチャンクコンテンツIDにより一意に識別される。

10

【0082】

また、記憶装置内における各チャンク（複製チャンク）は、その記憶装置におけるファイル名やパス名やinode番号やobjectID等のチャンクアクセスIDにより一意に識別される。また、本発明では、同一内容のチャンクはそれぞれ異なる記憶装置に配置されるので、各複製チャンクを記憶装置IDとチャンクアクセスIDとによって一意に識別することも可能である。ここで、あるチャンクと同一内容の複製チャンクに対して全記憶装置で共通のチャンクアクセスIDが割り当て可能であるならば、チャンクコンテンツIDがチャンクアクセスIDを兼ねることができるともできる。その場合、あるチャンクコンテンツIDで識別されるチャンクの複製チャンクを記憶装置IDにより一意に識別することもできる。

20

【0083】

まず、コンテンツIDとチャンクコンテンツIDとの間に相関をもたせない場合の管理方法を説明する。チャンク管理部3は、管理対象となるコンテンツについて、そのコンテンツの分割管理情報として、図3に示すような分割管理情報をコンテンツIDと対応づけてチャンク管理DB7に記憶させることによって、そのコンテンツをチャンクとして管理する。

【0084】

30

図3は、分割管理情報の例を示す説明図である。図3に示すように、チャンク管理部3は、あるコンテンツの分割管理情報として、そのコンテンツのコンテンツIDと対応づけて、そのコンテンツをチャンクに分割した順に、そのチャンクのチャンク参照情報を記憶してもよい。例えば、コンテンツ内において先頭から1番目のチャンクのチャンク参照情報は、その分割管理情報において1番目のエントリとして登録される。なお、チャンク管理部3は、このような分割管理情報を、コンテンツ毎にチャンク管理DB7に記憶すればよい。

【0085】

ここで、コンテンツ内における各チャンクのサイズが固定長である場合、各チャンクのチャンク参照情報は、そのチャンクのチャンクコンテンツIDを含んでいればよい。例えば、チャンクのサイズがcであるとすると、コンテンツ上のオフセットアドレスがaであるデータが格納されるチャンクは、分割管理情報において $a \div c$ 番目（小数点以下切捨て）のエントリとして登録されているチャンク参照情報に含まれるチャンクコンテンツIDで示されるチャンクであることがわかる。なお、実際のデータの格納先は、そのチャンクコンテンツIDに基づき、そのチャンクの複製管理情報を参照することによって得られる。なお、チャンクサイズは、1つの共通する値として分割管理情報に登録すればよい。

40

【0086】

また、コンテンツ内において各チャンクのサイズがそれぞれ異なる場合、各チャンクのチャンク参照情報は、そのチャンクのチャンクコンテンツIDとチャンクサイズとを含んでいればよい。例えば、コンテンツ内においてj番目のチャンクのサイズが c_j であると

50

すると、コンテンツ上のオフセットアドレスが a であるデータが格納されるチャンクは、以下に示す式 (1) となる i 番目のエントリとして登録されているチャンク参照情報に含まれるチャンクコンテンツ ID で示されるチャンクであることがわかる。

【 0 0 8 7 】

【 数 1 】

$$\sum_{j=0}^i c_j - c_i \leq a < \sum_{j=0}^i c_j \quad \dots \text{式(1)}$$

10

【 0 0 8 8 】

次に、コンテンツ ID とチャンクコンテンツ ID との間に相関をもたせる場合の管理方法を説明する。この方法では、コンテンツ ID とコンテンツ内におけるチャンクの位置とによりそのチャンクのチャンクコンテンツ ID を特定する。従って、各チャンクのチャンク参照情報には、チャンクコンテンツ ID を含めなくてもよい。チャンクコンテンツ ID の構成として、例えば、上位数 b i t をコンテンツ ID とし、残りの下位 b i t をそのチャンクがコンテンツ内において何番目のチャンクかを示す番号とすればよい。

【 0 0 8 9 】

例えば、コンテンツ内における各チャンクのサイズが固定長である場合、各チャンクのチャンク参照情報は省略してもよい。すなわち、分割管理情報を、そのコンテンツにおけるチャンクサイズを格納する表として実現すればよい。なお、全コンテンツでチャンクのサイズが固定長である場合には、分割管理情報をコンテンツ毎に持たせる必要はなく、システムにおけるチャンクサイズを格納する 1 つのデータとして実現すればよい。また、コンテンツ内における各チャンクのサイズが可変長である場合には、各チャンクのチャンク参照情報として、チャンクサイズだけを分割管理情報に登録されればよい。すなわち、分割管理情報を、各チャンクのサイズを格納する表として実現すればよい。なお、コンテンツ上のオフセットアドレスで示されるデータが格納されるチャンクの把握方法は、上記で示したエントリの位置をチャンクの位置としてチャンクコンテンツ ID を特定すればよい。

20

【 0 0 9 0 】

次に、チャンク管理部 3 でのチャンクの配置先の管理方法について説明する。チャンク管理部 3 は、管理対象となるコンテンツを構成する各チャンクについて、図 4 に示すような複製管理情報をチャンクコンテンツ ID と対応づけてチャンク管理 DB 7 に記憶させることによって各チャンクの配置先を管理する。

30

【 0 0 9 1 】

図 4 は、複製管理情報の一例を示す説明図である。図 4 に示すように、チャンク管理部 3 は、あるチャンクの複製管理情報として、そのチャンクのチャンクコンテンツ ID と対応づけて、そのチャンクと同一内容のチャンク (すなわち、複製チャンク) を記憶している記憶装置の記憶装置 ID を記憶すればよい。なお、記憶装置 ID だけでなく、その記憶装置内におけるアドレス情報や記憶装置固有の識別子名 (チャンクアクセス ID) などである。また、複製管理情報は、そのチャンクの複製数を含んでいてもよい。なお、チャンク管理部 3 は、このような複製管理情報を、各コンテンツを構成する各チャンク毎にチャンク管理 DB 7 に記憶すればよい。

40

【 0 0 9 2 】

チャンク管理部 3 は、コンテンツ内における複数のチャンクを連結して 1 つのチャンクとして管理することも可能である。チャンクの連結方法としては、記憶装置内においてチャンクを連結する方法と、コンテンツ内においてチャンクを連結する方法とがある。

【 0 0 9 3 】

記憶装置内においてチャンクを連結する方法とは、同じコンテンツを構成するチャンク

50

のうち同じ記憶装置に記憶される複数のチャンクに対し、1つのチャンクアクセスIDを割り当てる方法である。例えば、コンテンツをチャンクに分割した上で、そのチャンクを記憶装置に記憶させる際に、複製毎に異なるパターンで配置決定した後で同じ記憶装置に入るチャンクをまとめればよい。図5は、チャンクの連結結果の一例を示す説明図である。図5では、あるコンテンツを12個に分割してできるチャンク(チャンク0~11)を複製数=2として各記憶装置に記憶させる場合において、記憶装置ID=0~3に対して1通りの複製チャンクを割り当て、記憶装置ID=10~13に対してもう1通りの複製チャンクを割り当てた場合の例を示している。図5に示す例では、複製毎にチャンクの配置先とする記憶装置と割り当てパターンを替えることで、同一内容のチャンクが異なる記憶装置に格納されるように、かつ各記憶装置間の相関が低くなるように配置先を決定している。このように配置決定された場合、チャンク管理部3は、例えば、記憶装置ID=0を割当先とするチャンク0, 4, 8を連結した連結チャンクを記憶装置内における1つのチャンクとして管理する。ここで、連結チャンクとは、コンテンツを分割してできたチャンクの中から、2以上のチャンクを連結してできたチャンクをいう。同様に、例えば、記憶装置ID=10を割当先とするチャンク0, 7, 10を連結した連結チャンクを記憶装置内における1つのチャンクとして管理する。

10

【0094】

具体的には、複製管理情報に含まれるチャンクの配置先情報として、そのチャンクと同一内容のチャンクについて、そのチャンクを含む連結チャンクを記憶している記憶装置の記憶装置IDと、その記憶装置内でその連結チャンクにアクセスするための情報(チャンクアクセスID)と、そのチャンクがその連結チャンクのどこに位置するかを示す情報(例えば、連結順序を示す情報やオフセットアドレス等)とを記憶すればよい。このように、記憶装置に記憶させる際にチャンクを連結する方法を用いれば、記憶装置内におけるチャンク数とコンテンツの分割数とを切り離して管理することができ、管理すべきチャンクアクセスIDの量を減らすことができる。

20

【0095】

また、コンテンツ内においてチャンクを連結する方法とは、コンテンツ内における非連続なチャンクに対し、1つのチャンクコンテンツIDを割り当てる方法である。例えば、コンテンツをチャンクに分割した後で、巡回的なパターンで配置決定された場合に同じ記憶装置に入るチャンクを1つにまとめてもよい。図6は、チャンクの連結結果の一例を示す説明図である。図6では、あるコンテンツを12個に分割してできるチャンク(チャンク0~11)に対し、3を区切りに巡回させた場合の例を示している。このような場合、チャンク管理部3は、例えば、記憶装置ID=1が配置先となるチャンク0, 3, 9を連結した連結チャンクをコンテンツ内における1つのチャンクとして管理する。同様に、例えば、記憶装置ID=2が配置先となるチャンク1, 4, 10を連結した連結チャンクをコンテンツ内における1つのチャンクとして管理する。

30

【0096】

具体的には、分割管理情報として、コンテンツ内における連結チャンク数分のチャンク参照情報を記憶するようにし、各連結チャンクのチャンク参照情報に、各連結チャンクに対し割り当てたチャンクコンテンツIDと、その連結チャンクを構成する各チャンクの位置情報(コンテンツ上のオフセットアドレスとサイズ等)とを含めればよい。このように、コンテンツ内においてチャンクを連結する方法を用いれば、チャンクサイズの決定に際しデータの連続性を考慮しなくてもよいので、例えば、そのコンテンツに対し1度にアクセスされる可能性の高いデータサイズ分よりも大きいサイズを最終的なチャンクサイズとすることも可能である。また、システムで管理するチャンクコンテンツIDおよびチャンクアクセスIDの量を減らすこともできる。

40

【0097】

次に、本実施の形態によるデータ分割管理システムの動作について説明する。まず、ホスト1からコンテンツを読み出す動作について説明する。図7は、本実施の形態によるデータ分割管理システムのコンテンツ読み出し動作の一例を示すフローチャートである。本

50

実施の形態では、読み出しを行うホスト1が、チャンク管理部3に対し、読み出したいコンテンツのコンテンツIDとコンテンツ内の読み出し開始位置とサイズとを指定してアドレス変換を要求し、その応答として、そのコンテンツの該当部分を構成するチャンクの格納先を示す情報を得て、その情報に基づき記憶装置にアクセスする場合を例にとって説明する。

【0098】

ホスト1で読み込みを行う際には、ホスト1は、チャンク管理部3に、コンテンツIDとコンテンツ内の読み出し開始位置を示すオフセットアドレスとサイズと読み出しであることを指定して、アドレス変換を依頼する(ステップS100)。ホスト1は、例えば、通信ネットワークを介してチャンク管理部3を備えたサーバ装置に接続されている場合には、そのサーバ装置に、コンテンツIDとコンテンツ内オフセットアドレスと読み出しサイズとを含む読み出し用アドレス変換要求メッセージを送信すればよい。そして、チャンク管理部3からの応答を待つ(ステップS101)。なお、チャンク管理部3におけるアドレス変換動作については後述する。

10

【0099】

応答が帰ってきたら、ホスト1は、その応答の先頭のエントリを読み出し対象エントリとする(ステップS102)。ここで、読み出し用のアドレス変換要求に対するチャンク管理部からの応答のフォーマットは、読み出し対象のデータを構成するチャンク毎に1つのエントリとして格納されているものとする。また、各エントリは、例えば、そのチャンクを格納している記憶装置の記憶装置IDと、記憶装置内でのアクセス情報(例えば、チャンクアクセスID)と、そのチャンク内での読み出し開始位置と、そのチャンク内での読み出しサイズとで構成される。また、ここでは、各エントリが、応答のデータの一部に、そのコンテンツにおけるオフセットアドレスの順に格納されている場合を例にとって説明する。なお、複数のチャンクを連結して1つのチャンクとする場合には、連結前のチャンク毎に1つのエントリとしてもよい。

20

【0100】

図8は、読み出し用アドレス変換要求に対する応答フォーマットの一例を示す説明図である。図8に示す例では、読み出し対象のデータを構成するチャンク毎に、1つのエントリとして、そのチャンクと同一内容のチャンク(すなわち、複製チャンク)のいずれかの格納先である記憶装置の記憶装置IDと、その記憶装置内でのチャンクアクセスIDと、そのチャンク内での読み出し開始位置(チャンク内オフセットアドレス)とそのチャンク内での読み出しサイズとを格納する例を示している。図8に示す例では、まず、エントリ1が、読み出し対象エントリとしてホスト1によって処理される。

30

【0101】

ホスト1は、読み出し対象エントリが示すチャンク(以下、単に対象チャンクという。)が未割り当ての場合にはステップS106に移行し、割り当てられていた場合にはステップS104に移行する(ステップS103)。ここで、チャンクが未割り当てとは、そのチャンクが記憶装置に書き込まれていないことを示している。なお、対象チャンクが未割り当てか否かは、例えば、記憶装置IDやチャンクアクセスIDに無効な値が格納されているか否かで判断すればよい。

40

【0102】

対象チャンクが割り当てられていた場合、ホスト1は、読み出し対象エントリ中の記憶装置IDで特定される記憶装置2に対して、読み出し対象エントリ中のチャンクアクセスID、チャンク内オフセットアドレス、読み出しサイズを指定し読み出しを要求する(ステップS104)。そして、記憶装置2からの応答を待つ(ステップS105)。

【0103】

ホスト1は、チャンク管理部3から通知された全エントリ(全チャンク)について、上記読み出し処理を行う。すなわち、記憶装置2から応答が帰ってきたら、チャンク管理部3からのアドレス変換応答に含まれる全エントリで示されるチャンクの読み出しが全て完了したか否かを判定し、全て完了していたら処理を終了し、完了していなければステップ

50

S 1 0 7に移行する(ステップS 1 0 6)。ステップS 1 0 7では、まだ読み出しを行っていないチャンクを読み出すために、次のエントリを読み出し対象エントリとし、ステップS 1 0 3に移行する。

【 0 1 0 4 】

図 8 に示す例では、エントリ 1 の次に、エントリ 2 が読み出し対象エントリとされ、最終的にエントリ n までが、読み出し対象エントリとしてホスト 1 によって処理されることとなる。本例では、各エントリがそのコンテンツにおけるオフセットアドレスの順に格納されているので、ホスト 1 は、読み出し対象エントリに従い、各チャンクの該当データを読み出して順次結合していけば、所望のコンテンツを得ることができる。なお、各エントリを、そのコンテンツにおけるオフセットアドレスの順に格納しない場合には、各エントリに、そのチャンクがコンテンツにおけるどの位置のデータであるかを示す情報を含めればよい。

10

【 0 1 0 5 】

ここで、各チャンクをシーケンシャルに読み込む例を示したが、応答を待たずに並行に処理することも可能である。このようにすると処理時間を短縮することができる。なお、後述の配置先決定動作において説明するように、連続したチャンクをそれぞれ異なる記憶装置に配置させている場合には、ホスト 1 側で特別に意識しなくても、各記憶装置に並列にアクセスできるので、容易に処理時間を短縮することができる。

【 0 1 0 6 】

次に、ホスト 1 からコンテンツを書き込む動作について説明する。図 9 は、本データ分割管理システムのコンテンツ書き込み動作の一例を示すフローチャートである。本データ分割管理システムでは、書き込みを行うホスト 1 が、チャンク管理部 3 に対し、書き込みたいコンテンツとコンテンツ内の書き込み開始位置とサイズとを指定してアドレス変換を要求し、その応答として、そのコンテンツの該当部分を構成するチャンク(およびその複製)の格納先を示す情報を得て、その情報に基づき記憶装置にアクセスする場合を例にとって説明する。

20

【 0 1 0 7 】

ホスト 1 で書き込みを行う際には、ホスト 1 は、チャンク管理部 3 に、コンテンツ ID とコンテンツ内の書き込み開始位置を示すオフセットアドレスと書き込みサイズと書き込みであることを指定して、アドレス変換を依頼する(ステップS 2 0 0)。ホスト 1 は、例えば、通信ネットワークを介してチャンク管理部 3 を備えたサーバ装置に接続されている場合には、そのサーバ装置に、コンテンツ ID とコンテンツ内オフセットアドレスと書き込みサイズとを含む書き込み用アドレス変換要求メッセージを送信すればよい。そして、チャンク管理部 3 からの応答を待つ(ステップS 2 0 1)。なお、チャンク管理部 3 が、新規なコンテンツに対するコンテンツ ID の割り当てを行ってもよい。そのような場合、ホスト 1 は未割り当てを示すコンテンツ ID を指定すればよい。また、チャンク管理部 3 におけるアドレス変換動作については後述する。なお、本発明において、同一内容のチャンクは、チャンク管理部 3 によって異なる記憶装置に配置されるべく配置先決定される。

30

【 0 1 0 8 】

応答が帰ってきたら、ホスト 1 は、その応答の先頭のエントリを書き込み対象エントリとする(ステップS 2 0 2)。ここで、書き込み用のアドレス変換要求に対するチャンク管理部 3 からの応答のフォーマットは、書き込み対象のデータを構成するチャンク毎に 1 つのエントリとして格納されているものとする。また、各エントリは、例えば、そのチャンク内での書き込みサイズと、そのチャンク内での書き込み開始位置と、そのチャンクの複製数に応じた各複製チャンクの格納先を示す情報とで格納される。複製チャンクの格納先を示す情報は、具体的には、その複製チャンクを格納する記憶装置の記憶装置 ID と、その記憶装置内でのチャンクアクセス ID でよい。また、本例では、各エントリが、応答のデータの一部に、そのコンテンツにおけるオフセットアドレスの順に格納されている場合を例にとって説明する。なお、複数のチャンクを連結して 1 つのチャンクとする場合に

40

50

は、連結前のチャンク毎に1つのエントリとしてもよい。

【0109】

図10は、書き込み用アドレス変換要求に対する応答フォーマットの一例を示す説明図である。図10に示す例では、書き込み対象のデータを構成するチャンク毎に1つのエントリとし、そのチャンク内での書き込みサイズと、そのチャンク内での書き込み開始位置と、そのチャンクの複製数と、そのチャンクの配置先情報として、そのチャンクの複製数に応じた、各複製チャンクの格納先となる記憶装置の記憶装置IDとその記憶装置内でのアクセス情報(チャンクアクセスID)とを格納する例を示している。図10に示す例では、まず、エントリ1が、書き込み対象エントリとしてホスト1によって処理される。なお、全記憶装置で共通のチャンクアクセスIDが定義されている場合には1つのエントリにつき1つのチャンクアクセスIDが格納されていけばよい。また、同じ内容のチャンクの複製チャンクであっても、各複製チャンクでチャンク内での書き込み開始位置が異なる場合(例えば、記憶装置内においてチャンクを連結する場合)には、複製チャンク毎にその複製チャンク内での書き込み開始位置(チャンク内オフセットアドレス)を格納すればよい。

10

【0110】

次に、ホスト1は、書き込み対象エントリ中に、複製チャンクの書き込み先として先頭に示された記憶装置IDで特定される記憶装置2を書き込み対象記憶装置として選ぶ(ステップS203)。以下、書き込み対象記憶装置に書き込むチャンクのことを対象複製チャンクという。

20

【0111】

書き込み対象記憶装置を特定すると、ホスト1は、書き込み対象エントリ中の書き込みサイズと、チャンク内オフセットアドレスと、対象複製チャンクについてのチャンクアクセスIDとを指定し、その対象複製チャンクとして書き込むべきデータとともに、書き込み対象記憶装置に対して書き込みを要求する(ステップS204)。なお、本例では、各エントリがそのコンテンツにおけるオフセットアドレスの順に格納されているので、ホスト1は、書き込むべきデータとして、エントリ毎に、コンテンツ内における書き込み開始位置からそのチャンク内での書き込みサイズ分のデータを順次指定していけばよい。そして、記憶装置2からの応答を待つ(ステップS205)。

【0112】

ホスト1は、記憶装置2から応答が帰ってきたら、書き込み対象エントリ中に各複製チャンクの格納先として示された全記憶装置に対し書き込みが完了したか否かを判定し、全て完了していればステップS208へ、完了していなければステップS207へ移行する(ステップS206)。ステップS207では、まだ書き込みを行っていない記憶装置に書き込むために、現書き込み対象エントリ中に、複製チャンクの格納先として次に示された記憶装置IDで特定される記憶装置2を書き込み対象記憶装置とし、ステップS204に移行する。

30

【0113】

また、ステップS208では、チャンク管理部3からのアドレス変換応答に含まれる全エントリについて書き込み処理が完了していたか否かを判定し、全て完了していたら処理を終了し、完了していなければステップS209に移行する。ステップS209では、まだ書き込みを行っていないエントリについて書き込み処理を行うために、次のエントリを書き込み対象エントリとし、ステップS203に移行する。図10に示す例では、エントリ1で示される全複製チャンクの書き込み処理が完了した後に、エントリ2が書き込み対象チャンクとされ、エントリ2で示される各複製チャンクの書き込み処理が行われることとなる。最終的に、エントリnまでが書き込み対象エントリとされ、エントリnで示される各複製チャンクの書き込み処理が行われることとなる。

40

【0114】

このようにして、同じ内容のチャンクを複数の記憶装置に書き込むことができるので、データが冗長化され、データの可用性および安全性が向上する。また、書き込み動作につ

50

いても、各チャンクを、応答を待たずに並列に処理することが可能である。このようにすると処理時間を短縮することができる。また、ホスト側で特別に意識しなくても、各記憶装置に並列にアクセスできるので、容易に処理時間を短縮することができる。なお、後述の配置先決定動作において説明するように、同一内容のチャンクが異なる記憶装置に格納されるよう配置先決定されるので、その配置先に基づいて各記憶装置にチャンクを書き込むことにより、読み出す際にアクセスする記憶装置を分散させることができ、同時に読み出す際のアクセス性（スループットや応答時間）も向上する。

【0115】

なお、上記例では、コンテンツ内におけるチャンク毎に1つのエン트리として、その中に各複製チャンクの格納先についての情報を格納する応答フォーマットの例を示したが、これに限らず、例えば、複製チャンク毎に1つのエン트리として、その中に、格納先についての情報と、その複製チャンクがコンテンツにおけるどの位置のデータであるかを示す情報とを格納させてもよい。そのような場合には、ホスト1は、各エントリで示される複製チャンクについての書き込み処理を順次行っていけばよい。

10

【0116】

記憶装置に対するチャンクのアクセス処理は、ブロックデバイスへのアクセス処理や、NASにおけるファイルへのアクセス処理や、OSDにおけるオブジェクトへのアクセス処理など一般的なアクセス処理である。例えば、記憶装置がブロックデバイスによって実現される場合には、記憶装置内のチャンクに対応する、ブロックデバイス中の開始オフセットアドレスとサイズを指定してアクセスすればよい。記憶装置内のチャンクに対応する開始オフセットアドレスは、チャンクが固定長であれば、何番目のチャンクであるかによって算出すればよい。また、可変長の場合には、ブロックデバイス内のチャンクごとのオフセットアドレスやサイズのリストを定義しておき、これらに基づき算出すればよい。なお、チャンクの一部に対しアクセスする場合には、ブロックデバイス内のオフセットアドレスにチャンク内オフセットアドレスを加算すればよい。

20

【0117】

また、例えば、記憶装置がNAS等のファイルベースの記憶システムによって実現される場合には、記憶装置内のチャンクに対応するファイルを指定してアクセスすればよい。なお、チャンクの一部に対しアクセスする場合には、チャンクに対応するファイル内のオフセットアドレスにチャンク内オフセットアドレスを指定してアクセスすればよい。また、例えば、記憶装置がOSDによって実現される場合には、記憶装置内のチャンクに対応するオブジェクトのリードまたはライト処理を実行すればよい。チャンクの一部に対しアクセスする場合には、オブジェクト内のLBAにチャンク内オフセットアドレスを指定したリードまたはライト処理を実行すればよい。なおNASやOSDにおいて書き込み時に対象チャンクが存在していなかった場合には、ファイルやオブジェクトの作成処理も同時に行う。

30

【0118】

次に、チャンク管理部3でのアドレス変換処理について説明する。既に説明したように、アドレス変換処理は、コンテンツIDとコンテンツ内の読み出し/書き込み開始位置と読み出し/書き込みサイズを指定して依頼される。

40

【0119】

図11は、チャンク管理部3におけるアドレス変換処理の処理フローの一例を示すフローチャートである。チャンク管理部3は、まず、書き込みか読み出しのどちらが指示されたか調べ（ステップS301）、書き込みであれば書き込みアドレス変換処理（ステップS302）を行い、読み出しであれば読み出しアドレス変換処理（ステップS303）を行い、応答を返し（ステップS304）、終了する。

【0120】

図12は、読み出しアドレス変換処理の処理フローの一例を示すフローチャートである。図12に示すように、チャンク管理部3は、まず、アドレス変換の対象として指定されたコンテンツIDと対応づけられた分割管理情報を読み込む（ステップS401）。チャ

50

ンク管理部 3 は、例えば、指定されたコンテンツ ID と対応づけられた分割管理情報をチャンク管理 DB 7 から検索して読み込めばよい。ここで、チャンク管理部 3 は、指定されたコンテンツ ID と対応づけられた分割管理情報がチャンク管理 DB 7 中に存在していればステップ S 4 0 4 に移行し、存在していなければステップ S 4 0 3 に移行する（ステップ S 4 0 2）。

【 0 1 2 1 】

ステップ S 4 0 3 では、チャンク管理部 3 は、チャンクが未割り当てであることを示す値を応答のエントリに設定し、処理を終了する。

【 0 1 2 2 】

また、ステップ S 4 0 4 では、チャンク管理部 3 は、指定された読み出し開始位置と読み出しサイズ、およびステップ S 4 0 1 で読み込んだ分割管理情報によって特定される各チャンクのコンテンツ内におけるオフセットアドレスとチャンクサイズとから、アドレス変換対象領域をチャンクに分割する。チャンク管理部 3 は、例えば、読み込んだ分割管理情報に含まれる各チャンクのチャンク参照情報群に基づいて、各チャンクのオフセットアドレスとサイズとを特定し、指定された読み出し開始位置と読み出しサイズ、および特定した各チャンクのコンテンツ内におけるオフセットアドレスとチャンクサイズから、読み出し対象として指定されたデータがそのコンテンツ内においてどのチャンクに該当するかを算出することによって、アドレス変換の対象とするデータに該当するチャンクの範囲をチャンクコンテンツ ID により特定する。例えば、読み出し対象データが、コンテンツ内のオフセットアドレスの順でみた場合に、何番目のチャンクから何番目のチャンクまでの情報によって構成されているかを計算することによって、変換対象とするチャンクの範囲を特定すればよい。

【 0 1 2 3 】

変換対象とするチャンクの範囲の特定方法として、チャンクが固定長である場合には、指定された読み出し開始位置 a と読み出しサイズ s 、およびチャンクサイズ c とに基づいて、先頭のチャンクおよび末尾のチャンクが何番目（0 基準）かは、先頭チャンクを e 番目、末尾チャンクを f 番目とすると、それぞれ次のように求まる。

【 0 1 2 4 】

$$e = a \div c \text{ (小数点以下切り捨て)} \quad \dots \text{式(2)}$$

$$f = (a + s) \div c \text{ (小数点以下切り捨て)} \quad \dots \text{式(3)}$$

【 0 1 2 5 】

また、 e 番目の先頭チャンク中の読み出し開始位置 g 、および f 番目の末尾チャンク中の読み出し終了位置 h は、次のように求まる。

【 0 1 2 6 】

$$g = a - c (e - 1) \quad \dots \text{式(4)}$$

$$h = (a + s) - c (f - 1) \quad \dots \text{式(5)}$$

【 0 1 2 7 】

また、チャンクサイズが可変長の場合、 j 番目のチャンクサイズが c_j であるとする、先頭のチャンクが何番目かは以下の式(6)を満たす e によって求まる。また、末尾のチャンクが何番目かは以下の式(7)を満たす f によって求まる。

【 0 1 2 8 】

【 数 2 】

$$\sum_{j=0}^e c_j - c_e \leq a < \sum_{j=0}^e c_j \quad \dots \text{式(6)}$$

$$\sum_{j=0}^f c_j - c_f \leq a + s < \sum_{j=0}^f c_j \quad \dots \text{式(7)}$$

10

20

30

40

【 0 1 2 9 】

また、e 番目の先頭チャンク中の読み出し開始位置 g、および f 番目の末尾チャンク中の読み出し終了位置 h は、以下の式 (8) および式 (9) によって求まる。

【 0 1 3 0 】

【数 3】

$$g = a - \sum_{j=0}^e c_j - c_e \quad \cdots \text{式(8)}$$

$$h = a + s - \sum_{j=0}^f c_j - c_f \quad \cdots \text{式(9)}$$

10

【 0 1 3 1 】

変換対象とするチャンクの範囲を特定すると、チャンク管理部 3 は、先頭のチャンクを対象チャンクとして設定するとともに、その対象チャンクについてのエントリを応答の先頭のエントリに設定する (ステップ S 4 0 5)。そして、対象チャンクの複製管理情報を読み込む (ステップ S 4 0 6)。チャンク管理部 3 は、例えば、対象チャンクのチャンクコンテンツ ID と対応づけられた複製管理情報をチャンク管理 DB 7 中から検索して読み込めばよい。ここで、チャンク管理部 3 は、分割管理情報において対象チャンクにチャンクコンテンツ ID が割り当てられており、さらに複製管理情報において、対象チャンクの複製チャンク (対象チャンクと同一内容の記憶装置におけるチャンク) が記憶装置に割り当てられていればステップ S 4 0 9 に移行し、そうでなければステップ S 4 0 8 に移行する (ステップ S 4 0 7)。チャンク管理部 3 は、例えば、取得した複製管理情報において、対象チャンクのチャンクコンテンツ ID に対し、少なくとも 1 つ以上の記憶装置 ID やチャンクアクセス ID が割り当てられているかを確認すればよい。

20

【 0 1 3 2 】

ステップ S 4 0 8 では、対象チャンクについての応答のエントリに、チャンクが未割り当てであることを示す情報を設定し、ステップ S 4 1 2 に移行する。

【 0 1 3 3 】

また、ステップ S 4 0 9 では、対象チャンクの複製管理情報に含まれるチャンクの配置先情報を参照し、そして、ステップ S 4 0 9 で配置先情報に基づいて、ホスト 1 にアクセスさせる記憶装置を決定する (ステップ S 4 1 0)。

30

【 0 1 3 4 】

アクセス先の決定としては、例えば、乱数により決定する方法がある。例えば、所定の乱数を発生させ、チャンクの格納先の数 (複製数) で除算した余りに応じてアクセス先を決定してもよい。また、例えば、複製チャンクを記憶している各記憶装置を順番に使用させるラウンドロビン方式によってアクセス先を決定してもよい。また、例えば、複製チャンクを記憶している各記憶装置のその時点の負荷から、最も負荷が低い記憶装置をアクセス先として決定してもよい。なお、各記憶装置の負荷は、例えば、チャンク管理部 3 や、各記憶装置が有する制御装置が、単位時間当たりの記憶装置に対する I/O 数や転送データ量、CPU 利用率を検出することによって、判断すればよい。

40

【 0 1 3 5 】

アクセス先が決定すると、チャンク管理部 3 は、そのアクセス先を示す情報を応答の該当エントリに設定する (ステップ S 4 1 1)。例えば、チャンク管理部 3 は、対象チャンクについての応答のエントリに、アクセス先として決定した記憶装置の記憶装置 ID と、必要に応じてその記憶装置内でのアクセス情報 (チャンクアクセス ID) とを設定する。

【 0 1 3 6 】

次に、チャンク管理部 3 は、ステップ S 4 1 2 において、対象チャンクについての応答のエントリに、そのチャンク内での読み出し開始位置 (チャンク内オフセットアドレス)

50

と読み出しサイズとを設定する。なお、この動作は、対象チャンクが未割り当てであった場合にも行う。

【 0 1 3 7 】

チャンク内オフセットアドレスおよび読み出しサイズは、例えば、既に説明した変換対象とするチャンクの範囲の把握動作において求めたチャンクサイズ c (または c_j) や先頭チャンクの読み出し開始位置 g と末尾チャンクの読み出し終了位置 h とによって求められる。

【 0 1 3 8 】

具体的には、対象チャンクが先頭チャンクである場合、チャンク内オフセットアドレスは g , 読み出しサイズは $c - g$ となる。ここで、 c は対象チャンクのチャンクサイズを示す。また、対象チャンクが末尾チャンクである場合、チャンク内オフセットアドレスは 0 , 読み出しサイズは h となる。また、対象チャンクが先頭チャンクおよび末尾チャンク以外のチャンクである場合、チャンク内オフセットアドレスは 0 , 読み出しサイズは c となる。なお、先頭チャンク = 末尾チャンク ($e = f$) の場合、チャンク内オフセットアドレスは g , 読み出しサイズは $h - g$ となる。

10

【 0 1 3 9 】

なお、複数のチャンクを連結して1つのチャンクとする場合であって、連結前のチャンク毎に1つのエントリとする場合には、各複製チャンクのチャンク内オフセットアドレスには、連結チャンクにおける対象チャンクの位置を加味する必要がある。

【 0 1 4 0 】

20

対象チャンクについてのアドレス変換 (エントリの設定) が終了すると、チャンク管理部 3 は、変換対象とするチャンクの範囲内でアドレス変換が終わっていないチャンクがなければ処理を終了し、あればステップ S 4 1 4 に移行する (ステップ S 4 1 3)。ステップ S 4 1 4 では、対象チャンクを次のチャンクに設定し、ステップ S 4 0 6 に移行する。

【 0 1 4 1 】

なお、ステップ S 4 0 1 において、チャンクサイズの取得は、チャンクサイズがシステムに対し予め定められている場合には省略してもよい。また、コンテンツ毎に可変である場合には、例えば、そのコンテンツの分割管理情報に含まれているチャンクサイズを参照すればよい。また、各チャンク毎に可変である場合には、例えば、各チャンクのチャンク参照情報に含まれているチャンクサイズを参照すればよい。

30

【 0 1 4 2 】

また、本例では、ステップ S 4 0 3 において、全対象データが未割り当てであった場合には、全体を1つのチャンクとして未割り当てを示すエントリの応答として返す例を示している。

【 0 1 4 3 】

以上のように、同一内容のチャンクを複数の記憶装置に記憶した上で、アクセス先を分散させることによって、複数のホストから同一のチャンクが読み出されるような場合であっても、読み出し負荷を分散させることができる。

【 0 1 4 4 】

なお、本例では、チャンク管理部 3 がアクセス先を決定する例を示したが、チャンク管理部 3 では、1つのエントリに各複製チャンクの格納先を示す情報群を格納してホスト 1 に返信するようにし、ホスト 1 側でアクセス先とする記憶装置を決定してもよい。なお、ホスト側におけるアクセス先の決定方法としては、例えば、乱数を用いたり、ホスト 1 に予め定められた値を用いて該当チャンクを記憶した複数の記憶装置から1つの記憶装置を選べばよい。

40

【 0 1 4 5 】

例えば、チャンク管理部 3 は、ステップ S 4 1 0 においてアクセス先を決定せずに、対象チャンクの複製管理情報で示されるその対象チャンクの全配置先を、その対象チャンクについての応答のエントリに設定する。図 1 3 および図 1 4 は、ホスト 1 側でアクセス先とする記憶装置を決定する場合の応答フォーマットの例を示す説明図である。チャンク管

50

理部 3 は、例えば、図 1 3 に示すように、1 つの対象チャンクにつき 1 つのエントリとして、そのチャンク内での読み出し開始位置（チャンク内オフセットアドレス）と読み出しサイズと複製数と、全配置先情報として複製数分の記憶装置 ID およびチャンクアクセス ID とを格納すればよい。

【 0 1 4 6 】

なお、図 1 4 は、記憶装置内のチャンクアクセス ID が記憶装置間で共通である場合の応答フォーマットの例である。図 1 4 に示す例では、1 つの対象チャンクにつき 1 つのエントリとして、そのチャンク内での読み出し開始位置（チャンク内オフセットアドレス）と読み出しサイズと複製数とそのチャンク的全複製チャンクに対し割り当てられた共通のチャンクアクセス ID と、複製数分の記憶装置 ID とを格納する例を示している。また、
10
同じ内容のチャンクであっても、記憶装置内においてそのチャンク内での読み出し開始位置が異なる場合（例えば、記憶装置内においてチャンクを連結する場合）には、複製チャンク毎にその記憶装置内におけるチャンク内での読み出し開始位置（チャンク内オフセットアドレス）を格納すればよい。

【 0 1 4 7 】

そして、ホスト 1 は、ステップ S 1 0 4 において記憶装置 2 に読み出しを要求する前に、読み出し対象エントリ中に全配置先情報として示されている複数の記憶装置 ID から 1 つの記憶装置 ID を選び出すことによって、アクセス先とする記憶装置を決定する。この際、ホスト 1 は、乱数あるいは、ホストに規定された数字を用いてもよい。

【 0 1 4 8 】

また、図 1 5 は、書き込みアドレス変換処理の処理フローの一例を示すフローチャートである。図 1 5 に示すように、チャンク管理部 3 は、まず、アドレス変換の対象として指定されたコンテンツ ID と対応づけられた分割管理情報を読み込む（ステップ S 5 0 1）。チャンク管理部 3 は、例えば、指定されたコンテンツ ID と対応づけられた分割管理情報をチャンク管理 DB 7 から検索して読み込めばよい。ここで、チャンク管理部 3 は、指定されたコンテンツ ID と対応づけられた分割管理情報がチャンク管理 DB 7 中に存在していればステップ S 5 1 1 に移行し、存在していなければステップ S 5 0 3 に移行する（ステップ S 5 0 2）。なお、未割り当てを示すコンテンツ ID が指定された場合、チャンク管理部 3 は、新たにコンテンツ ID を割り当てて、ステップ S 5 0 3 に移行すればよい。
20
30

【 0 1 4 9 】

ステップ S 5 0 3 では、新規コンテンツ用に、チャンク管理 DB 7 中に、そのコンテンツ ID と対応づけた分割管理情報を記憶するための領域を確保し、その分割管理情報を初期化する。チャンク管理部 3 は、例えば、チャンクサイズを決定し、決定したチャンクサイズによって定まるチャンクの分割数に応じて、各チャンクのチャンク参照情報を生成し、分割管理情報として登録することによって初期化すればよい。チャンクサイズは、システムとして規定値を持たせてもよいし、コンテンツ毎に設定できるようにしてもよいし、コンテンツ内のチャンク毎に設定できるようにしてもよい。また、1 回の書き込み毎にチャンクを分割してその都度可変とする方法もある。なお、各チャンクのチャンク参照情報には、チャンクサイズに応じて算出されるチャンクのオフセットアドレス等を登録しても
40
よい。なお、各チャンクのチャンクコンテンツ ID は、そのチャンクの複製管理情報が作成されるまでに登録されていればよく、この時点では未割り当てを示す情報を登録してもよいし、ここで割り当ててその値を登録してもよい。

【 0 1 5 0 】

次いで、チャンク管理部 3 は、指定された書き込み開始位置と書き込みサイズ、およびステップ S 5 0 3 の初期化の際に決定されたチャンクサイズならびにチャンク参照情報から、アドレス変換対象領域をチャンクに分割する（ステップ S 5 0 4）。具体的には、チャンク管理部 3 は、指定された書き込み開始位置と書き込みサイズ、および各チャンクのコンテンツ内におけるオフセットアドレスとチャンクサイズとから、書き込み対象として指定されたデータがそのコンテンツ内においてどのチャンクに該当するかを算出すること
50

によって、アドレス変換の対象とするチャンクの範囲をチャンクコンテンツIDにより特定する。

【0151】

アドレス変換の対象とするチャンクの範囲を特定すると、チャンク管理部3は、先頭のチャンクを対象チャンクとして設定するとともに、その対象チャンクについてのエントリを応答の先頭のエントリに設定する(ステップS505)。そして、対象チャンクの複製チャンク(対象チャンクと同一内容の記憶装置におけるチャンク)を記憶装置に割り当てるチャンク作成処理を行う(ステップS506)。なお、チャンク作成処理の詳細については後述する。

【0152】

チャンク作成処理が完了すると、対象チャンクの配置先情報を、応答の該当エントリに設定する(ステップS507)。チャンク管理部3は、例えば、チャンク作成処理によってチャンク管理DB7中に作成された、対象チャンクのチャンクコンテンツIDと対応づけられた複製管理情報を読み込み、その複製管理情報で配置先として示される複数の記憶装置IDとアクセス情報(チャンクアクセスID)とを応答のエントリに設定する。また、チャンク管理部3は、対象チャンクについての応答のエントリに、そのチャンク内での書き込み開始位置(チャンク内オフセットアドレス)と書き込みサイズとを設定する(ステップS508)。なお、チャンク内オフセットアドレスおよび書き込みサイズの算出方法は、読み出しアドレス変換処理におけるチャンク内オフセットアドレスおよび読み出しサイズの算出方法と同様である。

【0153】

対象チャンクについてのアドレス変換(エントリの設定)が終了すると、チャンク管理部3は、変換対象とするチャンクの範囲内でアドレス変換が終わっていないチャンクがなければ処理を終了し、あればステップS510に移行する(ステップS509)。ステップS510では、対象チャンクを次のチャンクに設定し、ステップS506に移行する。

【0154】

また、チャンク管理部3は、指定されたコンテンツIDと対応づけられた分割管理情報が存在していた場合には(ステップS502のYes)、指定された書き込み開始位置と書き込みサイズ、および読み込んだ分割管理情報によって特定される各チャンクのコンテンツ内におけるオフセットアドレスとチャンクサイズとから、アドレス変換対象領域をチャンクに分割する(ステップS511)。具体的には、チャンク管理部3は、ステップS504と同様に、書き込み対象として指定されたデータがそのコンテンツ内においてどのチャンクに該当するかを算出することによって、アドレス変換の対象とするチャンクの範囲をチャンクコンテンツIDにより特定する。

【0155】

アドレス変換の対象とするチャンクの範囲を特定すると、チャンク管理部3は、先頭のチャンクを対象チャンクとして設定するとともに、その対象チャンクについてのエントリを応答の先頭のエントリに設定する(ステップS512)。

【0156】

対象チャンクが設定されると、チャンク管理部3は、まず、対象チャンクの複製管理情報を読み込む(ステップS513)。ここで、チャンク管理部3は、分割管理情報において、対象チャンクにチャンクコンテンツIDが割り当てられており、さらに複製管理情報において、そのチャンクに対し配置先である記憶装置が割り当てられていればステップS516に移行し、そうでなければステップS515に移行する(ステップS514)。

【0157】

ステップS515では、対象チャンクに対し、配置先として複製数分の記憶装置を割り当てるチャンク作成処理を行う。チャンク作成処理は、ステップS506と同様である。なお、チャンク管理部3は、チャンク作成処理が完了すると、チャンク作成処理によって作成された対象チャンクの複製管理情報を読み込む。

【0158】

10

20

30

40

50

また、ステップ S 5 1 6 では、チャンク管理 D B 7 に登録されている対象チャンクの複製管理情報を読み込む。チャンク管理部 3 は、例えば、対象チャンクのチャンクコンテンツ ID と対応づけられた複製管理情報をチャンク管理 D B 7 中から検索して読み込めばよい。

【 0 1 5 9 】

そして、対象チャンクの複製管理情報に含まれるチャンクの配置先情報を参照して（ステップ S 5 1 7）、配置先として示されている複数の記憶装置 ID とアクセス情報（チャンクアクセス ID）とを、応答のエントリに設定する（ステップ S 5 1 8）。また、チャンク管理部 3 は、対象チャンクについての応答のエントリに、そのチャンク内での書き込み開始位置（チャンク内オフセットアドレス）と書き込みサイズとを設定する（ステップ S 5 1 9）。ここで、複製管理情報において既に対象チャンクの配置先として記憶装置が割り当てられていた場合には、再書き込みを行わせないようにするために、応答のエントリに配置先情報を設定しないようにすることも可能である。なお、この場合においても、ホスト 1 に次のチャンクのオフセットアドレスを知らせるために書き込みサイズは設定しておく。

10

【 0 1 6 0 】

対象チャンクについてのアドレス変換（エントリの設定）が終了すると、チャンク管理部 3 は、変換対象とするチャンクの範囲内でアドレス変換が終わっていないチャンクがなければ処理を終了し、あればステップ S 5 2 0 に移行する（ステップ S 5 1 9）。ステップ S 5 3 0 では、対象チャンクを次のチャンクに設定し、ステップ S 5 1 3 に移行する。

20

【 0 1 6 1 】

ここで、ステップ S 5 0 4、S 5 1 1 における変換対象とするチャンクの範囲の特定方法は、読み出しアドレス変換処理における方法と同様である。また、ステップ S 5 0 8、S 5 1 8 で設定するチャンク内オフセットアドレスおよび書き込みサイズの算出方法についても、読み出しアドレス変換処理における算出方法と同様である。

【 0 1 6 2 】

次に、チャンク作成処理について説明する。図 1 6 は、チャンク管理部 3 におけるチャンク作成処理の処理フローの一例を示すフローチャートである。図 1 6 に示すように、まず、チャンク管理部 3 は、対象チャンクの配置先を管理するために、チャンク管理 D B 7 中に、その対象チャンクについての複製管理情報を記憶するための領域を確保する（ステップ S 6 0 1）。チャンク管理部 3 は、対象チャンクにチャンクコンテンツ ID が割り当てられていなければ、チャンクコンテンツ ID を割り当て、そのチャンクコンテンツ ID と対応づけた複製管理情報を記憶するための領域を確保し、その複製管理情報を初期化する。チャンク管理部 3 は、例えば、複製数を決定し、決定した複製数を複製管理情報に登録することによって初期化すればよい。チャンクの複製数はシステムの規定の値としてもよいし、複製数計画部 6 によって決定させてもよい。なお、チャンクの配置先情報には、未割り当てを示す情報を登録すればよい。

30

【 0 1 6 3 】

次に、対象チャンクの配置先とする記憶装置を決定する（ステップ S 6 0 2）。なお、配置先の決定方法については後述する。そして、決定した配置先を示す情報を、その対象チャンクのチャンクコンテンツ ID と対応づけた複製管理情報に記録する（ステップ S 6 0 3）。

40

【 0 1 6 4 】

複製数分の配置先を決定していればステップ S 6 0 5 に移行し、まだ決定されていなければステップ S 6 0 2 に移行する（ステップ S 6 0 4）。最後に、対象チャンクのチャンクコンテンツ ID をステップ S 6 0 1 で割り当てていれば、その割り当てたチャンクコンテンツ ID を、その対象チャンクを含むコンテンツの分割管理情報の該当チャンクのチャンク参照情報に記憶し、終了する（ステップ S 6 0 5）。本例では、実際に記憶装置へのチャンクの書き込みは、配置先を通知されたホスト 1 側で行われるが、例えば、後述の複製チャンク数追加処理で示すように、チャンク管理部 3 の指示に応じて複製処理部 5 が行

50

うようにしてもよい。

【0165】

次に、記憶装置を追加した際の処理について説明する。記憶装置を追加した旨が通知されると、チャンク管理部3は、少なくともチャンク管理DB7中の記憶装置情報にその記憶装置の情報を登録する。チャンク管理部3は、記憶装置情報に登録することによって、以降のチャンク作成処理で、追加された記憶装置をチャンクの割り当て先として認識させる。なお、チャンク管理部3は、記憶装置が追加されたことを契機にして、チャンクの再配置を行ってもよい。

【0166】

次に、システムに障害が発生した際の処理について説明する。例えば、障害検出部4は、障害が発生し使用不能になった記憶装置を検出すると、その記憶装置の記憶装置IDをチャンク管理部3に通知する。障害検出部4は、例えば、記憶装置が正常動作をしていることを外部に知らせるために送出している信号を監視し、一定時間以上その送信が確認されないときに使用不能を検出してもよいし、例えば、システム管理者の操作に応じて、障害対応動作が必要な記憶装置IDを入力することによって、使用不能な記憶装置を検出してもよい。使用不能の旨が通知されると、チャンク管理部3は、次に示すような障害対応処理を行う。

【0167】

図17は、チャンク管理部3における障害対応処理の処理フローの一例を示すフローチャートである。図17に示すように、チャンク管理部3は、まず、通知された記憶装置IDで示される記憶装置をチャンクの割り当て対象から外す(ステップS701)。チャンク管理部3は、例えば、チャンク管理DB7中に記憶装置情報として登録されている情報のうち、その記憶装置IDと対応づけられた稼働状況を停止中とすればよい。

【0168】

次いで、本システムが管理しているチャンクコンテンツID(チャンク管理DB7に登録されているチャンクコンテンツID)のうち、未検査のチャンクを選び、そのチャンクの複製管理情報を読み込む(ステップS702)。そして、その複製管理情報において、割り当て対象から外した記憶装置IDが、チャンクの配置先として登録されているか否かを確認し、その記憶装置IDが登録されていた場合ステップS704に、そうでなければステップS705に移行する(ステップS703)。

【0169】

ステップS704では、複製管理情報において、チャンクの配置先として登録されているその装置IDを削除する。例えば、チャンク管理部3は、読み込んだ複製管理情報においてチャンクの配置先情報として示されている、割り当て対象から外した装置IDを、未割り当てを示す情報に変更し、変更した複製管理情報をチャンク管理DB7に記録する。また、ステップS705では、チャンク管理DB7に登録されている全チャンクコンテンツIDについて検査が完了していれば処理を終了し、完了していなければステップS702に移行する。

【0170】

このように、チャンク管理部3により、障害等により使用不能となった記憶装置はチャンク管理DB7中の複製管理情報から削除されるので、その後の読み出しアドレス変換処理においてアクセス先として決定されることはなく、ホスト1が使用不能となった記憶装置にアクセスすることはなくなる。また、チャンク管理DB7中の記憶装置情報にも停止中である旨が登録されることにより、チャンク作成処理において配置先とする記憶装置の対象からも外されるので、使用できない記憶装置がチャンクの配置先として割り当てられることもない。従って、データの保全性および可用性が向上する。

【0171】

また、チャンク管理部3は、障害対応処理の一環として、複製数を維持するための処理を行ってもよい。なお、複製数を維持するための処理は、例えば、障害対応動作を行った後、システムで規定した時間が経過する、または、負荷がシステムで規定した値を下回っ

10

20

30

40

50

たことを契機に自動的に行うようにしてもよい。また、複製数を維持するための処理は、障害対応動作の一環として行うに限らず、記憶装置 2 が追加される等のシステムで規定した契機や、ユーザからの指示を契機に行うことも可能である。

【 0 1 7 2 】

図 1 8 は、複製数を維持するための処理の処理フローを示すフローチャートである。図 1 8 に示すように、チャンク管理部 3 は、まず、チャンク管理 DB 7 に記録されている各チャンクの複製管理情報から、そのチャンクの複製数と、実際に記憶装置に割り当てられている複製チャンクの数とが一致しないものを検索する（ステップ S 1 1 0 1）。ここで、チャンクの複製数は、システムで規定した値やその複製管理情報に登録されている値によって認識すればよい。また、実際に記憶装置に割り当てられている複製チャンクの数は、その複製管理情報に、チャンクの配置先として登録されている記憶装置 ID の数によって認識すればよい。ステップ S 1 1 0 1 において、一致しないものが見つかった場合にはステップ S 1 1 0 3 に移行し、一致しないものが見つからなかった場合には処理を終了する（ステップ S 1 1 0 2）。

10

【 0 1 7 3 】

ステップ S 1 1 0 3 では、検索された複製管理情報に基づき、当該チャンクの複製数と実際に記憶装置に割り当てられている複製チャンクの数とを比較し、その差が、複製数のほうが多い場合にはステップ S 1 1 0 5 に移行し、逆に複製数のほうが少ない場合にはステップ S 1 1 0 4 に移行する。

【 0 1 7 4 】

ステップ S 1 1 0 4 では、当該チャンクの複製数と実際に記憶装置に割り当てられている複製チャンクの数との差をチャンクの追加数として、当該チャンクについて複製チャンク追加処理を行う。なお、複製チャンク追加処理については後述するが、チャンク管理部 3 はコピー元およびコピー先とする記憶装置を決定するにとどめ、実際の記憶装置間のチャンクのコピー処理は複製処理部 5 に行わせることが好ましい。このようにすることによって、チャンク管理部 3 において、障害回復や再配置のための処理による通常のチャンクアクセスにかかる処理への負担が増大しないようにする。

20

【 0 1 7 5 】

また、ステップ S 1 1 0 5 では、当該チャンクの複製数と実際に記憶装置に割り当てられている複製チャンクの数との差をチャンクの削除数として、当該チャンクについて複製チャンク削除処理を行う。なお、複製チャンク削除処理についても同様に、チャンク管理部 3 は削除対象とする記憶装置を決定するにとどめ、実際の記憶装置からのチャンクの削除は複製処理部 5 に行わせることが好ましい。

30

【 0 1 7 6 】

当該チャンクについての複製チャンク追加処理または複製チャンク削除処理が完了すると、次の対象チャンクを検索するため、ステップ S 1 1 0 1 に移行する。

【 0 1 7 7 】

なお、複製チャンク追加処理または複製チャンク削除処理は、複数のチャンクについて並行に実行してもよい。さらに、同時にコピー元およびコピー先として処理を行う記憶装置 2 や削除対象として処理を行う記憶装置 2 が異なれば処理時間が短縮される。処理時間が短縮されれば、複製数が一致していない期間が短縮される。特に、障害により使用不能になった記憶装置が割当先として削除された場合には、複製チャンクの数の減少によって冗長度が減少していることになるので、その期間が短縮されれば、信頼度が向上する。

40

【 0 1 7 8 】

また、本例では、チャンク数追加処理またはチャンク数追加処理だけを複製処理部 5 に行わせる例を示したが、この複製数を維持するための処理全体を複製処理部 5 に行わせ、チャンク管理部 3 は、複製処理部 3 からの依頼に応じてコピー元およびコピー先チャンクや削除対象チャンクを決定するだけでもよい。

【 0 1 7 9 】

また、この複製数を維持するための処理は、障害回復のために行われるだけでなく、複

50

複製数計画部 6 によって複製数に変更された場合にも行われる。例えば、複製数計画部 6 は、各コンテンツのアクセス状況に応じて複製数を変更してもよい。複製数計画部 6 は、例えば、各コンテンツまたはコンテンツ内の各チャンクに求められる複製数の算出を行い、その時点での複製数と算出した複製数とが異なる場合に、チャンク管理 DB 7 中の該当チャンクの複製管理情報における複製数を更新して、そのチャンクについて、チャンク管理部 3 に複製数を維持するための処理を行わせてもよい。または、複製数計画部 6 が、複製数の変更対象をチャンクとして意識しない場合（コンテンツを単位にする場合）や、コンテンツ内における複製数のチャンクを対象とする場合には、算出した複製数とコンテンツ ID と、必要に応じてそのコンテンツ内における対象チャンクの範囲とを指定して、チャンク管理部 3 に後述の複製数変更処理を行わせてもよい。対象チャンクの範囲の指定は、例えば、コンテンツ内におけるオフセットアドレスとサイズで指定すればよい。

10

【 0 1 8 0 】

複製数の算出を行うタイミングとしては、コンテンツへのその時点でのアクセス数がシステムで規定した閾値を超えた場合あるいは閾値を下回った場合、コンテンツのアクセス予想数がシステムで規定した閾値を超えた場合あるいは閾値を下回った場合などがある。また、記憶装置 2 - 1 ~ n 中の未使用記憶容量がシステムで規定した閾値を超えた場合あるいは閾値を下回った場合でもよい。

【 0 1 8 1 】

図 19 は、チャンク管理部 3 における複製数変更処理の処理フローを示すフローチャートである。複製数変更処理は、複製数計画部 6 から複製数の変更が指示された際に行われる。図 19 に示すように、チャンク管理部 3 は、複製数計画部 6 から指定されたコンテンツ ID に対応する分割管理情報をチャンク管理 DB 7 から読み込む（ステップ S 8 0 1）。ここで、対象チャンクの範囲としてオフセットアドレスとサイズとが指定されている場合にはステップ S 8 0 4 に移行し、そうでなければステップ S 8 0 3 に移行する（ステップ S 8 0 2）。

20

【 0 1 8 2 】

ステップ S 8 0 3 では、そのコンテンツ内の全チャンクを複製数の変更対象として設定する。一方、ステップ S 8 0 4 では、指定されたオフセットアドレスとサイズからチャンクの範囲を計算によって把握し、把握したチャンクを複製数の変更対象として設定する。チャンクの範囲の特定方法は、読み出しアドレス変換処理のステップ S 4 0 4 において説明した方法と同様である。

30

【 0 1 8 3 】

変更対象とするチャンクの範囲を特定すると、変更対象中のチャンクからまだ未処理のチャンクを選ぶ（ステップ S 8 0 5）。ここで、未処理のチャンクがない場合には処理を終了し、あった場合にはステップ S 8 0 7 に移行する（ステップ S 8 0 6）。

【 0 1 8 4 】

ステップ S 8 0 7 では、選択したチャンクに対応した複製管理情報をチャンク管理 DB 7 から読み込む。次いで、指定された複製数と複製管理情報中の複製数とを比較し、一致した場合にはステップ S 8 0 5 へ、または、指定された複製数のほうが小さかった場合にはステップ S 8 1 0 へ、逆に指定された複製数のほうが大きかった場合にはステップ S 8 0 9 へ移行する（ステップ S 8 0 8）。

40

【 0 1 8 5 】

ステップ S 8 0 9 では、指定された複製数と複製管理情報中の複製数との差をチャンクの追加数として、選択したチャンクについて複製チャンク追加処理を行う。また、ステップ S 8 1 0 では、指定された複製数と複製管理情報中の複製数との差をチャンクの削除数として、選択したチャンクについて複製チャンク削除処理を行う。選択したチャンクについての複製チャンク追加処理または複製チャンク削除処理が完了すると、次のチャンクを選択するため、ステップ S 8 0 5 に移行する。

【 0 1 8 6 】

なお、ステップ S 8 0 5 以降の処理は並行に実行してもよい。このようにすることで処

50

理時間は短縮される。さらに、同時にコピー元およびコピー先として処理を行う記憶装置 2 や削除の処理を行う記憶装置 2 が異なれば処理時間が短縮される。

【 0 1 8 7 】

次に、複製チャンク削除処理について説明する。図 2 0 は、チャンク管理部 3 における複製チャンク削除処理の処理フローの一例を示すフローチャートである。なお、本例では、複製チャンク削除処理は、削除対象となるチャンクのチャンクコンテンツ ID と現時点での複製管理情報と削除数とが指定されて行われる。図 2 0 に示すように、チャンク管理部 3 は、まず、対象チャンク（コンテンツ内におけるチャンク）の複製管理情報によって示される、対象チャンクの複製チャンクとして現時点で実際に記憶装置が割り当てられているチャンクの中から、削除対象とする複製チャンクを決定する（ステップ S 9 0 1）。各複製チャンクは異なる記憶装置に格納されているため、ここでは、削除対象とする記憶装置を決定することとなる。

10

【 0 1 8 8 】

削除対象の記憶装置が決定すると、チャンク管理部 3 は、対象チャンクの複製管理情報から、削除対象とする記憶装置の記憶装置 ID を削除する（ステップ S 9 0 2）。すなわち、対象チャンクの複製管理情報において、その対象チャンクの配置先情報として示されている記憶装置 ID のうち削除対象とした記憶装置の記憶装置 ID を、未割り当てを示す情報に変更する。そして、変更した複製管理情報をチャンク管理 DB 7 に反映させる（ステップ S 9 0 3）。

【 0 1 8 9 】

次いで、複製処理部 5 に対し、削除対象の記憶装置内における対象チャンク（削除対象として決定した複製チャンク）の削除処理を指示する（ステップ S 9 0 4）。例えば、削除対象とする記憶装置の記憶装置 ID と、削除対象とする複製チャンクのチャンクアクセス ID とを指定して、対象記憶装置内におけるチャンクの削除処理を指示すればよい。そして、複製処理部 5 からの削除完了の通知を待つ（ステップ S 9 0 5）。

20

【 0 1 9 0 】

指定された削除数分、削除処理が完了した場合には処理を終了し、指定された削除数に満たない場合はステップ S 9 0 1 へ移行する（ステップ S 9 0 6）。

【 0 1 9 1 】

なお、本例では、記憶装置内におけるチャンクの削除処理をシーケンシャルに行う例を示しているが、削除対象とする複製チャンクを削除数分決定した上で、複数の記憶装置に対し並行に削除処理を行わせるようにしてもよい。このようにすることで、処理時間が短縮される。なお、記憶装置 2 でのチャンク削除処理は、それぞれ異なる記憶装置 2 で行われるため並列処理の処理効率がよい。

30

【 0 1 9 2 】

また、ステップ S 9 0 5 でチャンクの削除処理の完了待ちを行ったが、完了を待たずに処理を進めてもよい。このようにすることで、さらに処理時間が短縮される。本例では、チャンク管理 DB 7 の更新後に記憶装置 2 に対するチャンクの削除を行っているので、ホスト 1 がチャンク削除処理中のチャンクへアクセスする可能性を低くでき、ホストからのアクセスの安定性が向上する。

40

【 0 1 9 3 】

なお、削除処理を並列に行うと、処理を行う記憶装置 2 の負荷となる場合がある。そのような場合には、記憶装置 2 毎に同時に発行する削除処理数の上限を設け、それを越える場合には削除処理の発行を待つようにしてもよい。このようにすることで、記憶装置 2 の他の処理に使用できる処理能力を保障することができる。すなわちアクセス性能が安定する。

【 0 1 9 4 】

次に、複製処理部 5 での記憶装置内におけるチャンクの削除処理を説明する。複製処理部 5 には、記憶装置 ID とチャンクアクセス ID とが指定され、対象記憶装置内における対象チャンクを削除する旨が指示される。複製処理部 5 は、指定された記憶装置 ID によ

50

り特定される記憶装置 2 - 1 ~ n のいずれかに対し、指定されたチャンクアクセス ID で示されるチャンクの削除要求を発行する。チャンクの削除要求は、具体的には、ブロックデバイスに対するデータ削除命令であったり、ファイルシステムに対するファイルの削除命令であったり、OSD に対するオブジェクトの削除命令である。

【 0 1 9 5 】

次に、複製チャンク追加処理（冗長度増加処理）について説明する。図 2 1 は、チャンク管理部 3 における複製チャンク追加処理の処理フローの一例を示すフローチャートである。なお、複製チャンク追加処理は、追加対象となるチャンクのチャンクコンテンツ ID と現時点での複製管理情報と追加数とが指定されて行われることとする。図 2 1 に示すように、チャンク管理部 3 は、まず、追加対象のチャンク（コンテンツ内におけるチャンク）の複製管理情報に基づき、新たな配置先を決定する（ステップ S 1 0 0 1）。ステップ S 1 0 0 1 では、後述する配置先決定方法に従って、新たな配置先とする記憶装置をコピー先記憶装置として決定する。次いで、ステップ S 1 0 0 1 で割り当てた記憶装置 2 にデータをコピーする際のコピー元記憶装置を決定する（ステップ S 1 0 0 2）。なお、ステップ S 1 0 0 1、ステップ S 1 0 0 2 においては、具体的には、コピー先の記憶装置 ID とその記憶装置内でのチャンクアクセス ID とを決定したり、コピー元の記憶装置 ID とその記憶装置内でのチャンクアクセス ID とを決定する。

10

【 0 1 9 6 】

次いで、複製処理部 5 に対し、記憶装置間での対象チャンク（コピー元およびコピー先として決定した複製チャンク）のコピー処理を指示する（ステップ S 1 0 0 3）。例えば、コピー先とする記憶装置 ID およびチャンクアクセス ID と、コピー元とする記憶装置 ID およびチャンクアクセス ID とを指定して、対象記憶装置間でのチャンクのコピー処理を指示すればよい。なお、コピー先記憶装置でのチャンクアクセス ID が、コピー処理の結果割り当てられる場合には、この時点では指定しなくてもよい。そして、複製処理部 5 からのコピー完了の通知を待って（ステップ S 1 0 0 4）、複製管理情報の対象チャンクの配置先情報に、追加した複製チャンクの配置先を示す情報（ここでは、記憶装置 ID およびチャンクアクセス ID）を追加する（ステップ S 1 0 0 5）。そして、変更した複製管理情報をチャンク管理 DB 7 に反映させる（ステップ S 1 0 0 6）。

20

【 0 1 9 7 】

指定された増加数分、コピー処理が完了した場合には処理を終了し、指定された増加数に満たない場合は S 1 0 0 1 へ移行する（ステップ S 1 0 0 7）。

30

【 0 1 9 8 】

なお、本例では、記憶装置間でのチャンクのコピー処理をシーケンシャルに処理を行っているが、コピー元およびコピー先対象とする複製チャンクを追加数分決定した上で、複数の記憶装置に対し並行にコピー処理を行わせるようにしてもよい。このようにすることで、処理時間が短縮される。なお、記憶装置 2 でのチャンクのコピー処理は、コピー先が異なるだけでなく、コピー元もそれぞれ異なる記憶装置 2 とすることができるので、並列処理の処理効率はよい。

【 0 1 9 9 】

また、本例では、ステップ S 1 0 0 4 でチャンクのコピー処理の完了を待った上で、チャンク管理 DB 7 の更新を行っているので、ホスト 1 がチャンクコピー処理中のチャンクへアクセスする可能性を低くでき、ホストからのアクセスの安定性が向上する。

40

【 0 2 0 0 】

次に、複製処理部 5 での記憶装置間のチャンクのコピー処理を説明する。複製処理部 5 には、コピー元の記憶装置 ID およびその記憶装置内でのチャンクアクセス ID と、コピー先の記憶装置 ID およびその記憶装置内でのチャンクアクセス ID とが指定され、対象記憶装置間で対象チャンクをコピーする旨が指示される。複製処理部 5 は、指定されたコピー元記憶装置 ID により特定される記憶装置に対し、指定されたコピー元チャンクアクセス ID で示されるチャンクの読み出し要求を発行する。そして、対象チャンクのデータが読み出されると、指定されたコピー元記憶装置 ID により特定される記憶装置に対し、

50

指定されたコピー先チャンクアクセスIDで示されるチャンクへの読み出したデータの書き込み要求を発行する。なお、コピー先記憶装置に該当チャンクがない場合、チャンクの作成を行う。チャンクの作成はファイルシステムのファイルの作成、あるいはOSDでのオブジェクトの作成と同様である。

【0201】

また、本例では、データを複製処理部5がいったん読み出してからコピーしたが、ストレージのサードパーティコピー機能と同様に、複製処理部5を介さずに、コピー元記憶装置とコピー先記憶装置間で直接チャンクのコピーを行ってもよい。

【0202】

次に、チャンクの配置先決定方法について説明する。なお、チャンクの配置先決定方法は、図16に示すチャンク作成処理のステップS602や、図21に示す複製チャンク追加処理のステップS1001において、チャンク割り当て処理として実施される。

10

【0203】

まず、同一チャンクの重複配置を防止するための配置先決定方法について説明する。図22は、チャンク管理部3におけるこの方法でのチャンク割り当て処理の処理フローの一例を示すフローチャートである。なお、本例では、チャンク割り当て処理は、割り当て対象となるチャンクのチャンクコンテンツIDと現時点での複製管理情報とが指定されて行われる。なお、チャンク割り当て処理における割当数は1である。図22に示すように、まず、チャンク管理部3は、割り当て可能な全記憶装置をリストアップする(ステップS1201)。割り当て可能な記憶装置は、例えば、チャンク管理DB7中に登録されている記憶装置情報の稼働状況に基づいて、現時点で稼働中の記憶装置の記憶装置IDを集合させた割り当て候補リストを作成すればよい。

20

【0204】

次いで、対象チャンクの複製管理情報で示される、対象チャンクの複製チャンクに既に割り当てられている記憶装置を、割り当て候補リストから削除する(ステップS1202)。ここで、割り当て候補リストには、現在稼働中であって対象チャンクの複製チャンクに対しまだ割り当てられていない記憶装置がリストアップされていることとなる。最後に、割り当て対象リストから、例えば乱数を用いて1の記憶装置を選び、選んだ記憶装置を配置先として決定する(ステップS1203)。

【0205】

このようにすることで、同じ内容のチャンクが異なる記憶装置に配置される。従って、一つの記憶装置に障害が発生しても他の記憶装置に記憶したチャンクにアクセスすることができるので、データの保全性・可用性が向上する。また、コンテンツ内におけるチャンクと同じ内容のチャンク(同じチャンクコンテンツIDのチャンク)が複数あることで、チャンクへのアクセスのスケラビリティが向上する。

30

【0206】

なお、ステップS1203では乱数により配置先を決定したが、例えば、各記憶装置が記憶しているチャンク数や未使用データ領域をテーブルで保持しておき、記憶しているチャンクの数最も少ない記憶装置や、未使用データ領域が最も多い記憶装置に割り当てる方法がある。また、各記憶装置の負荷状態がわかる場合には、最も負荷が低い記憶装置に割り当てる方法もある。

40

【0207】

なお、チャンク管理部3は、ステップS120において割り当て候補リストから全ての記憶装置が削除される場合を考慮するならば、各々の記憶装置に均等化されるように配置先を決定すればよい。具体的には、ステップS1202において、割り当て候補リストから記憶装置を削除する前に、その割り当て候補リストに対象チャンクの割当数を登録するようにし、その後、割当数が最も少ない記憶装置以外の記憶装置を割り当て候補から削除すればよい。

【0208】

次に、コンテンツ内において連続するチャンクの重複配置を防止するための配置先決定

50

方法について説明する。図 2 3 は、チャンク管理部 3 におけるこの方法でのチャンク割り当て処理の処理フローの一例を示すフローチャートである。本例でも、チャンク割り当て処理は、割り当て対象となるチャンクのチャンクコンテンツ ID と現時点での複製管理情報とが指定されて行われる。なお、チャンク割り当て処理における割当数は 1 である。なお、図 2 3 に示すステップ S 1 3 0 1 , S 1 3 0 2 については、図 2 2 におけるステップ S 1 2 0 1 , S 1 2 0 2 と同様であるため、説明省略している。

【 0 2 0 9 】

本例では、さらに、対象チャンクのチャンクコンテンツ ID に対応づけられた分割管理情報から、対象チャンクの前後 N 個以内のチャンクを調べ、それらチャンクの複製管理情報で示される、対象チャンクの前後 N 個以内のチャンクの複製チャンクに既に割り当てられている記憶装置を、割り当て候補リストから削除する（ステップ S 1 3 0 3 ）。ここで、割り当て候補リストには、現在稼働中であって対象チャンクの複製チャンクおよびコンテンツ内においてその対象チャンクの前後 N 個以内のチャンクの複製チャンクに対しまだ割り当てられていない記憶装置がリストアップされていることとなる。最後に、割り当て対象リストから 1 の記憶装置を選び、選んだ記憶装置を配置先として決定する（ステップ S 1 3 0 4 ）。なお、割り当て対象リストからの選択動作については、ステップ S 1 2 0 3 と同様である。

【 0 2 1 0 】

このようにすることで、同じ内容のチャンクが異なる記憶装置に配置されるだけでなく、コンテンツ中の前後 N 個以内のチャンクと同じ内容のチャンクとも異なる記憶装置に配置される。従って、例えば、チャンクアクセスがシーケンシャルに行われる場合にも、各アクセスが異なる記憶装置に分散される。すなわち、シーケンシャルアクセス時のスループットが向上する。

【 0 2 1 1 】

なお、前後 N 個以内の N は、1 つであってもよいが、例えば、バッファを使用して先読みを行うようなコンテンツを扱う場合には、先読みで同時にアクセスする量に相当するチャンクが含まれるように規定することが好ましい。そうすることで、先読みアクセスにおいてもスループットの向上とレスポンスタイムの短縮が期待される。

【 0 2 1 2 】

次に、記憶装置間の相関を平均化するための配置先決定方法について説明する。、例えば、記憶装置内のチャンクの配置内容が、他の記憶装置内のチャンクの配置内容とを一致させないように配置決定する方法である。チャンク管理部 3 は、2 つの記憶装置間で共有する同じ内容の複製チャンクの数に共有数とした場合に、各々の記憶装置の組み合わせにおける共有数が均等化されるように、複製チャンクの配置先を決定すればよい。このようにすることで、記憶装置内のチャンクの配置パターンが均等化され、ある記憶装置が障害等により使用できなくなった際に他の稼働できる記憶装置への代替による負荷が均等化する。

【 0 2 1 3 】

実現方法の 1 つを説明する。チャンク管理 DB 7 に記憶装置間の共通度を格納するテーブルを用意する。

【 0 2 1 4 】

このテーブルの操作について説明する。このテーブルは、記憶装置数 $n \times n$ の二次元配列テーブルであって、2 記憶装置間で共有するチャンク数（以下、単に共有数という。）をその 2 記憶装置の記憶装置 ID を添字とする配列要素に示すテーブルである。以下、このテーブルを配列 A と呼ぶ。配列 A の各配列要素の初期値は全て 0 とする。なお、ここでは、記憶装置 ID には 0 からの整数が割り当てられているものとする。チャンク管理部 3 は、同一内容の複数のチャンクを一度に配置決定した場合、同一内容のチャンクの配置先となった記憶装置 ID の組み合わせを添字とする配列要素で示される共有数を + 1 する。例えば、記憶装置 2 - a , c , k に同一内容のチャンクを割り当てた場合には、配列 A [a] [c] , A [c] [a] , A [a] [k] , A [k] [a] , A [c] [k] , A [

10

20

30

40

50

k][c]に、1を加算する。

【0215】

また、チャンク管理部3は、あるチャンクを記憶装置から削除した場合には、そのチャンクの配置先であった記憶装置IDと、現時点でそのチャンクと同一内容のチャンクが割り当てられている他の記憶装置IDとの組み合わせを添字とする配列要素で示される共有数を-1する。例えば、記憶装置2-a, c, kに割り当てられているチャンクのうち、記憶装置2-kに割り当てられているチャンクを削除した場合には、配列A[a][k], A[k][a], A[c][k], A[k][c]から1を減算する。

【0216】

また、チャンク管理部3は、あるチャンクを記憶装置に追加した場合には、そのチャンクの配置先となった記憶装置IDと、現時点でそのチャンクと同一内容のチャンクが割り当てられている他の記憶装置IDとの組み合わせを添字とする配列要素で示される共有数を+1する。例えば、記憶装置2-a, c, kに配置されているチャンクと同じ内容のチャンクを新たに記憶装置2-hに割り当てた場合、配列A[a][h], A[c][h], A[k][h], A[h][a], A[h][c], A[h][k]に1を加算する。

【0217】

チャンク管理部3が新たなチャンクの配置先を決定する場合には、このように扱われるテーブルを用いて、そのチャンクと同一内容のチャンクが既に割り当てられている記憶装置と最も相関の低い記憶装置を配置先として決定する。具体的には、既に記憶装置に割り当てられているチャンクと同一のチャンクの配置を決定する場合には、そのチャンクが割り当てられている記憶装置IDを添字の1つとして固定して配列Aの各配列要素を見た場合に、共有数が最も小さくなる添字を示す記憶装置IDを配置先として決定する。なお、既に記憶装置に割り当てられているチャンクが複数ある場合には、各記憶装置IDを添字の1つにそれぞれ固定した複数の配列Aの和として見ればよい。例えば、既に記憶装置2-a, cに配置されているチャンクと同一内容のチャンクの配置先を決定する場合には、配列A[a][i]+配列A[c][i]が最小となるiを示す記憶装置2-iを配置先として決定すればよい。

【0218】

なお、チャンクの割り当て処理としては、例えば、図22に示すステップS1203や図23に示すステップS1304において、割り当て候補リストにある記憶装置の中から、上記iを満たす記憶装置を選べばよい。なお、複数の記憶装置でiを満たす場合には、乱数等により配置先を決定すればよい。

【0219】

次に、チャンク管理部3における読み出し元記憶装置の決定方法について説明する。なお、読み出し元記憶装置決定方法は、図21に示す冗長度増加処理のステップS1002におけるコピー元複製チャンクの決定動作として実施される。チャンク管理部3は、各記憶装置で行っているチャンクのコピー処理の数をテーブルにより保持しているものとする。図21に示す例では、ステップS1003でコピー元記憶装置のコピー処理数を+1し、ステップS1004が完了するとその記憶装置のコピー処理数を-1すればよい。

【0220】

そして、チャンク管理部3は、読み出し元を決定する際に、上記テーブルを参照して、追加対象チャンクが割り当てられている記憶装置の中で最もコピー数が少ない記憶装置をコピー元記憶装置に決定すればよい。なお、選択候補の記憶装置のコピー数がシステムで規定した数よりも大きかった場合に、コピー数が規定数未満になるまでコピー処理を待つ方法もある。このようにすることで、コピーに使用する帯域を制限でき、コピー処理中の記憶装置で他のアクセスに使用する帯域を確保することができる。よって、アクセスの安定性が保証される。

【0221】

また、上記例では、コピー元記憶装置についてコピー数を管理しているが、コピー先の記憶装置についてもコピー数を管理し、コピー先記憶装置のコピー数がシステムで規定し

10

20

30

40

50

た数よりも大きかった場合に、コピー処理を待つ方法もある。このようにすることで、コピー先の記憶装置で他の処理に使用できる処理能力を保障することができる。よって、アクセス性能が安定する。

【0222】

また、障害等により使用できなくなった記憶装置の代替として新規に記憶装置を追加した場合など、割り当て先を代替用に追加した記憶装置に固定する方法もある。更に代替記憶装置にコピーを行う期間中は、読み出し対象としないなど代替記憶装置をコピー以外のアクセスの対象外とするとともに、他の記憶装置からのコピー処理を並行に行う方法がある。このようにすると、代替記憶装置へのコピー処理については、代替記憶装置のバンド幅全てが使用できるため、コピー処理の時間を短縮できる。すなわち、冗長度が減少している時間が短縮され、可用性が向上する。また、代替記憶装置が短い時間で復帰すれば、使用不能となった記憶装置に割り当てられたチャンクアクセスの代替による負荷がかかる記憶装置へのアクセス性能も安定する。

10

【0223】

次に、チャンク管理部3における削除対象記憶装置の決定方法について説明する。なお、削除対象記憶装置の決定方法は、図20に示す複製チャンク数削除処理のステップS901における削除対象の複製チャンクの決定動作として実施される。

【0224】

チャンク管理部3は、例えば、対象チャンクの複製管理情報を参照し、対象チャンクの複製チャンクが割り当てられている記憶装置の中から乱数により1の記憶装置を選択してもいい。また、例えば、各記憶装置が記憶しているチャンク数や未使用データ領域をテーブルで保持しておき、記憶しているチャンクの数最も多い記憶装置や、未使用データ領域が最も少ない記憶装置を削除対象とする方法もある。また、各記憶装置の負荷状態がわかる場合には、最も負荷が高い記憶装置を削除対象とする方法もある。

20

【0225】

次に、複製数計画部6におけるチャンクの複製数の決定方法について説明する。複製数は、可用性・保全性から決定される可用性面での複製数と、スループット・レスポンスタイム等から決定されるアクセス面での複製数とがあるが、ここでは、両者の大きい方の値を複製数として決定する。

【0226】

可用性面での複製数は、記憶装置が障害等によりアクセスできなくなった際に、障害で使えなくなった記憶装置に格納されていたチャンクと同じチャンクが他のいずれかの記憶装置でアクセス可能であることを保障するための数である。

30

【0227】

さらに、障害発生後にコピーを作成して複製数を維持する処理を行う場合、そのコピーが完了するまでの間に他の記憶装置にも障害が発生する可能性を考慮して決定される。コピー完了までの時間が長くなればなるほど、可用性面の複製数は多くなる。すなわち、記憶装置の容量とコピーに使用できる帯域と記憶装置の信頼度から算出される値である。可用性面での複製数は記憶装置が1つ使用できなくなっても冗長度が0とならないよう、3以上と定める。

40

【0228】

アクセス面での複製数は、コンテンツの要求予測量と、各記憶装置の処理能力と、システムで規定した安全係数とに基づいて決定される。例えば、そのコンテンツの必要スループットを各記憶装置のスループットで割ったものにシステムで規定した安全係数をかけたものである。ここで、安全係数とは、予想以上に要求が増えた場合に備えた安全係数の意と、複数のコンテンツで記憶装置を共有することによるそのコンテンツでの記憶装置の占有率の意を含めた係数である。安全係数は、1以下の値を設定する。コンテンツの必要スループットは、コンテンツのエンコードのビットレートと要求ユーザ数との積によって求めてもよい。なお、記憶装置の平均スループットや最低スループットを用いてもよい。または、配置決定した記憶装置のスループットの和をとって、必要スループットに足りな

50

ったら複製数を追加するという方法もある。

【0229】

また、コンテンツの先頭やコンテンツにチャプタが付いている場合のチャプタの先頭に位置するチャンクの複製数を多くする方法もある。ビデオのようなメディアコンテンツのアクセスではバッファを用いてアクセスの安定性を確保するが、コンテンツの先頭やチャプターの直後はバッファによる安定化ができないため、例えば、コンテンツの先頭やチャプターの直後に該当するデータを含むチャンクについては、他のチャンクよりも複製数を多くすることで再生の安定化が図れる。

【0230】

次に、複製チャンクの再配置について説明する。複製チャンクの再配置を行う契機としては、例えば、ユーザが指示した場合や、記憶装置が追加された場合、記憶装置が削除された場合、ある記憶装置の記憶容量の使用率またはbusy率が閾値(80%, 90%等)を超えた場合、ある記憶装置の使用率(またはbusy率)とシステム全体の平均使用率(またはbusy率)との差が閾値を超えた場合などが考えられる。

10

【0231】

例えば、記憶装置が追加された場合、まず、チャンク管理部3が、チャンク管理DB7に記憶されている各チャンクの複製管理情報において、追加された記憶装置以外の記憶装置が配置先として割り当てられているチャンク(複製チャンク)をランダムに抽出する。そして、抽出されたチャンクを追加された記憶装置に格納できるか否かを判断する。ここで、格納できないチャンクの例としては、追加された記憶装置を既に格納先としたチャンクと同一内容のチャンクや、追加された記憶装置を既に格納先としてチャンクの前後N個以内のチャンクや、追加された記憶装置に記憶することである記憶装置間の相関が高くなるようなチャンクである。

20

【0232】

抽出したチャンクが追加された記憶装置に格納できる場合には、そのチャンクを対象チャンクとして新たな記憶装置に移動させる。格納できない場合には、他のチャンクを再度抽出する。チャンクの移動方法としては、例えば、複製処理部5に対象チャンクのコピー処理を行わせたのち、その対象チャンクの配置先情報を、元の記憶装置から新たな記憶装置に変更してチャンク管理DB7に反映させる。そして、複製処理部5に元の記憶装置からの対象チャンクの削除処理を行わせればよい。以上の動作を、例えば、システムで規定されているチャンク数M分だけ繰り返し行うようにすればよい。なお、チャンク数Mは、システムに予め設定された値を用いてもよいし、再配置を行う際にユーザが指定した値を用いてもよい。また、追加された記憶装置の容量に応じて決定する方法もある。

30

【0233】

このようにすることで、追加された記憶装置に即座にチャンクが割り当てられ、追加された記憶装置が使用されるようになる。また、追加された記憶装置に対してもチャンクの配置先決定ルールが適用されるので、可用性面での複製数が低くなることも、アクセス面での複製数が低くなることも、記憶装置間の相関が高くなることもない。なお、上記で示したチャンクの移動方法によれば、移動中のチャンクにアクセスされる可能性も低い。

【0234】

また、記憶装置が削除された場合の再配置動作は、使用不能になった記憶装置に対する障害対処理における複製数を維持するための処理と同様である。

40

【0235】

また、それ以外の場合の再配置動作としては、例えば、チャンク管理部3が、任意の記憶装置または使用率やbusy率において高負荷が検出された記憶装置から、移動対象とするチャンクをランダムに抽出する。次に、移動対象チャンクを抽出した記憶装置以外の記憶装置から、移動先とする記憶装置をランダムに選択する。そして、移動対象チャンクを移動先として選択した記憶装置に格納できるか否かを判断する。

【0236】

抽出した移動対象チャンクが移動先として選択した記憶装置に格納できる場合には、そ

50

のチャンクを選択した記憶装置に移動させる。格納できない場合には、移動先とする記憶装置を再度選択する。再度選択した記憶装置にも格納できない場合や全記憶装置に格納できない場合には、移動対象チャンクを再度抽出する。なお、チャンクの移動方法は、既に説明した方法と同様である。

【0237】

このようにすることで、システム全体としての記憶装置間の相関を更に低くすることができる。さらに、再配置に伴って可用性面での複製数が低くなることも、アクセス面での複製数が低くなることもない。

【0238】

以上のように、本実施の形態によれば、コンテンツを複数のチャンクに分割し、さらにそのチャンクを複製してできる複製チャンクを、所定の配置先決定ルールに従って複数の記憶装置に分散させて格納することによって、データ分散格納における特性・性能の向上を図ることができる。

10

【0239】

例えば、チャンク管理部3が、同一内容のチャンクがそれぞれ異なる記憶装置に格納されるように各チャンクの配置を決定することによって、記憶装置が障害により使用できない場合であってもその記憶装置に格納されているチャンクと同一内容のチャンクを格納している他の記憶装置でそのチャンクアクセスが可能とあり、可用性を向上させることができる。また、チャンク管理部3が複数チャンクのアクセス先を決定すれば、ホストからの読み出しを確実に分散させることができ、アクセス性を向上させることができる。

20

【0240】

また、例えば、チャンク管理部3が、コンテンツ内において連続するチャンクと同一内容のチャンクがそれぞれ異なる記憶装置に格納されるように各チャンクの配置を決定することによって、さらにシーケンシャルアクセスでのアクセス性を向上させることができる。

【0241】

また、例えば、チャンク管理部3が、上記条件に加えて、記憶装置間の相関を低くするように配置決定することによって、障害等により使用できなくなった記憶装置に格納されていたチャンクへのアクセスを、他の記憶装置全体に平均的に分散することができる。従って、障害発生時のアクセス性低下の度合いを低減させることができる。また、障害発生後に複製数を維持させる場合であっても、そのためのコピー処理におけるコピー元が他の記憶装置全体に分散されるので、並列処理を可能とし、回復に要する時間を短縮することができる。すなわち、障害発生後の可用性低下の度合いを低減させることができる。なお、連続するチャンクと同一内容の複製チャンクとの記憶装置の共有を防ぎつつ、記憶装置間の相関を低くするように配置決定することも可能である。

30

【0242】

また、チャンクの記憶装置群への配置管理をチャンク管理部3で行うことにより、コンテンツが形成する論理データ空間と記憶装置群が形成する物理データ空間とのマッピングが柔軟になり、記憶装置の追加に対して柔軟に対応できる。

【0243】

また、チャンク管理部3が障害対応動作として使用できなくなった記憶装置をアクセス対象外からはずすことによって、ホストが障害によりアクセスできなくなった記憶装置にアクセスすることがなくなる。すなわち障害発生時のアクセス失敗の可能性を低減することができる。

40

【0244】

また、複製数計画部6が、コンテンツへのアクセス要求の変化に即して、コンテンツの複製数を変更することによって、変化するコンテンツへのデマンドに対処できかつ、記憶容量の無駄を抑えることが可能となる。すなわち、変化するコンテンツのデマンドに適したアクセス性を提供しながら容量効率を向上できる。なお、複製数計画部6がコンテンツ内におけるチャンクの位置に応じて複製数を決定することによって、コンテンツの位置によ

50

り異なるデマンドに対応できる。すなわち、変化するコンテンツのデマンドに適したアクセス性を提供しながら容量効率を向上できる。

【産業上の利用可能性】

【0245】

本発明の活用例として、ストリーミング配信サーバがある。特に、不特定多数のユーザに向けて多量のコンテンツをサービス対象とするストリーミング配信サーバのように、多くの記憶装置を制御対象とするデータ分散格納システムに好適に適用可能である。

【図面の簡単な説明】

【0246】

【図1】本発明によるデータ分散格納システムの構成例を示すブロック図である。 10

【図2】本実施の形態によるデータ分散格納システムの構成例を示すブロック図である。

【図3】分割管理情報の例を示す説明図である。

【図4】複製管理情報の一例を示す説明図である。

【図5】チャンクの連結結果の一例を示す説明図である。

【図6】チャンクの連結結果の一例を示す説明図である。

【図7】コンテンツ読み出し動作の一例を示すフローチャートである。

【図8】読み出し用アドレス変換要求に対する応答フォーマットの一例を示す説明図である。

【図9】コンテンツ書き込み動作の一例を示すフローチャートである。

【図10】書き込み用アドレス変換要求に対する応答フォーマットの一例を示す説明図である。 20

【図11】チャンク管理部3におけるアドレス変換処理の処理フローの一例を示すフローチャートである。

【図12】読み出しアドレス変換処理の処理フローの一例を示すフローチャートである。

【図13】読み出し用アドレス変換要求に対する応答フォーマットの他の例を示す説明図である。

【図14】読み出し用アドレス変換要求に対する応答フォーマットの他の例を示す説明図である。

【図15】書き込みアドレス変換処理の処理フローの一例を示すフローチャートである。

【図16】チャンク管理部3におけるチャンク作成処理の処理フローの一例を示すフローチャートである。 30

【図17】チャンク管理部3における障害対応処理の処理フローの一例を示すフローチャートである。

【図18】複製数を維持するための処理の処理フローを示すフローチャートである。

【図19】チャンク管理部3における複製数変更処理の処理フローを示すフローチャートである。

【図20】チャンク管理部3における複製チャンク削除処理の処理フローの一例を示すフローチャートである。

【図21】チャンク管理部3における複製チャンク追加処理の処理フローの一例を示すフローチャートである。 40

【図22】チャンク管理部3におけるチャンク割り当て処理の処理フローの一例を示すフローチャートである。

【図23】チャンク管理部3におけるチャンク割り当て処理の処理フローの一例を示すフローチャートである。

【符号の説明】

【0247】

1 ホスト

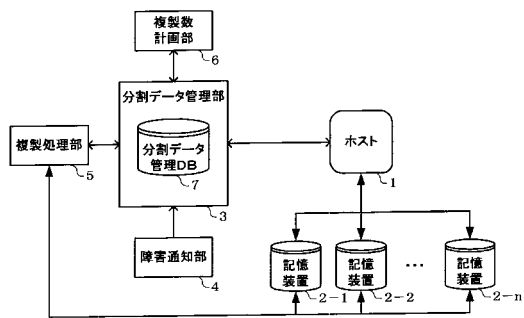
2 - 1 ~ n 記憶装置

3 分割データ管理部、チャンク管理部

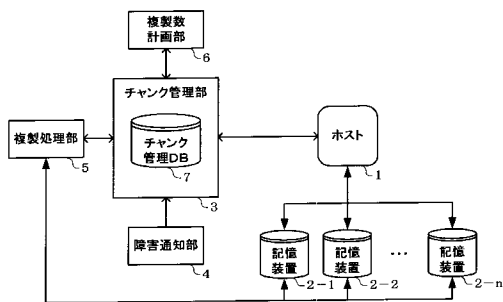
4 障害通知部 50

- 5 複製処理部
- 6 複製数計画部
- 7 分割データ管理DB、チャンク管理DB

【図1】



【図2】



【図3】

チャンク参照情報1	チャンク参照情報2	チャンク参照情報3	...	チャンク参照情報n
-----------	-----------	-----------	-----	-----------

【図4】

記憶装置ID1	記憶装置ID2	記憶装置ID3	...	記憶装置IDn
---------	---------	---------	-----	---------

【図5】

コンテンツファイル内におけるチャンク

0	1	2	3	4	5	6	7	8	9	10	11
---	---	---	---	---	---	---	---	---	---	----	----

記憶装置ID=0のチャンク

0	4	8
---	---	---

記憶装置ID=10のチャンク

0	7	10
---	---	----

記憶装置ID=1のチャンク

1	5	9
---	---	---

記憶装置ID=11のチャンク

1	4	11
---	---	----

記憶装置ID=2のチャンク

2	6	10
---	---	----

記憶装置ID=12のチャンク

2	5	8
---	---	---

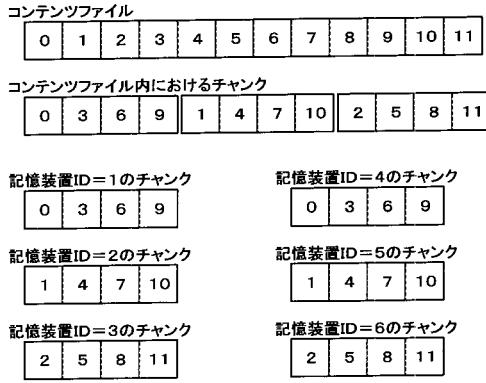
記憶装置ID=3のチャンク

3	7	11
---	---	----

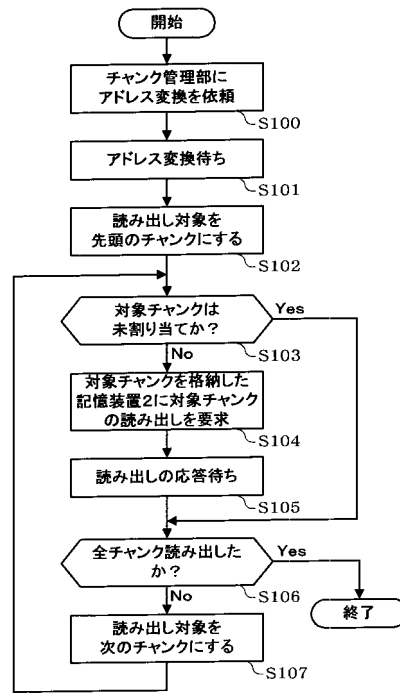
記憶装置ID=13のチャンク

3	6	9
---	---	---

【図6】



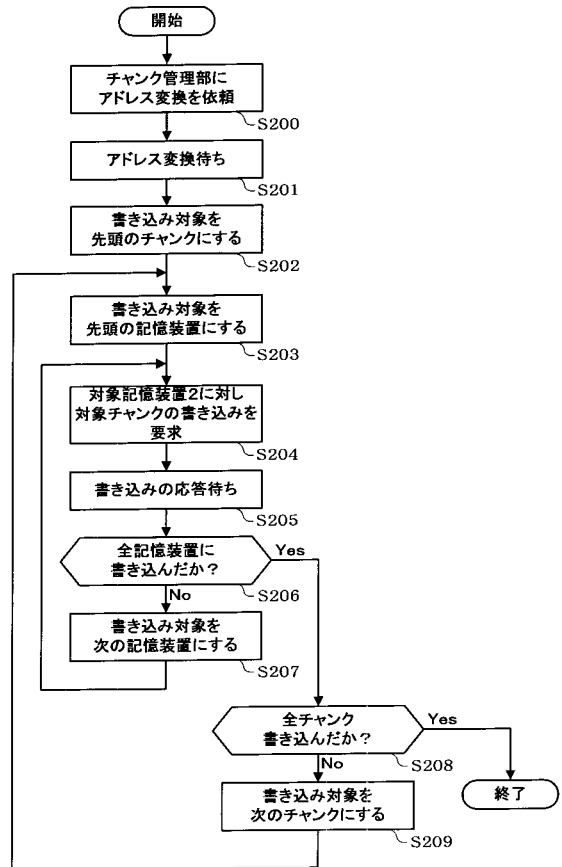
【図7】



【図8】

エン트리1:	記憶装置ID	チャンクアクセスID	チャンク内オフセットアドレス	読み出しサイズ
エン트리2:	記憶装置ID	チャンクアクセスID	チャンク内オフセットアドレス	読み出しサイズ
	:	:	:	:
エン트리n:	記憶装置ID	チャンクアクセスID	チャンク内オフセットアドレス	読み出しサイズ

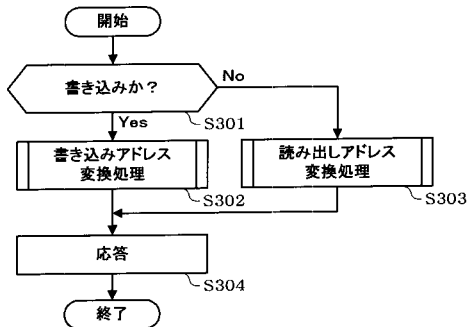
【図9】



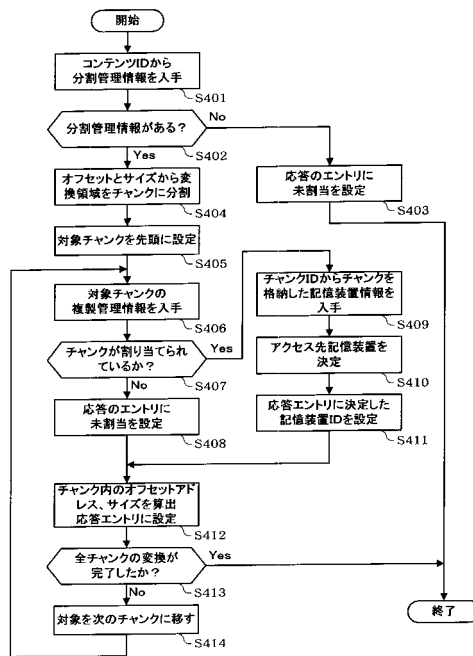
【図10】

エン트리1:	書き込み サイズ	チャンク内 オフセットアドレス	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID
エン트리2:	書き込み サイズ	チャンク内 オフセットアドレス	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID
:								
エン트리n:	書き込み サイズ	チャンク内 オフセットアドレス	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID

【図11】



【図12】



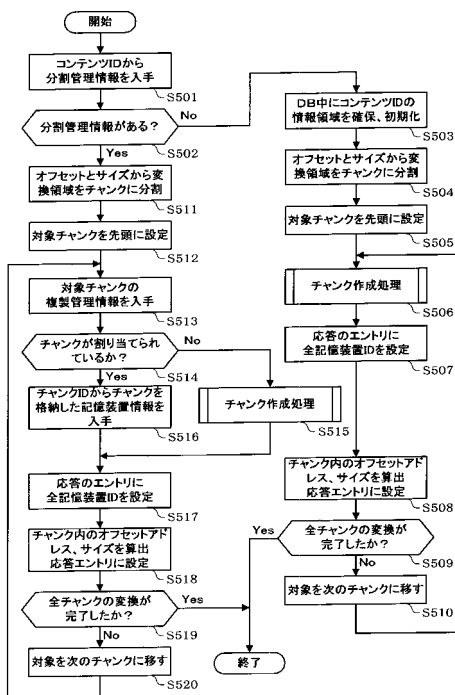
【図13】

エン트리1:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID
エン트리2:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID
:								
エン트리n:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	記憶 装置ID	チャンク アクセスID	...	記憶 装置ID	チャンク アクセスID

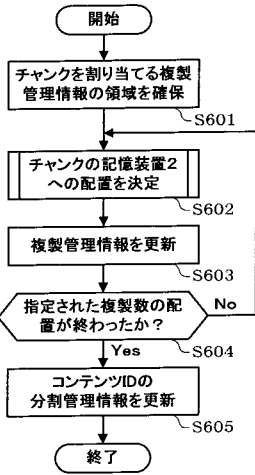
【図14】

エン트리1:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	チャンク アクセスID	記憶 装置ID	...	記憶 装置ID
エン트리2:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	チャンク アクセスID	記憶 装置ID	...	記憶 装置ID
:							
エン트리n:	チャンク内 オフセットアドレス	読み出し サイズ	複製 数	チャンク アクセスID	記憶 装置ID	...	記憶 装置ID

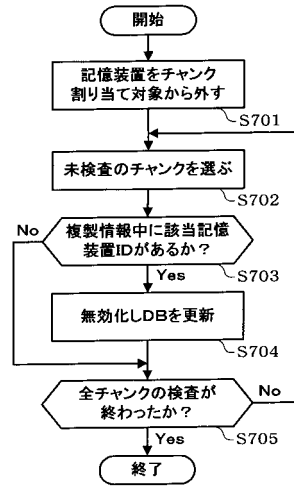
【図15】



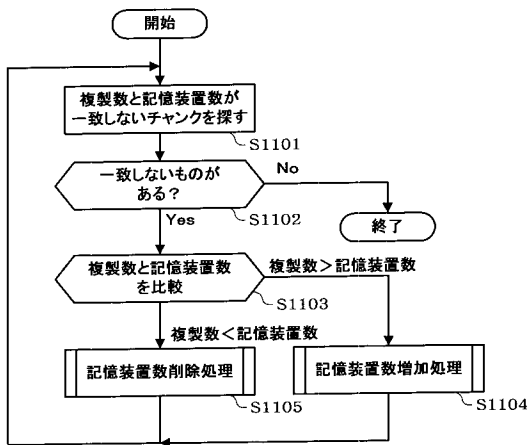
【図16】



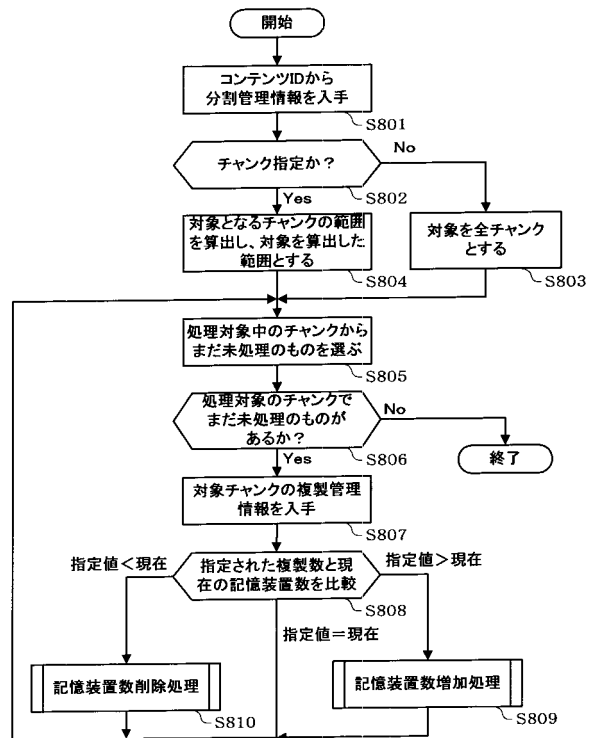
【図17】



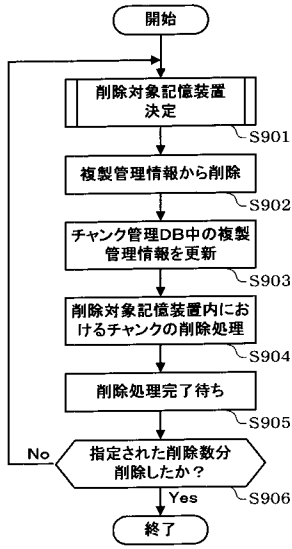
【図18】



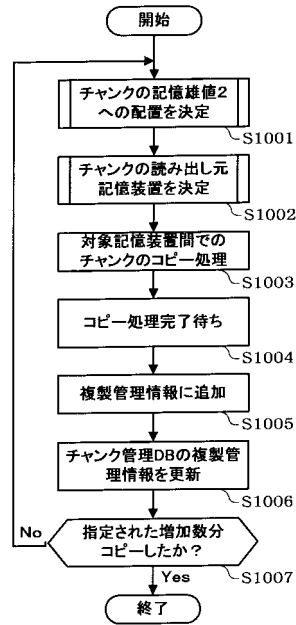
【図19】



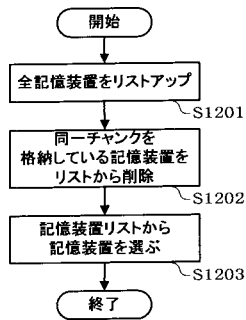
【図20】



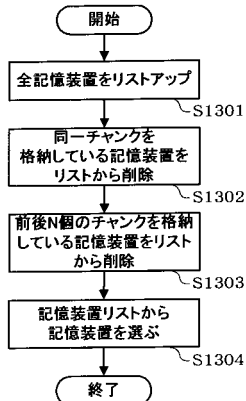
【図21】



【図22】



【図23】



フロントページの続き

- (56)参考文献 特開平07-141232(JP,A)
特開2006-350470(JP,A)
特開平09-091186(JP,A)
特開2004-046352(JP,A)
特開2005-071247(JP,A)
特開2001-243099(JP,A)
藤田 憲治, ディスク・アレイ導入の勧め, 日経バイト, 日本, 日経BP社, 1999年 4月
22日, 第190号, pp.104~117
IBM ServeRAID ユーザーズ・リファレンス 第4版, 日本, 日本アイ・ビー・エム
株式会社, 2007年 1月31日, 第4版, pp.24~26

(58)調査した分野(Int.Cl., DB名)

G06F 12/00
G06F 3/06
JSTPlus(JDreamIII)