

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6005566号
(P6005566)

(45) 発行日 平成28年10月12日 (2016. 10. 12)

(24) 登録日 平成28年9月16日 (2016. 9. 16)

(51) Int. Cl.	F I
G 0 6 F 12/16 (2006. 01)	G 0 6 F 12/16 3 1 0 Q
G 0 6 F 13/10 (2006. 01)	G 0 6 F 13/10 3 4 0 A
G 0 6 F 3/06 (2006. 01)	G 0 6 F 3/06 3 0 1 K

請求項の数 56 (全 105 頁)

(21) 出願番号	特願2013-55548 (P2013-55548)	(73) 特許権者	000003078
(22) 出願日	平成25年3月18日 (2013. 3. 18)		株式会社東芝
(65) 公開番号	特開2014-182503 (P2014-182503A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成26年9月29日 (2014. 9. 29)	(74) 代理人	110002147
審査請求日	平成27年2月10日 (2015. 2. 10)		特許業務法人酒井国際特許事務所
		(72) 発明者	橋本 大輔
			東京都港区芝浦一丁目1番1号 株式会社東芝内
		(72) 発明者	永井 宏一
			東京都港区芝浦一丁目1番1号 株式会社東芝内
		(72) 発明者	渡部 孝紀
			東京都港区芝浦一丁目1番1号 株式会社東芝内

最終頁に続く

(54) 【発明の名称】 情報処理システム、制御プログラムおよび情報処理装置

(57) 【特許請求の範囲】

【請求項 1】

リードおよびライト可能な第1のメモリを有する第1の記憶部と、引き継ぎ履歴格納部が記憶され、リードおよびライト可能な第2のメモリを有する第2の記憶部とを接続可能な情報処理装置であって、

前記第1の記憶部から取得した信頼性情報に基づき、前記第1の記憶部のディスクステータスを判定する第1の処理と、

前記第1の処理により前記第1の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第2の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第2の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第1の記憶部からデータを読み出す第2の処理と、

前記第2の処理による書き込みに伴って、前記第2の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第3の処理と、

前記第1の処理により前記第1の記憶部のディスクステータスが保護状態として認識されなかった場合、読み出しは前記引き継ぎ履歴格納部を読み出すこと無く前記第1の記憶部からデータを読み出す第4の処理と、

を実行する制御部

を備えることを特徴とする情報処理装置。

【請求項 2】

10

20

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 1 に記載の情報処理装置。

【請求項 3】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 1 に記載の情報処理装置。

【請求項 4】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうちの一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 1 に記載の情報処理装置。

【請求項 5】

前記制御部は、

前記第 1 の記憶部に書き込まれているデータのうち前記第 2 の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履歴格納部に格納されている場合に、前記第 1 の記憶部の保護状態を終了すると判定することを特徴とする請求項 1 乃至 4 の何れか一つに記載の情報処理装置。

【請求項 6】

前記制御部は、

前記第 3 の処理では、

前記引き継ぎ履歴格納部に第 1 のアドレスを記録する場合に、前記第 1 のアドレスと重複または連続する第 2 のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第 1 のアドレスと前記第 2 のアドレスを合併した第 3 のアドレスを前記引き継ぎ履歴格納部に記録し、前記第 2 のアドレスを無効化することを特徴とする請求項 1 乃至 5 の何れか一つに記載の情報処理装置。

【請求項 7】

前記制御部は、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、前記引き継ぎ履歴格納部を消去することを特徴とする請求項 1 乃至 6 の何れか一つに記載の情報処理装置。

【請求項 8】

リードおよびライト可能な不揮発性の半導体メモリである第 1 のメモリを有する第 1 の記憶部と、引き継ぎ履歴格納部が記憶され、リードおよびライト可能な不揮発性の半導体メモリである第 2 のメモリを有する第 2 の記憶部と、引き継ぎ履歴格納部とを接続可能な情報処理装置であって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出す第 2 の処理と、

前記第 2 の処理による書き込みに伴って、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第 3 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、読み出しは前記引き継ぎ履歴格納部を読み出すこと無く前記第 1 の記憶部からデータを読み出す第 4 の処理と、

を実行する制御部

を備えることを特徴とする情報処理装置。

【請求項 9】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 8 に記載の情報処理装置。

【請求項 10】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 8 に記載の情報処理装置。

【請求項 11】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうちの一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 8 に記載の情報処理装置。

【請求項 12】

前記制御部は、

前記第 1 の記憶部に書き込まれているデータのうち前記第 2 の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履歴格納部に格納されている場合に、前記第 1 の記憶部の保護状態を終了すると判定することを特徴とする請求項 8 乃至 11 の何れか一つに記載の情報処理装置。

【請求項 13】

前記制御部は、

前記第 3 の処理では、

前記引き継ぎ履歴格納部に第 1 のアドレスを記録する場合に、前記第 1 のアドレスと重複または連続する第 2 のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第 1 のアドレスと前記第 2 のアドレスを合併した第 3 のアドレスを前記引き継

10

20

30

40

50

ぎ履歴格納部に記録し、前記第 2 のアドレスを無効化することを特徴とする請求項 8 乃至 12 の何れか一つに記載の情報処理装置。

【請求項 14】

前記制御部は、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、前記引き継ぎ履歴格納部を消去することを特徴とする請求項 8 乃至 13 の何れか一つに記載の情報処理装置。

【請求項 15】

リードおよびライト可能な第 1 のメモリを有する第 1 の記憶部と、リードおよびライト可能な第 2 のメモリを有する第 2 の記憶部とを接続可能な情報処理装置であって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態でないと判定された場合、前記第 1 の記憶部に書き込みを実行し、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 2 の記憶部に書き込みを実行する第 2 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された後に前記第 1 の記憶部に対する読み出しが発生した場合は、前記第 1 の記憶部からデータを読み出し、前記読み出されたデータを前記第 2 の記憶部に書き込む第 3 の処理と、

を実行する制御部

を備えることを特徴とする情報処理装置。

【請求項 16】

前記制御部は、

前記第 3 の処理では、前記第 1 の記憶部からデータを読み出す処理と、前記読み出されたデータを前記第 2 の記憶部に書き込む処理を並行して実行することを特徴とする請求項 15 に記載の情報処理装置。

【請求項 17】

メタデータを格納するメタデータ格納部に接続され、

前記制御部は、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態でないと判定された場合、前記第 1 の記憶部に書き込まれたデータのアドレスと前記第 1 の記憶部の識別情報との紐付けを前記メタデータに記録し、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態と判定された場合、前記第 2 の記憶部に書き込まれたデータのアドレスと前記第 2 の記憶部の識別情報との紐付けを前記メタデータに記録し、

読み出し処理において、読み出し対象アドレスと前記第 1 の記憶部の識別情報との紐付けが前記メタデータに記録されている場合は、前記第 1 の記憶部からデータを読み出し、読み出し対象アドレスと前記第 2 の記憶部との紐付けが前記メタデータに記録されている場合は、前記第 2 の記憶部からデータを読み出すことを特徴とする請求項 15 乃至 16 の何れか一つに記載の情報処理装置。

【請求項 18】

引き継ぎ履歴格納部に接続され、

前記制御部は、

前記第 3 の処理では、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された場合、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録し、

読み出し処理において、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出すことを特徴とする請求項 15 乃至 16 の何れか一つに記載の情報処理装置。

【請求項 19】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 15に記載の情報処理装置。

【請求項 20】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

10

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 15に記載の情報処理装置。

【請求項 21】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうちの一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

20

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 15に記載の情報処理装置。

【請求項 22】

リードおよびライト可能な第 1 のメモリを有する第 1 の記憶部と、リードおよびライト可能な第 2 のメモリを有する第 2 の記憶部と、引き継ぎ履歴格納部とを接続可能な情報処理装置であって、

30

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出す第 2 の処理と、

前記第 2 の処理による書き込みに伴って、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第 3 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された後に、前記第 1 の記憶部に格納されたデータの削除処理命令を受信した場合に、前記削除処理命令による削除対象データのアドレスを前記引き継ぎ履歴格納部に記録する第 4 の処理

40

を実行する制御部

を備えることを特徴とする情報処理装置。

【請求項 23】

前記制御部は、

前記第 4 の処理では、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された後に、前記第 1 または第 2 の記憶部に格納されたデータの削除処理命令を受信した場合に

50

、前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスが記録されている場合は、前記受信された削除処理命令を前記第２の記憶部に送信し、前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスを記録せず、

前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスが記録されていない場合は、前記削除処理命令による削除対象データのアドレスを前記引き継ぎ履歴格納部に記録する

ことを特徴とする請求項２２に記載の情報処理装置。

【請求項２４】

前記制御部は、

前記第１の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第１の記憶部の寿命到達条件が成立したと判定したときに、前記第１の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項２２に記載の情報処理装置。

10

【請求項２５】

前記制御部は、

前記第１の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第１の記憶部の寿命到達条件が成立したと判定したときに、前記第１の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第１の処理により前記第１の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第２の記憶部を前記第１の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第１の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項２２に記載の情報処理装置。

20

【請求項２６】

リードおよびライト可能な第３のメモリを有する１乃至複数の第３の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第１および第２の記憶部のうち的一方に存在するように、前記第１および前記第２の記憶部での記憶データのアドレスを管理し、

前記第１の記憶部のディスクステータスが保護状態として認識された場合、前記第１の記憶部と前記第２の記憶部を同一の論理単位として認識し、

30

前記第１の記憶部のディスクステータスが保護状態として認識された場合、前記１乃至複数の第３の記憶部と前記論理単位とがＲＡＩＤを構成するように制御することを特徴とする請求項２２に記載の情報処理装置。

【請求項２７】

前記制御部は、

前記第１の記憶部に書き込まれているデータのうち前記第２の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履歴格納部に格納されている場合に、前記第１の記憶部の保護状態を終了すると判定することを特徴とする請求項２２乃至２６の何れか一つに記載の情報処理装置。

【請求項２８】

前記制御部は、

前記第３の処理では、

前記引き継ぎ履歴格納部に第１のアドレスを記録する場合に、前記第１のアドレスと重複または連続する第２のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第１のアドレスと前記第２のアドレスを合併した第３のアドレスを前記引き継ぎ履歴格納部に記録し、前記第２のアドレスを無効化することを特徴とする請求項２２乃至２７の何れか一つに記載の情報処理装置。

40

【請求項２９】

リードおよびライト可能な第１のメモリを有する第１の記憶部と、引き継ぎ履歴格納部が記憶され、リードおよびライト可能な第２のメモリを有する第２の記憶部とを接続可能

50

な情報処理装置にロードされる制御プログラムであって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出す第 2 の処理と、

前記第 2 の処理による書き込みに伴って、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第 3 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、読み出しは前記引き継ぎ履歴格納部を読み出すこと無く前記第 1 の記憶部からデータを読み出す第 4 の処理と、

を前記情報処理装置に実行させるように構成された制御プログラム。

【請求項 3 0】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 2 9 に記載の制御プログラム。

【請求項 3 1】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 2 9 に記載の制御プログラム。

【請求項 3 2】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうちの一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 2 9 に記載の制御プログラム。

【請求項 3 3】

前記制御部は、

前記第 1 の記憶部に書き込まれているデータのうち前記第 2 の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履歴格納部に格納されている場合に、前記第 1 の記憶部の保護状態を終了すると判定することを特徴とする請求項 2 9 乃至 3 2 の何れか一つに記載の制御プログラム。

【請求項 3 4】

前記制御部は、

前記第 3 の処理では、

前記引き継ぎ履歴格納部に第 1 のアドレスを記録する場合に、前記第 1 のアドレスと重

10

20

30

40

50

複または連続する第 2 のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第 1 のアドレスと前記第 2 のアドレスを合併した第 3 のアドレスを前記引き継ぎ履歴格納部に記録し、前記第 2 のアドレスを無効化することを特徴とする請求項 2 9 乃至 3 3 の何れか一つに記載の制御プログラム。

【請求項 3 5】

前記制御部は、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、前記引き継ぎ履歴格納部を消去することを特徴とする請求項 2 9 乃至 3 4 の何れか一つに記載の制御プログラム。

【請求項 3 6】

リードおよびライト可能な不揮発性の半導体メモリである第 1 のメモリを有する第 1 の記憶部と、引き継ぎ履歴格納部が記憶され、リードおよびライト可能な不揮発性の半導体メモリである第 2 のメモリを有する第 2 の記憶部と、引き継ぎ履歴格納部とを接続可能な情報処理装置にロードされる制御プログラムであって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出す第 2 の処理と、

前記第 2 の処理による書き込みに伴って、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第 3 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、読み出しは前記引き継ぎ履歴格納部を読み出すこと無く前記第 1 の記憶部からデータを読み出す第 4 の処理と、

を前記情報処理装置に実行させるように構成された制御プログラム。

【請求項 3 7】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 3 6 に記載の制御プログラム。

【請求項 3 8】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 3 6 に記載の制御プログラム。

【請求項 3 9】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうち的一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至

10

20

30

40

50

複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 3 6 に記載の制御プログラム。

【請求項 4 0】

前記制御部は、

前記第 1 の記憶部に書き込まれているデータのうち前記第 2 の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履歴格納部に格納されている場合に、前記第 1 の記憶部の保護状態を終了すると判定することを特徴とする請求項 3 6 乃至 3 9 の何れか一つに記載の制御プログラム。

【請求項 4 1】

前記制御部は、

前記第 3 の処理では、

前記引き継ぎ履歴格納部に第 1 のアドレスを記録する場合に、前記第 1 のアドレスと重複または連続する第 2 のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第 1 のアドレスと前記第 2 のアドレスを合併した第 3 のアドレスを前記引き継ぎ履歴格納部に記録し、前記第 2 のアドレスを無効化することを特徴とする請求項 3 6 乃至 4 0 の何れか一つに記載の制御プログラム。

【請求項 4 2】

前記制御部は、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識されなかった場合、前記引き継ぎ履歴格納部を消去することを特徴とする請求項 3 6 乃至 4 1 の何れか一つに記載の制御プログラム。

【請求項 4 3】

リードおよびライト可能な第 1 のメモリを有する第 1 の記憶部と、リードおよびライト可能な第 2 のメモリを有する第 2 の記憶部とを接続可能な情報処理装置にロードされる制御プログラムであって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態でないと判定された場合、前記第 1 の記憶部に書き込みを実行し、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 2 の記憶部に書き込みを実行する第 2 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された後に前記第 1 の記憶部に対する読み出しが発生した場合は、前記第 1 の記憶部からデータを読み出し、前記読み出されたデータを前記第 2 の記憶部に書き込む第 3 の処理と、を前記情報処理装置に実行させるように構成された制御プログラム。

【請求項 4 4】

前記制御部は、

前記第 3 の処理では、前記第 1 の記憶部からデータを読み出す処理と、前記読み出されたデータを前記第 2 の記憶部に書き込む処理を並行して実行することを特徴とする請求項 4 3 に記載の制御プログラム。

【請求項 4 5】

メタデータを格納するメタデータ格納部に接続され、

前記制御部は、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態でないと判定された場合、前記第 1 の記憶部に書き込まれたデータのアドレスと前記第 1 の記憶部の識別情報との紐付けを前記メタデータに記録し、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記保護状態と判定された場合、前記第 2 の記憶部に書き込まれたデータのアドレスと前記第 2 の記憶部の識別情報との紐付けを前記メタデータに記録し、

読み出し処理において、読み出し対象アドレスと前記第 1 の記憶部の識別情報との紐付けが前記メタデータに記録されている場合は、前記第 1 の記憶部からデータを読み出し、

10

20

30

40

50

読み出し対象アドレスと前記第 2 の記憶部との紐付けが前記メタデータに記録されている場合は、前記第 2 の記憶部からデータを読み出すことを特徴とする請求項 4 3 乃至 4 4 の何れか一つに記載の制御プログラム。

【請求項 4 6】

引き継ぎ履歴格納部に接続され、

前記制御部は、

前記第 3 の処理では、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された場合、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録し、

10

読み出し処理において、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記第 1 の記憶部からデータを読み出すことを特徴とする請求項 4 3 乃至 4 4 の何れか一つに記載の制御プログラム。

【請求項 4 7】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 4 3 に記載の制御プログラム。

20

【請求項 4 8】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 4 3 に記載の制御プログラム。

【請求項 4 9】

30

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうちの一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 4 3 に記載の制御プログラム。

40

【請求項 5 0】

リードおよびライト可能な第 1 のメモリを有する第 1 の記憶部と、リードおよびライト可能な第 2 のメモリを有する第 2 の記憶部と、引き継ぎ履歴格納部とを接続可能な情報処理装置にロードされる制御プログラムであって、

前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されている場合は、前記第 2 の記憶部からデータを読み出し、前記引き継ぎ履歴格納部に読み出し対象アドレスが記録されていない場合は、前記

50

第 1 の記憶部からデータを読み出す第 2 の処理と、

前記第 2 の処理による書き込みに伴って、前記第 2 の記憶部に書き込まれたデータのアドレスを前記引き継ぎ履歴格納部に記録する第 3 の処理と、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された後に、前記第 1 の記憶部に格納されたデータの削除処理命令を受信した場合に、前記削除処理命令による削除対象データのアドレスを前記引き継ぎ履歴格納部に記録する第 4 の処理と、

を前記情報処理装置に実行させるように構成された制御プログラム。

【請求項 5 1】

前記制御部は、

前記第 4 の処理では、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態と判定された後に、前記第 1 または第 2 の記憶部に格納されたデータの削除処理命令を受信した場合に、前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスが記録されている場合は、前記受信された削除処理命令を前記第 2 の記憶部に送信し、前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスを記録せず、

前記引き継ぎ履歴格納部に前記削除処理命令による削除対象データのアドレスが記録されていない場合は、前記削除処理命令による削除対象データのアドレスを前記引き継ぎ履歴格納部に記録する

ことを特徴とする請求項 5 0 に記載の制御プログラム。

【請求項 5 2】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを前記保護状態であると判定することを特徴とする請求項 5 0 に記載の制御プログラム。

【請求項 5 3】

前記制御部は、

前記第 1 の処理では、前記信頼性情報と閾値との比較を行い、この比較に基づいて前記第 1 の記憶部の寿命到達条件が成立したと判定したときに、前記第 1 の記憶部のディスクステータスを信頼性劣化状態であると判定し、

前記第 1 の処理により前記第 1 の記憶部のディスクステータスが前記信頼性劣化状態として判定された場合、前記第 2 の記憶部を前記第 1 の記憶部に格納されているデータのデータ引き継ぎ先状態として認識し、前記第 1 の記憶部のディスクステータスを前記保護状態に変更することを特徴とする請求項 5 0 に記載の制御プログラム。

【請求項 5 4】

リードおよびライト可能な第 3 のメモリを有する 1 乃至複数の第 3 の記憶部に接続可能であり、

前記制御部は、

同一アドレスのデータが前記第 1 および第 2 の記憶部のうち的一方に存在するように、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを管理し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記第 1 の記憶部と前記第 2 の記憶部を同一の論理単位として認識し、

前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、前記 1 乃至複数の第 3 の記憶部と前記論理単位とが R A I D を構成するように制御することを特徴とする請求項 5 0 に記載の制御プログラム。

【請求項 5 5】

前記制御部は、

前記第 1 の記憶部に書き込まれているデータのうち前記第 2 の記憶部へ引き継ぎ対象とするデータを選び、前記引き継ぎ対象のデータに対応するアドレスが全て前記引き継ぎ履

10

20

30

40

50

歴格納部に格納されている場合に、前記第 1 の記憶部の保護状態を終了すると判定することを特徴とする請求項 5 0 乃至 5 4 の何れか一つに記載の制御プログラム。

【請求項 5 6】

前記制御部は、

前記第 3 の処理では、

前記引き継ぎ履歴格納部に第 1 のアドレスを記録する場合に、前記第 1 のアドレスと重複または連続する第 2 のアドレスがすでに前記引き継ぎ履歴格納部に格納されている場合には、前記第 1 のアドレスと前記第 2 のアドレスを合併した第 3 のアドレスを前記引き継ぎ履歴格納部に記録し、前記第 2 のアドレスを無効化することを特徴とする請求項 5 0 乃至 5 5 の何れか一つに記載の制御プログラム。

10

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明の実施形態は、第 1 の記憶部のデータを第 2 の記憶部に引き継ぐための情報処理システムおよび制御プログラムに関する。

【背景技術】

【0 0 0 2】

フラッシュメモリを用いてパーソナルコンピュータやサーバの二次記憶装置を構成する場合、誤りが多いなどの理由によって記憶領域として使用できなくなる不良ブロック(defective block)や、読み出せなくなる領域(不良領域)などが生じる場合がある。このような不良ブロック数や不良領域数が上限値を超えた場合には、新たな不良ブロックや不良領域を登録することができないので、ライト要求のあったデータをフラッシュメモリへ書き込むことを保障できない。このため、不良ブロック数や不良領域数が上限値を超えた場合には、フラッシュメモリに空き容量があるにも関わらず、突然データの書き込みが不能になるといった問題があった。

20

【0 0 0 3】

そこで、記憶装置の劣化をモニタし、記憶装置の寿命到達より前に、情報処理装置から記憶装置へのデータの書き込みを制限する方法がある。この方法により、記憶装置の寿命到達により前にユーザデータを他の記憶装置にバックアップして引き継ぎ、記憶装置の寿命到達によるデータ損失を未然に防ぐことが可能である。

30

【0 0 0 4】

しかしながら、この方法では、ユーザデータを他の記憶装置にバックアップする作業があり、煩雑である。また、情報処理装置によってバックアップ元の記憶装置へのデータ書き込みが制限されるので、情報処理装置にロードされている各種アプリケーションプログラムの動作が制限され、アプリケーションプログラムの処理速度が低下する可能性がある。

【先行技術文献】

【特許文献】

【0 0 0 5】

【特許文献 1】米国特許出願公開第 2 0 1 2 / 0 2 8 4 4 5 3 号明細書

40

【特許文献 2】米国特許出願公開第 2 0 1 2 / 0 2 4 6 3 8 8 号明細書

【特許文献 3】米国特許出願公開第 2 0 0 9 / 0 2 2 2 6 1 7 号明細書

【特許文献 4】米国特許第 8 , 2 3 0 , 1 6 4 号明細書

【特許文献 5】米国特許第 7 , 8 6 1 , 1 2 2 号明細書

【特許文献 6】米国特許第 7 , 8 4 0 , 8 3 6 号明細書

【特許文献 7】米国特許出願公開第 2 0 0 7 / 0 2 1 4 1 8 7 号明細書

【特許文献 8】米国特許第 7 , 6 0 6 , 9 4 6 号明細書

【特許文献 9】米国特許第 6 , 5 2 9 , 9 9 2 号明細書

【特許文献 1 0】米国特許第 7 , 2 2 2 , 2 5 9 号明細書

【特許文献 1 1】米国特許第 7 , 4 2 4 , 6 3 9 号明細書

50

【特許文献 1 2】米国特許第 7,739,544 号明細書
【特許文献 1 3】米国特許第 7,525,749 号明細書
【特許文献 1 4】米国特許第 8,037,380 号明細書
【特許文献 1 5】米国特許出願公開第 2010/0325199 号明細書
【特許文献 1 6】米国特許出願公開第 2011/0239083 号明細書
【特許文献 1 7】特許第 5052376 号公報
【特許文献 1 8】特開 2010-225021 号公報
【特許文献 1 9】特許第 3,565,687 号公報
【特許文献 2 0】米国特許第 8,176,367 号明細書
【特許文献 2 1】米国特許出願公開第 2011/0197045 号明細書
【特許文献 2 2】特許第 4,643,456 号公報
【特許文献 2 3】特許第 4,764,490 号公報
【特許文献 2 4】特開 2011-209879 号公報
【特許文献 2 5】国際公開第 2013/027642 号パンフレット

10

【非特許文献】

【0006】

【非特許文献 1】Information technology ATA/ATAPI Command Set-3 (ACS-3), d2161r4-ATAATAPI_Command_Set_-_3.pdf, Sep.04.2012, <http://www.t13.org/>

【非特許文献 2】NVM Express Revision 1.1, Oct.11.2012, NVM_Express_1_1.pdf, <http://www.nvmexpress.org/>

20

【非特許文献 3】Serial ATA International Organization: Serial ATA Revision 3.1 Gold Revision, Jul.18.2011, <http://www.serialata.org/>

【非特許文献 4】PCI Express Base Specification Revision 3.0, Nov.10.2010, <http://www.pcisig.com/>

【非特許文献 5】Serial ATA Technical Proposal: SATA31_TPR_C108 Title: Device Sleep, <http://www.serialata.org/>

【非特許文献 6】SCSI Primary Commands-4 (SPC-4), INCITS T10/1731-D, Revision 36e, <http://www.t10.org/>

【非特許文献 7】喜連川優 編著、「よくわかるストレージネットワーキング」、第 1 版、株式会社オーム社、2011年5月20日、p.6-9、p67-93

30

【発明の概要】

【発明が解決しようとする課題】

【0007】

本発明の一つの実施形態は、ユーザによるバックアップ作業を行うことなくデータ引き継ぎをなし得る情報処理システムおよび制御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】

【0008】

本発明の一つの実施形態によれば、情報処理システムは、リードおよびライト可能な第 1 のメモリを有する第 1 の記憶部と、リードおよびライト可能な第 2 のメモリを有する第 2 の記憶部と、前記第 1 の記憶部と前記第 2 の記憶部に接続可能な情報処理装置とを備える。前記情報処理装置は、前記第 1 の記憶部から取得した信頼性情報に基づき、前記第 1 の記憶部のディスクステータスを判定する第 1 の処理と、前記第 1 の処理により前記第 1 の記憶部のディスクステータスが保護状態として認識された場合、書き込みは前記第 2 の記憶部に対し実行し、読み出しは前記第 1 および前記第 2 の記憶部に対し実行する第 2 の処理と、前記第 2 の処理による書き込みに伴って、前記第 1 および前記第 2 の記憶部での記憶データのアドレスを更新する第 3 の処理とを実行する制御部を備えることを特徴とする。

40

【図面の簡単な説明】

【0009】

【図 1】図 1 は、第 1 の実施形態の情報処理システムの機能構成例を示すブロック図であ

50

る。

【図 2】図 2 は、制御プログラムを S S D に保存する場合の情報処理システムの機能構成例を示すブロック図である。

【図 3】図 3 は、制御プログラムをその他の外部記憶装置に保存する場合の情報処理システムの機能構成例を示すブロック図である。

【図 4】図 4 は、制御プログラムを W E B からインストールする場合の情報処理システムの機能構成例を示すブロック図である。

【図 5】図 5 は、制御プログラムを光学ドライブからインストールする場合の情報処理システムの機能構成例を示すブロック図である。

【図 6】図 6 は、制御プログラムを U S B メモリからインストールする場合の情報処理システムの機能構成例を示すブロック図である。

10

【図 7】図 7 は、情報処理装置の階層的機能構成例を示すブロック図である。

【図 8】図 8 は、情報処理システムの外観的構成を示す図である。

【図 9】図 9 は、情報処理システムの他の外観的構成を示す図である。

【図 10】図 10 は、N A N D メモリチップの機能構成例を示すブロック図である。

【図 11】図 11 は、N A N D メモリチップに含まれる 1 個のプレーンの構成例を示す回路図である。

【図 12】図 12 は、4 値データ記憶方式での閾値分布を示す図である。

【図 13】図 13 は、S S D の機能構成例を示すブロック図である。

【図 14】図 14 は、S S D の管理情報を示す図である。

20

【図 15】図 15 は、L B A と S S D の管理単位との関係を示す図である。

【図 16】図 16 は、L B A から物理アドレスを特定する手順を示すフローチャートである。

【図 17】図 17 は、S S D の読み出し動作例を示すフローチャートである。

【図 18】図 18 は、S S D の読み出し動作例を示すフローチャートである。

【図 19】図 19 は、S S D の書き込み動作例を示すフローチャートである。

【図 20】図 20 は、S S D の書き込み動作例を示すフローチャートである。

【図 21】図 21 は、N A N D 整理を行う場合の動作手順を示すフローチャートである。

【図 22】図 22 は、削除通知を受信したときの S S D の動作例を示すフローチャートである。

30

【図 23】図 23 は、エラー発生時の S S D の動作例を示すフローチャートである。

【図 24】図 24 は、統計情報取得処理を示すフローチャートである。

【図 25】図 25 は、パッド論理セクタテーブルを示す図である。

【図 26】図 26 は、パッドクラスタテーブルを示す図である。

【図 27】図 27 は、統計情報の例としての S M A R T 情報を示す図である。

【図 28】図 28 は、統計情報の Raw Value と S S D の不良率の関係を示すグラフである。

。

【図 29】図 29 は、コマンド応答エラー処理を示すフローチャートである。

【図 30】図 30 は、他のコマンド応答エラー処理を示すフローチャートである。

【図 31】図 31 は、寿命到達時処理を行うより以前における情報処理装置が管理するデータの構成例を示す図である。

40

【図 32】図 32 は、Boot Loader の内部データを示す図である。

【図 33】図 33 は、メタデータの構成例を示す図である。

【図 34】図 34 は、アプリケーションプログラムが O S に対して論理ドライブへのアクセス要求を行った場合のフローチャートである。

【図 35】図 35 は、情報処理装置が起動した時などの制御プログラムの処理手順を示すフローチャートである。

【図 36】図 36 は、記憶部のライフサイクルを示す状態遷移図である。

【図 37】図 37 は、記憶部のライフサイクルを示す状態遷移図である。

【図 38】図 38 は、データ引き継ぎ元記憶部及びデータ引き継ぎ先記憶部が接続された

50

状態での情報処理装置が管理するデータの構成例を示す図である。

【図 3 9】図 3 9 は、記憶部のライフサイクルを示す状態遷移図である。

【図 4 0】図 4 0 は、寿命到達時処理において、制御プログラムが行う処理を示すフローチャートである。

【図 4 1】図 4 1 は、論理ドライブステータステーブルを示す図である。

【図 4 2】図 4 2 は、アプリケーションプログラムから O S に書き込み要求が送信された時の、O S の処理手順を示すフローチャートである。

【図 4 3】図 4 3 は、アプリケーションプログラムから O S に削除要求が送信された時の、O S の処理手順を示すフローチャートである。

【図 4 4】図 4 4 は、アプリケーションプログラムから O S に読み出し要求が送信された時の、O S の処理手順を示すフローチャートである。

10

【図 4 5】図 4 5 は、バックグラウンドバックアップの処理手順を示すフローチャートである。

【図 4 6】図 4 6 は、データ引き継ぎ完了時の制御プログラムの動作手順を示すフローチャートである。

【図 4 7】図 4 7 は、データ引き継ぎ状態における論理ドライブからの読み出しを概念的に示す図である。

【図 4 8】図 4 8 は、データ引き継ぎ状態における論理ドライブへの書き込みを概念的に示す図である。

【図 4 9】図 4 9 は、第 2 の実施形態の情報処理システムの機能構成例を示すブロック図である。

20

【図 5 0】図 5 0 は、引き継ぎ履歴を示す図である。

【図 5 1】図 5 1 は、引き継ぎ履歴に対する履歴の書き込み例を示す図である。

【図 5 2】図 5 2 は、寿命到達時処理を示すフローチャートである。

【図 5 3】図 5 3 は、論理ドライブからの読み出しの動作手順を示すフローチャートである。

【図 5 4】図 5 4 は、アプリケーションプログラムから O S に L B A データ削除要求が送信された時の、O S の処理手順を示す図である。

【図 5 5】図 5 5 は、論理ドライブへの書き込みの動作手順を示すフローチャートである。

30

【図 5 6】図 5 6 は、引き継ぎ履歴を用いたデータ引き継ぎ状況の監視手順を示すフローチャートである。

【図 5 7】図 5 7 は、記憶部が寿命到達する際の論理ドライブのステータスの遷移を示す図である。

【図 5 8】図 5 8 は、第 3 の実施形態の情報処理システムの機能構成例を示すブロック図である。

【図 5 9】図 5 9 は、通常状態のアレイ L B A (A L B A) と記憶部 L B A (S L B A) の変換方法を示す図である。

【図 6 0】図 6 0 は、引き継ぎ状態遷移直後のアレイ L B A (A L B A) と記憶部 L B A (S L B A) の変換方法を示す図である。

40

【図 6 1】図 6 1 は、引き継ぎ履歴を示す図である。

【図 6 2】図 6 2 は、寿命到達時処理を示すフローチャートである。

【図 6 3】図 6 3 は、R A I D コントローラが起動した時などの制御部の処理手順を示すフローチャートである。

【図 6 4】図 6 4 は、読み出し要求を実行する際の、R A I D コントローラの処理手順を示すフローチャートである。

【図 6 5】図 6 5 は、読み出し要求を実行する際の、R A I D コントローラの他の処理手順を示すフローチャートである。

【図 6 6】図 6 6 は、書き込みし要求を実行する際の、R A I D コントローラの他の処理手順を示すフローチャートである。

50

【図 6 7】図 6 7 は、書き込み処理の具体例を示す図である。

【図 6 8】図 6 8 は、引き継ぎ履歴を用いたデータ引き継ぎ状況の監視手順を示すフローチャートである。

【図 6 9】図 6 9 は、各記憶部のステータスの遷移を示す図である。

【図 7 0】図 7 0 は、1 台のデータ引き継ぎ状態と他の 1 台の故障が同時に発生した場合の例を示す図である。

【図 7 1】図 7 1 は、2 台のデータ引き継ぎ状態と他の 1 台の故障が同時に発生した場合の例を示す図である。

【図 7 2】図 7 2 は、1 台のデータ引き継ぎ状態において、他の 1 台の記憶部で ECC 訂正不能エラーが発生した場合の例を示す図である。

10

【図 7 3】図 7 3 は、第 4 の実施形態の情報処理システムの機能構成例を示すブロック図である。

【図 7 4】図 7 4 は、第 4 の実施形態の他の情報処理システムの機能構成例を示すブロック図である。

【図 7 5】図 7 5 は、ステータステーブルを示す図である。

【図 7 6】図 7 6 は、第 5 の実施形態の情報処理システムの機能構成例を示すブロック図である。

【図 7 7】図 7 7 は、第 6 の実施形態の情報処理システムの機能構成例を示すブロック図である。

【図 7 8】図 7 8 は、第 6 の実施形態の情報処理システムの他の機能構成例を示すブロック図である。

20

【図 7 9】図 7 9 は、制御部の通常状態から信頼性劣化モードへの遷移操作を示すフローチャートである。

【図 8 0】図 8 0 は、制御部が CPU から記憶部の記憶部情報要求を受信した場合の処理手順の例を示すフローチャートである。

【図 8 1】図 8 1 は、制御部が CPU から記憶部の記憶部情報要求を受信した場合の他の処理手順の例を示すフローチャートである。

【図 8 2】図 8 2 は、情報処理装置でのドライブ表示画面例を示す図である。

【図 8 3】図 8 3 は、情報処理装置での他のドライブ表示画面例を示す図である。

【図 8 4】図 8 4 は、情報処理装置での他の表示画面例を示す図である。

30

【図 8 5】図 8 5 は、第 6 の実施形態の情報処理システムの他の機能構成例を示すブロック図である。

【発明を実施するための形態】

【0010】

以下に添付図面を参照して、実施形態にかかる情報処理システムおよび制御プログラムを詳細に説明する。なお、これら実施形態により本発明が限定されるものではない。

【0011】

(第 1 の実施形態)

(システムの構成)

図 1 に情報処理システムの一例としてのコンピュータシステムの第 1 の実施形態の構成を示す。情報処理システム 1 は、情報処理装置 1 1 1 と、1 乃至複数の記憶部と、情報処理装置 1 1 1 と記憶部とを接続するインタフェース 1 9 とを含む。本実施形態では、記憶部として、不揮発性記憶装置である SSD (Solid State Drive) を用いるが、たとえば、ハードディスクドライブ (HDD)、ハイブリッドディスクドライブ、SD カード、USB メモリ、NAND 型フラッシュメモリチップ、磁気テープなど、他の記憶装置であってもよい。1 つの記憶装置に含まれる複数の記憶領域をそれぞれ別の記憶部として用いてもよい。また、本実施形態では、インタフェース 1 9 として、SATA (Serial Advanced Technology Attachment) インタフェースを用いた場合について説明するが、PCI Express (Peripheral Component Interconnect Express, PCIe)、USB (Universal Serial Bus)、SAS (Serial Attached SCSI)、Thunderbolt (登録商標)、イーサネ

40

50

ット（登録商標）、ファイバーチャネルなどが使用可能である。CPU（制御回路）5が情報処理装置111における中央演算処理装置であり、情報処理装置111における種々の演算及び制御はCPU5によって行われる。CPU5とチップセット7はたとえばDMI（Direct Media Interface）などのインタフェースで接続され、CPU5はチップセット7を介して記憶部2や、DVDドライブなどの光学ドライブ10の制御を行う。CPU5は主メモリ6の制御を行う。主メモリ6としては、たとえばDRAM(Dynamic Random Access Memory)やMRAM(Magnetoresistive Random Access Memory)やReRAM(Resistance Random Access Memory)やFeRAM(Ferroelectric Random Access Memory)を採用してよい。

【0012】

10

ユーザは、キーボード14やマウス15などの入力装置を通して情報処理装置111の制御を行い、キーボード14やマウス15からの信号は、たとえばUSB(Universal Serial Bus)コントローラ13及びチップセット7を介してCPU5で処理される。CPU5は、表示コントローラ8を介してディスプレイ（表示装置）9に画像データやテキストデータなどを送る。ユーザは、ディスプレイ9を介して情報処理装置111からの画像データやテキストデータなどを視認することができる。

【0013】

CPU5は、情報処理装置111の動作を制御するために設けられたプロセッサであり、たとえば記憶部2から主メモリ6にロードされるオペレーティングシステム(OS)100を実行する。更に、光学ドライブ10が、装填された光ディスクに対して読出し処理及び書込み処理の少なくとも1つの処理の実行を可能にした場合に、CPU5は、それらの処理を実行する。また、CPU5は、ROM11に格納されたUEFI(Unified Extensible Firmware Interface)ファームウェアやシステムBIOS(Basic Input/Output System)などを実行する。尚、UEFIファームウェアやシステムBIOSは、情報処理装置111内のハードウェア制御のためのプログラムである。その他、CPU5は、チップセット7を介してネットワークコントローラ12を制御する。ネットワークコントローラ12としては、たとえばLAN(Local Area Network)コントローラ、無線LANコントローラなどがある。

20

【0014】

主メモリ6は、プログラムやデータを一時的に記憶し、CPU5の作業用メモリとして機能する。主メモリ6は、OS100を格納する記憶領域6Aと、制御プログラム200を格納する記憶領域6Bを含んでいる。OS100は、たとえばWindows（登録商標）として一般的に知られているように、情報処理装置111の入出力装置を管理し、ディスクやメモリを管理し、ソフトウェアが情報処理装置111のハードウェアや記憶装置を利用可能にするための制御を行うなど、情報処理装置111全体を管理するプログラムである。本実施形態では、制御プログラム200はOS100の一部の場合である例を示すが、制御プログラム200はOS100とは独立したプログラムであってもよい。また、本実施形態では制御プログラム200はインタフェース19を介して主メモリ100に読み出されて実行されるプログラムとして説明するが、制御プログラム200はROM11に保存され、UEFIファームウェアやシステムBIOSの一部として動作してもよい。また制御プログラム200はハードウェアとして情報処理装置111内に実装されてもよい。主メモリ6には、メタデータ300と、論理ドライブステータステーブル450が格納されている。

30

40

【0015】

表示コントローラ8は、情報処理装置111のディスプレイ9を制御するビデオ再生コントローラである。チップセット7は、CPU5のローカルバスと接続するブリッジデバイスである。チップセット7は、各種ソフトウェア及びデータを格納する記憶装置である記憶部2や記憶部3を、インタフェース19を介して制御する。なお、記憶部2および記憶部3は、チップセット7を経由してCPU5に接続されてもよいし、CPU5に直接接続されてもよい。

50

【 0 0 1 6 】

情報処理装置 1 1 1 では、論理セクタ単位で記憶部 2 や記憶部 3 へのアクセスを行う。インタフェース 1 9 を介して、書き込みコマンド（ライトコマンド、ライト要求）、読み出しコマンド（リードコマンド、リード要求）、フラッシュコマンド等が記憶部 2 や記憶部 3 に入力される。

【 0 0 1 7 】

また、チップセット 7 は、ROM 1 1、光学ドライブ 1 0、ネットワークコントローラ 1 2、USB コントローラ 1 3 をアクセス制御するための機能も有している。USB コントローラ 1 3 にキーボード 1 4、マウス 1 5 が接続されている。

【 0 0 1 8 】

（制御プログラムの形態）

本実施形態では、制御プログラム 2 0 0 は、例えば、図 2 に示すように、情報処理装置 1 1 1 が電源オフになっているときは記憶部 2 の NAND 型フラッシュメモリ（NAND メモリ）1 6 の領域 1 6 B に格納されているが、情報処理装置 1 1 1 の起動時またはプログラム起動時に、NAND メモリ 1 6 の領域 1 6 B から主メモリ 6 上の領域 6 B にロードされる。一方、情報処理装置 1 1 1 に記憶部 2 や記憶部 3 とは別の記憶部 2 0 が接続されている場合など、複数の記憶部が接続されている場合は、図 3 に示すように、制御プログラム 2 0 0 は、記憶部 2 0 の領域 2 0 B に格納されており、情報処理装置 1 1 1 の起動時またはプログラム起動時に、領域 2 0 B から主メモリ 6 上の領域 6 B にロードされるようにしてもよい。特に、記憶部 2 0 が OS を格納するシステムドライブとして使用されており、記憶部 2 が文書や静止画データや動画データなどユーザデータを格納するデータドライブとして使用されている場合は、システムドライブとしての記憶部 2 0 には OS やアプリケーションプログラムを主として格納する記憶ドライブとして使用し、データドライブとしての記憶部 2 にはユーザデータを格納する記憶ドライブとして使用するというように、記憶部 2 と記憶部 2 0 の役割を明確に分ける観点で、システムドライブとしての記憶部 2 0 に制御プログラム 2 0 0 を格納することが望ましい。

【 0 0 1 9 】

制御プログラム 2 0 0 のセットアップをユーザが行う労力を省く観点では、例えば図 2 や図 3 に示したように、制御プログラム 2 0 0 が記憶部 2 や記憶部 2 0 に格納された状態で情報処理システム 1 が製造元から出荷され、店頭に並び、ユーザの手に渡ることが望ましい。一方、ユーザが制御プログラム 2 0 0 のインストールの可否を選択できるようにする観点、およびユーザに最新の制御プログラムを提供できるようにするという観点では、制御プログラム 2 0 0 は、WEB からのダウンロード、または DVD-ROM や USB メモリなど外部記憶媒体からのインストールにより記憶部 2 や記憶部 2 0 に格納できるようにすることが望ましい。

【 0 0 2 0 】

図 4 は WEB からのダウンロードの一例である。制御プログラム 2 0 0 は WEB サーバ 2 1 内の記憶媒体 2 2 の領域 2 2 B に格納されており、制御プログラム 2 0 0 はインターネットやローカルネットワークや無線 LAN などのネットワークを介して、たとえばネットワークコントローラ 1 2 を経由して記憶部 2 の NAND メモリ 1 6 上の領域 1 6 B にダウンロード（またはインストール）される。記憶媒体 2 2 は、たとえば SSD、HDD、ハイブリッドディスクドライブ、磁気テープなどが用いられる。なお、図 3 の場合は、制御プログラム 2 0 0 は、記憶部 2 0 上の領域 2 0 B にダウンロードまたはインストールされる。

【 0 0 2 1 】

図 5 は光学メディアからのインストールの一例である。制御プログラム 2 0 0 は DVD-ROM や CD-ROM や DVD-RW や Blu-ray（登録商標）Disc などの光学メディア 2 3 に格納されており、光学メディア 2 3 が光学ドライブ 1 0 にセットされることで、制御プログラム 2 0 0 は光学ドライブ 1 0 を介して記憶部 2 の NAND メモリ 1 6 上の領域 1 6 B（または領域 2 0 B）にインストールされる。

【 0 0 2 2 】

図 6 は U S B メモリからのインストールの一例である。制御プログラム 2 0 0 は U S B メモリ 2 4 の領域 2 4 B に格納されており、U S B メモリ 2 4 が U S B コントローラ 1 3 に接続されることで、制御プログラム 2 0 0 は U S B コントローラ 1 3 を介して記憶部 2 の N A N D メモリ 1 6 上の領域 1 6 B (または領域 2 0 B) にインストールされる。U S B メモリ 2 4 のかわりに、S D カードなどその他外部メモリであってもよいのはもちろんである。ユーザによる入手容易性の観点から、光学メディア 2 3 や U S B メモリ 2 4 は、情報処理システム 1 または記憶部 2 の出荷時に、付属品として情報処理システム 1 または記憶部 2 と一緒に梱包されて販売されることが望ましい。一方、光学メディア 2 3 や U S B メモリ 2 4 はソフトウェア商品として単独で販売されてもよいし、雑誌や書籍の付録として付属していてもよい。

10

【 0 0 2 3 】

(ソフトウェアの階層構造)

図 7 に、情報処理装置 1 1 1 のソフトウェアレベルでの階層構造を示す。主メモリ 6 上にロードされたアプリケーションプログラム 4 0 0 は、通常は記憶部 2 や記憶部 3 などの記憶部と直接通信せずに、主メモリ 6 にロードされた O S 1 0 0 を経由して記憶部と通信する。また、O S 1 0 0 は U E F I ファームウェアやシステム B I O S を経由して記憶部と通信する。後述する寿命到達時処理前において、O S 1 0 0 は記憶部 2 を論理単位としての論理ドライブ 4 (図 3 8 参照)として認識し、アプリケーションプログラム 4 0 0 に対して論理ドライブ 4 をアクセス可能な記憶ドライブとして通知する。アプリケーションプログラム 4 0 0 が論理ドライブ 4 に対しリード要求、ライト要求などの命令を送信する必要がある場合、アプリケーションプログラム 4 0 0 は、O S 1 0 0 に論理ドライブ 4 に対するファイル単位でのアクセス要求を送信する。O S 1 0 0 は主メモリ 6 に格納されたメタデータ 3 0 0 を参照し、アクセス要求のあったファイルに対応する記憶部 2 の論理アドレス (L B A : Logical Block Address) を特定し、記憶部 2 に対してインタフェース 1 9 を介して命令および L B A およびデータを送信する。記憶部 2 から応答が返ってきた場合、O S 1 0 0 はアプリケーションプログラム 4 0 0 に対して応答を返す。論理ドライブ 4 が後述する引き継ぎ状態になった時には、O S 1 0 0 は記憶部 2 および記憶部 3 を論理ドライブ 4 として認識する。

20

【 0 0 2 4 】

(情報処理装置の構成)

次に、情報処理システム 1 の構成例について説明する。図 8 は、情報処理システム 1 としてのデスクトップコンピュータの概略図である。デスクトップコンピュータは、コンピュータ本体 3 1、ディスプレイ 9、キーボード 1 4、及びマウス 1 5などを備えている。コンピュータ本体 3 1は、主要なハードウェアが搭載されたマザーボード 3 0、記憶部 2、及び電源装置 3 2などを備えている。記憶部 2 は、S A T A ケーブルを介してマザーボード 3 0 に物理的に接続され、マザーボード 3 0 上に実装されたチップセット 7 を介して、同じくマザーボード上に実装された C P U 5 に電氣的に接続されている。電源装置 3 2 は、デスクトップコンピュータで使用される各種電源を発生し、電源ケーブルを介してマザーボード 3 0 や記憶部 2 などに電源を供給する。記憶部 3 は、S A T A ケーブルを介してマザーボード 3 0 に物理的に接続可能であり、それによりマザーボード 3 0 上に実装されたチップセット 7 を介して、同じくマザーボード上に実装された C P U 5 に電氣的に接続される。

30

40

【 0 0 2 5 】

図 9 は、情報処理システム 1 としてのポータブルコンピュータの概略図である。ポータブルコンピュータは、コンピュータ本体 3 4、ディスプレイユニット 3 5 などから構成されている。ディスプレイユニット 3 5 は、例えば L C D (Liquid Crystal Display) で構成される表示装置 9 が組み込まれている。ディスプレイユニット 3 5 は、コンピュータ本体 3 4 に対し、この本体 3 4 の上面が露出される開放位置と本体 3 4 の上面を覆う閉塞位置との間を回動自由に取り付けられている。本体 3 4 は薄い箱形の筐体を有しており、そ

50

の上面には、電源スイッチ 36、キーボード 14、タッチパッド 33 等が配置されている。また、本体 34 も、デスクトップコンピュータと同様に、記憶部 2、マザーボード、及び電源装置などを備えている。

【0026】

本発明を適用する情報処理システム 1 としては、他に、スチルカメラ或いはビデオカメラなどの撮像装置などであってもよいし、タブレットコンピュータやスマートフォンやゲーム機器やカーナビゲーションシステムやプリンタ機器やスキャナ機器やサーバーシステムなどであってもよい。

【0027】

(記憶装置の構成 その 1)

次に、記憶部 2 の構成について説明する。本実施形態では、記憶部 2 の主な構成要素が NAND メモリ 16 である場合について説明する。本実施形態では、記憶部 2 と記憶部 3 が同一の構成である場合について説明するが、一方、記憶部 2 は主な構成要素を NAND メモリ 16 とする SSD であり、記憶部 3 は主な構成要素を磁気ディスクとするハードディスクドライブである場合など、記憶部 3 が記憶部 2 とは異なる構成であっても本発明は適用可能である。記憶部 2 と記憶部 3 が統計情報 65 を格納していることが望ましいが、記憶部 2 が統計情報 65 を格納しており、記憶部 3 は統計情報 65 を格納していない場合であっても本発明は適用可能である。データ引き継ぎ先である記憶部 3 の記憶可能容量は、データ引き継ぎ元である記憶部 2 の記憶可能容量と同じか、記憶部 2 の記憶可能容量よりも大きいことが望ましいが、記憶部 2 の記憶可能容量よりも小さい場合であっても本発明は適用可能である。

【0028】

(NAND メモリの構成)

図 10 に、本実施形態において記憶部 2 および記憶部 3 の構成要素として使用される NAND メモリ 16 を構成する NAND メモリチップ 80 の内部構成例を示す。NAND メモリ 16 は、1 以上の数の NAND メモリチップ 80 よりなる。NAND メモリチップ 80 は、複数のメモリセルがマトリクス状に配列されたメモリセルアレイを有する。メモリセルアレイを構成するメモリセルトランジスタは、半導体基板上に形成された積層ゲート構造を備えた MOSFET (Metal Oxide Semiconductor Field Effect Transistor) から構成される。積層ゲート構造は、半導体基板上にゲート絶縁膜を介在して形成された電荷蓄積層 (浮遊ゲート電極)、及び浮遊ゲート電極上にゲート間絶縁膜を介在して形成された制御ゲート電極を含んでいる。メモリセルトランジスタは、浮遊ゲート電極に蓄えられる電子の数に応じて閾値電圧が変化し、この閾値電圧の違いに応じてデータを記憶する。

【0029】

本実施形態では、個々のメモリセルが上位ページ及び下位ページを使用して 2 bit/cell の 4 値記憶方式の書き込み方式である場合について説明するが、個々のメモリセルが単一ページを使用しての 1 bit/cell の 2 値記憶方式の書き込み方式、または上位ページ及び中位ページ及び下位ページを使用しての 3 bit/cell の 8 値記憶方式の書き込み方式である場合、あるいは 4 bit/cell 以上の多値記憶方式の書き込み方式を採用する場合であっても本発明の本質は変わらない。また、メモリセルトランジスタは浮遊ゲート電極を有する構造に限らず、MONOS (Metal-Oxide-Nitride-Oxide-Silicon) 型など、電荷蓄積層としての窒化界面に電子をトラップさせることで閾値電圧を調整可能な構造であってもよい。MONOS 型のメモリセルトランジスタについても同様に、1 ビットを記憶するように構成されていてもよいし、多値を記憶するように構成されていてもよい。また、不揮発性記憶媒体として、米国特許 8,189,391 号明細書や米国特許出願公開第 2010/0207195 号明細書や米国特許出願公開第 2010/0254191 号明細書に記述されるような 3 次元的にメモリセルが配置された半導体記憶媒体であってもよい。

【0030】

図 10 に示すように、NAND メモリチップ 80 は、データを記憶するメモリセルをマ

10

20

30

40

50

トリックス状に配置してなるメモリセルアレイ 8 2 を備えている。メモリセルアレイ 8 2 は複数のビット線と複数のワード線と共通ソース線を含み、ビット線とワード線の交点に電氣的にデータを書き換え可能なメモリセルがマトリクス状に配置されている。このメモリセルアレイ 8 2 には、ビット線を制御するためのビット線制御回路 8 3、及びワード線電圧を制御するためのワード線制御回路 8 5 が接続されている。すなわち、ビット線制御回路 8 3 は、ビット線を介してメモリセルアレイ 8 2 中のメモリセルのデータを読み出す一方、ビット線を介してメモリセルアレイ 8 2 中のメモリセルに書き込み制御電圧を印加してメモリセルに書き込みを行う。

【 0 0 3 1 】

ビット線制御回路 8 3 には、カラムデコーダ 8 4、データ入出力バッファ 8 9 及びデータ入出力端子 8 8 が接続されている。メモリセルアレイ 8 2 から読み出されたメモリセルのデータは、ビット線制御回路 8 3、データ入出力バッファ 8 9 を介してデータ入出力端子 8 8 から外部へ出力される。また、外部からデータ入出力端子 8 8 に入力された書き込みデータは、データ入出力バッファ 8 9 を介して、カラムデコーダ 8 4 によってビット線制御回路 8 3 に入力され、指定されたメモリセルへの書き込みが行われる。

【 0 0 3 2 】

また、メモリセルアレイ 8 2、ビット線制御回路 8 3、カラムデコーダ 8 4、データ入出力バッファ 8 9、及びワード線制御回路 8 5 は、制御回路 8 6 に接続されている。制御回路 8 6 は、制御信号入力端子 8 7 に入力される制御信号に従い、メモリセルアレイ 8 2、ビット線制御回路 8 3、カラムデコーダ 8 4、データ入出力バッファ 8 9、及びワード線制御回路 8 5 を制御するための制御信号及び制御電圧を発生させる。NANDメモリチップ 8 0 のうち、メモリセルアレイ 8 2 以外の回路部分を NAND コントローラ (NANDC) 8 1 と呼ぶ。

【 0 0 3 3 】

図 1 1 は、図 1 0 に示すメモリセルアレイ 8 2 の構成を示している。メモリセルアレイ 8 2 は NAND セル型メモリセルアレイであり、複数の NAND セルを含んで構成されている。1 つの NAND セルは、直列接続されたメモリセルからなるメモリストリング MS と、その両端に接続される選択ゲート S 1、S 2 とにより構成されている。選択ゲート S 1 はビット線 BL に接続され、選択ゲート S 2 はソース線 SRC に接続されている。同一のロウに配置されたメモリセル MC の制御ゲートはワード線 WL 0 ~ WL m - 1 に共通接続されている。また、第 1 の選択ゲート S 1 はセレクト線 SGD に共通接続され、第 2 の選択ゲート S 2 はセレクト線 SGS に共通接続されている。

【 0 0 3 4 】

メモリセルアレイ 8 2 は、1 または複数のプレーンを含んでおり、プレーンは複数のブロックを含んでいる。各ブロックは、複数の NAND セルにより構成され、このブロック単位でデータが消去される。

【 0 0 3 5 】

また、1 つのワード線に接続された複数のメモリセルは、1 物理セクタを構成する。この物理セクタ毎にデータが書き込まれ、読み出される。この物理セクタは、後述する LBA の論理セクタとは無関係である。1 物理セクタには、2 bit/cell 書き込み方式 (4 値) の場合例えば 2 物理ページ (2 ページ) 分のデータが記憶される。一方、1 bit/cell 書き込み方式 (2 値) の場合は、1 物理セクタに例えば 1 物理ページ (1 ページ) 分のデータが記憶され、3 bit/cell 書き込み方式 (8 値) の場合、1 物理セクタに例えば 3 物理ページ (3 ページ) 分のデータが記憶される。

【 0 0 3 6 】

リード動作、プログラムベリファイ動作及びプログラム動作時において、後述の SSDC 4 1 から受信した例えば Row Address 等の物理アドレスに応じて、1 本のワード線が選択され、1 物理セクタが選択される。この物理セクタ内のページの切り替えは物理アドレスによって行われる。本実施形態では NAND メモリ 1 6 は 2 bit/cell 書き込み方式であり、SSDC 4 1 は物理セクタには上位ページ (Upper Page) と下位ページ (Lower Page

）の２ページが物理ページとして割り当てられているとして取扱い、それら全ページに対して物理アドレスが割り当てられている。

【 0 0 3 7 】

2 bit/cell の 4 値 N A N D メモリは、１つのメモリセルにおける閾値電圧が、４通りの分布を持ち得るように構成されている。図 1 2 は、４値 N A N D セル型フラッシュメモリのメモリセルに記憶される 2 ビットの 4 値データ（データ “ 1 1 ”、“ 0 1 ”、“ 1 0 ”、“ 0 0 ”）とメモリセルの閾値電圧分布との関係を示している。なお、図 1 2 において、V A 1 は、下位ページのみ書き込み済みで上位ページが未書き込みの物理セクタについて、２つのデータを読み出す場合に選択ワード線に印加される電圧であり、V A 1 V は、A 1 への書き込みを行う場合において、書き込みが完了したかどうかを確認するために印加されるベリファイ電圧を示している。

10

【 0 0 3 8 】

また、V A 2、V B 2、V C 2 は、下位ページと上位ページが書き込み済みの物理セクタについて、４つのデータを読み出す場合に選択ワード線に印加される電圧であり、V A 2 V、V B 2 V、V C 2 V は、各閾値電圧分布への書き込みを行う場合において、書き込みが完了したかどうかを確認するために印加されるベリファイ電圧を示している。また、V read 1、V read 2 は、データの読み出しを行う場合に、N A N D セル中の非選択メモリセルに対し印加され、その保持データに拘わらず当該非選択メモリセルを導通させる読み出し電圧を示している。さらに、V ev 1、V ev 2 は、メモリセルのデータを消去する場合において、その消去が完了したか否かを確認するためメモリセルに印加される消去ベリファイ電圧であり、負の値を有する。その大きさは、隣接メモリセルの干渉の影響を考慮して決定される。上述の各電圧の大小関係は、

20

$$V_{ev1} < V_{A1} < V_{A1V} < V_{read1}$$

$$V_{ev2} < V_{A2} < V_{A2V} < V_{B2} < V_{B2V} < V_{C2} < V_{C2V} < V_{read2}$$

である。

【 0 0 3 9 】

なお、消去ベリファイ電圧 V ev 1、V ev 2、V ev 3 は前述の通り負の値であるが、実際に消去ベリファイ動作においてメモリセル M C の制御ゲートに印加される電圧は、負の値ではなく、ゼロ又は正の値である。すなわち、実際の消去ベリファイ動作においては、メモリセル M C のバックゲートに正の電圧を与え、メモリセル M C の制御ゲートには、ゼロ又はバックゲート電圧より小さい正の値の電圧を印加している。換言すれば、消去ベリファイ電圧 V ev 1、V ev 2、V ev 3 は、等価的に負の値を有する電圧である。

30

【 0 0 4 0 】

ブロック消去後のメモリセルの閾値電圧分布 E R は、その上限値も負の値であり、データ “ 1 1 ” が割り当てられる。下位ページおよび上位ページ書き込み状態のデータ “ 1 1 ”、“ 0 1 ”、“ 1 0 ”、“ 0 0 ” のメモリセルは、それぞれ正の閾値電圧分布 E R 2、A 2、B 2、C 2 を有し（A 2、B 2、C 2 の下限値も正の値である）、データ “ 0 1 ” の閾値電圧分布 A 2 が最も電圧値が低く、データ “ 0 0 ” の閾値電圧分布 C 2 が最も電圧値が高く、各種閾値電圧分布の電圧値は A 2 < B 2 < C 2 の関係を有する。下位ページ書き込みかつ上位ページ未書き込み状態のデータ “ 1 0 ” のメモリセルは、正の閾値電圧分布 A 1 を有する（A 1 の下限値も正の値である）。なお、図 1 2 に示す閾値電圧分布はあくまでも一例であって、本発明はこれに限定されるものではない。例えば、図 1 2 では閾値電圧分布 A 2、B 2、C 2 は全て正の閾値電圧分布であるとして説明したが、閾値電圧分布 A 2 は負の電圧の分布であり、閾値電圧分布 B 2、C 2 が正の電圧の分布であるような場合も、本発明の範囲に含まれる。また、閾値電圧分布 E R 1・E R 2 は正の値であったとしても、本発明はこれに限定されるものではない。また、本実施形態では E R 2、A 2、B 2、C 2 のデータの対応関係がそれぞれ “ 1 1 ”、“ 0 1 ”、“ 1 0 ”、“ 0 0 ” であるとしているが、たとえばそれぞれ “ 1 1 ”、“ 0 1 ”、“ 0 0 ”、“ 1 0 ” であるような他の対応関係であってもよい。

40

【 0 0 4 1 】

50

1つのメモリセルの2ビットデータは、下位ページデータと上位ページデータからなり、下位ページデータと上位ページデータは別々の書き込み動作、つまり、2回の書き込み動作により、メモリセルに書き込まれる。データを“*@"と標記するとき、*は上位ページデータを、@は下位ページデータを表している。

【0042】

まず、下位ページデータの書き込みを、図12の1段目～2段目を参照して説明する。全てのメモリセルは、消去状態の閾値電圧分布 E_R を有し、データ“11”を記憶しているものとする。図12に示すように、下位ページデータの書き込みを行うと、メモリセルの閾値電圧分布 E_R は、下位ページデータの値(“1”、或いは“0”)に応じて、2つの閾値電圧分布(E_{R1} 、 A_1)に分けられる。下位ページデータの値が“1”の場合には、消去状態の閾値電圧分布 E_R を維持するので $E_{R1} = E_R$ であるが、 $E_{R1} > E_R$ であってもよい。

10

【0043】

一方、下位ページデータの値が“0”の場合には、メモリセルのトンネル酸化膜に高電界を印加し、フローティングゲート電極に電子を注入して、メモリセルの閾値電圧 V_{th} を所定量だけ上昇させる。具体的には、ペリファイ電位 $VA1V$ を設定し、このペリファイ電圧 $VA1V$ 以上の閾値電圧となるまで書き込み動作が繰り返される。その結果、メモリセルは、書き込み状態(データ“10”)に変化する。書き込み動作を所定回繰り返しても閾値電圧に到達しなかった場合(または閾値電圧に到達しないメモリセル数が所定値以上の場合)、当該物理ページに対する書き込みは「書き込みエラー」(Program Error, Program Fail)となる。

20

【0044】

次に、上位ページデータの書き込みを、図12の2段目～3段目を参照して説明する。上位ページデータの書き込みは、チップの外部から入力される書き込みデータ(上位ページデータ)と、メモリセルに既に書き込まれている下位ページデータとに基づいて行われる。

【0045】

即ち、図12の2段目～3段目に示すように、上位ページデータの値が“1”の場合には、メモリセルのトンネル酸化膜に高電界がかからないようにし、メモリセルの閾値電圧 V_{th} の上昇を防止する。その結果、データ“11”(消去状態の閾値電圧分布 E_{R1})のメモリセルは、データ“11”をそのまま維持し(E_{R2})、データ“10”(閾値電圧分布 A_1)のメモリセルは、データ“10”をそのまま維持する(B_2)。ただし、各分布間の電圧マージンを確保するという点で、上述のペリファイ電圧 $VA1V$ よりも大きい正のペリファイ電圧 $VB2V$ を用いて閾値電圧分布の下限値を調整し、これにより閾値電圧分布の幅を狭めた閾値電圧分布 B_2 を形成するのが望ましい。下限値調整を所定回繰り返しても閾値電圧に到達しなかった場合(または閾値電圧に到達しないメモリセル数が所定値以上の場合)、当該物理ページに対する書き込みは「書き込みエラー」となる。

30

【0046】

一方、上位ページデータの値が“0”の場合には、メモリセルのトンネル酸化膜に高電界を印加し、フローティングゲート電極に電子を注入して、メモリセルの閾値電圧 V_{th} を所定量だけ上昇させる。具体的には、ペリファイ電位 $VA2V$ 、 $VC2V$ を設定し、このペリファイ電圧 $VA1V$ 以上の閾値電圧となるまで書き込み動作が繰り返される。その結果、データ“11”(消去状態の閾値電圧分布 E_{R1})のメモリセルは、閾値電圧分布 A_2 のデータ“01”に変化し、データ“10”(A1)のメモリセルは、閾値電圧分布 C_2 のデータ“00”に変化する。このとき、ペリファイ電圧 $VA2V$ 、 $VC2V$ が用いられて、閾値電圧分布 A_2 、 C_2 の下限値が調整される。書き込み動作を所定回繰り返しても閾値電圧に到達しなかった場合(または閾値電圧に到達しないメモリセル数が所定値以上の場合)、当該物理ページに対する書き込みは「書き込みエラー」となる。

40

【0047】

一方、消去動作においては、消去ペリファイ電位 V_{ev} を設定し、このペリファイ電圧 V

50

ev以下の閾値電圧となるまで消去動作が繰り返される。その結果、メモリセルは、書き込み状態（データ“00”）に変化する。消去動作を所定回繰り返しても閾値電圧に到達しなかった場合（または閾値電圧に到達しないメモリセル数が所定値以上の場合）、当該物理ページに対する消去は「消去エラー」となる。

【0048】

以上が、一般的な4値記憶方式におけるデータ書き込み方式の一例である。3ビット以上の多ビット記憶方式においても、上記の動作に更に上位のページデータに応じ、閾値電圧分布を8通り以上に分割する動作が加わるのみであるので、基本的な動作は同様である。

【0049】

（記憶装置の構成 その2）

つぎに、記憶部2および記憶部3の構成例について説明する。本実施形態では、図13に示すように、SSDとしての記憶部2および記憶部3は、不揮発性半導体メモリとしてのNAND型フラッシュメモリ（以下NANDメモリと略す）16と、インタフェース19を介して情報処理装置111との信号の送受信を行うインタフェースコントローラ（IFC）42と、IFC42とNANDメモリ16との中間バッファとして機能するキャッシュメモリ（CM）46を有する半導体メモリとしてのRAM（Random Access Memory）40と、NANDメモリ16及びRAM40の管理、制御、及びインタフェースコントローラ42の制御を司るSSDコントローラ（SSDC）41と、これら構成要素を接続するバス43を備える。

【0050】

RAM40としては、たとえば、DRAM（Dynamic Random Access Memory）、SRAM（Static Random Access Memory）などの揮発性RAMや、FeRAM（Ferroelectric Random Access Memory）、MRAM（Magnetoresistive Random Access Memory）、PRAM（Phase Change Random Access Memory）、ReRAM（Resistance Random Access Memory）などの不揮発性RAMを採用することができる。RAM40はSSDC41に含まれてもよい。

【0051】

NANDメモリ16は、複数のNANDメモリチップ80からなり、情報処理装置111によって指定されたユーザデータを記憶したり、ユーザデータを管理する管理テーブルを記憶したり、RAM40で管理される管理情報をバックアップ用に記憶したりする。NANDメモリ16は、複数のメモリセルがマトリクス状に配列されたメモリセルアレイ82を有し、個々のメモリセルは上位ページ及び下位ページを使用して多値記憶が可能である。NANDメモリ16は、複数のメモリチップによって構成され、各メモリチップは、データ消去の単位であるブロックを複数配列して構成される。また、NANDメモリ16では、ページごとにデータの書き込み及びデータの読み出しが行われる。ブロックは、複数のページによって構成されている。

【0052】

RAM40は、情報処理装置111とNANDメモリ16間でのデータ転送用キャッシュとして機能するキャッシュメモリ（CM）46を有する。また、RAM40は、管理情報記憶用メモリ及び作業領域用メモリとして機能する。RAM40の領域40Aで管理される管理テーブルは、NANDメモリ16の領域40Mに記憶されている各種管理テーブルが記憶部2および記憶部3の起動時などに展開されたものであり、定期的、スタンバイコマンド受信時、フラッシュコマンド受信時、あるいは電源断時にNANDメモリ16の領域40Mに退避保存される。

【0053】

SSDC41は、NANDメモリ16に記憶されたシステムプログラム（ファームウェア）を実行するプロセッサと、各種ハードウェア回路などによってその機能が実現され、情報処理装置111からのライト要求、キャッシュフラッシュ要求、リード要求等の各種コマンドに対する情報処理装置111 - NANDメモリ16間のデータ転送制御、RAM

10

20

30

40

50

40及びNANDメモリ16に記憶された各種管理テーブルの更新・管理、NANDメモリ16に書き込むデータのECC符号化、NANDメモリ16から読み出したデータのECC復号化などを実行する。

【0054】

情報処理装置111は記憶部2に対し、リード要求またはライト要求を発行する際には、インタフェース19を介して論理アドレスとしてのLBAを入力する。LBAは、論理セクタ(サイズ:例えば512Byte)に対して0からの通し番号をつけた論理アドレスである。また、情報処理装置111は記憶部2に対し、リード要求またはライト要求を発行する際には、LBAと併せて、リード要求またはライト要求の対象となる論理セクタサイズを入力する。

10

【0055】

IFC42は、情報処理装置111からのリード要求、ライト要求、その他要求及びデータを受信し、受信した要求やデータをSSDC41に送信したり、SSDC41の制御によりRAM40にデータを送信したりする機能を持つ。

【0056】

図14に記憶部2および記憶部3で使用する管理情報44の構成例を示す。これらの管理情報44は、前述したように、NANDメモリ16の領域40Mで不揮発記憶されている。領域40Mで記憶された管理情報が記憶部2の起動時にRAM40の領域40Aに展開されて使用される。領域40Aの管理情報44は、定期的あるいは電源断時に領域40Mに退避保存される。RAM40がMRAMやFeRAMなどのような不揮発RAMである場合には、この管理情報44はRAM40にのみ記憶されるようにしてもよく、この場合はこの管理情報44はNANDメモリ16には記憶されない。NANDメモリ16への書き込み量を少なくするためには、管理情報44に記憶されるデータは、RAM40の領域40Aに記憶されているデータを圧縮したものであることが望ましい。また、NANDメモリ16への書き込み頻度を少なくするためには、管理情報44には、RAM40の領域40Aに記憶されている管理情報44の更新情報(差分情報)を追記するようにすることが望ましい。

20

【0057】

図14に示すように、管理情報は、フリーブロックテーブル(FBT)60と、パッドブロックテーブル(BBT)61と、アクティブブロックテーブル(ABT)62と、トラックテーブル(トラック単位の論物変換テーブル)63と、クラスタテーブル(クラスタ単位の論物変換テーブル)64、統計情報65を含む。

30

【0058】

LBAは、図15に示すように、論理セクタ(サイズ:例えば512Byte)に対して0からの通し番号をつけた論理アドレスである。本実施形態においては、記憶部2の論理アドレス(LBA)の管理単位として、LBAの下位($s+1$)ビット目から上位のビット列で構成されるクラスタアドレスと、LBAの下位($s+t+1$)ビットから上位のビット列で構成されるトラックアドレスとを定義する。すなわち、論理セクタは、情報処理装置111からの最小アクセス単位である。クラスタは、SSD内部で「小さなデータ」を管理する管理単位であり、論理セクタサイズの自然数倍がクラスタサイズとなるように定められる。また、トラックは、SSD内部で「大きなデータ」を管理する管理単位であり、クラスタサイズの2以上の自然数倍がトラックサイズとなるように定められる。したがって、トラックアドレスはLBAをトラックサイズで割ったものであり、トラック内アドレスはLBAをトラックサイズで割った余りであり、クラスタアドレスはLBAをクラスタサイズで割ったものであり、クラスタ内アドレスはLBAをクラスタサイズで割った余りである。以下の説明では、便宜上、トラックのサイズは1物理ブロックに記録可能なデータのサイズに等しい(物理ブロックにSSDC41で行うECC処理の冗長ビットが含まれる場合はこれを除いたサイズ)場合について述べ、クラスタのサイズは1物理ページに記録可能なデータのサイズに等しい(物理ページにSSDC41で行うECC処理の冗長ビットが含まれる場合はこれを除いたサイズ)場合について述べる。

40

50

【 0 0 5 9 】

フリーブロックテーブル (F B T) 6 0 は、N A N D メモリ 1 6 に書き込みを行うときに書き込み用に新規に割り当てることができる N A N D メモリの用途未割り当ての物理ブロック (フリーブロック : F B) のブロックアドレス (物理ブロック I D) を管理する。また、物理ブロック I D 毎に消去回数を管理しており、物理ブロックが消去されたとき、当該ブロックの消去回数をインクリメントする。

【 0 0 6 0 】

バッドブロックテーブル (B B T) 6 1 は、誤りが多いなど記憶領域として使用できない物理ブロックとしてのバッドブロック (B B) のブロック I D を管理する。F B T 6 0 と同様に物理ブロック I D 毎に消去回数を管理するようにしてもよい。

10

【 0 0 6 1 】

アクティブブロックテーブル (A B T) 6 2 は、用途が割り当てられた物理ブロックであるアクティブブロック (A B) を管理する。また、物理ブロック I D 毎に消去回数を管理しており、物理ブロックが消去されたとき、当該ブロックの消去回数をインクリメントする。

【 0 0 6 2 】

トラックテーブル 6 3 は、トラックアドレスと、このトラックアドレスに対応するトラックデータが記憶される物理ブロック I D との対応関係を管理する。

【 0 0 6 3 】

クラスタテーブル 6 4 は、クラスタアドレスと、このクラスタアドレスに対応するクラスタデータが記憶される物理ブロック I D と、このクラスタアドレスに対応するクラスタデータが記憶される物理ブロック内ページアドレスとの対応関係を管理する。

20

【 0 0 6 4 】

S S D C 4 1 は、統計情報 6 5 に、記憶部 2 の信頼性に関わる種々のパラメータ (X 0 1 ~ X 3 2) を信頼性情報として格納する (図 2 7 参照) 。

【 0 0 6 5 】

信頼性情報の例として用いられる統計情報 6 5 の値 (Raw Value) としては、バッドブロック数総計 (統計情報 X 0 1) 、バッド論理セクタ数総計 (統計情報 X 0 2) 、消去回数総計 (統計情報 X 0 3) 、消去回数平均値 (統計情報 X 0 4) 、N A N D メモリの書き込みエラー発生回数累積値 (統計情報 X 0 5) 、N A N D メモリの消去エラー発生回数累積値 (統計情報 X 0 6) 、読み出し論理セクタ数総計 (統計情報 X 0 7) 、書き込み論理セクタ数総計 (統計情報 X 0 8) 、E C C 訂正不能回数総計 (統計情報 X 0 9) 、リトライ読み出し発生回数 (統計情報 X 1 0) 、n ビット ~ m ビット E C C 訂正単位総計数 (統計情報 X 1 1) 、インタフェース 1 9 のデータ化けエラー発生回数 (統計情報 X 1 2) 、インタフェース 1 9 の通信速度ダウンシフト回数 (統計情報 X 1 3) 、インタフェース 1 9 のレーン数ダウンシフト回数 (統計情報 X 1 4) 、インタフェース 1 9 のエラー発生回数 (統計情報 X 1 5) 、R A M 4 0 のエラー発生回数 (統計情報 X 1 6) 、記憶部 2 の使用時間総計 (統計情報 X 1 7) 、起動回数 (統計情報 X 1 8) 、不正電源断発生回数 (統計情報 X 1 9) 、温度が推奨動作温度の最高値を上回った時間累計 (統計情報 X 2 0) 、温度が推奨動作温度の最低値を下回った時間累計 (統計情報 X 2 1) 、コマンドの応答時間最大値 (統計情報 X 2 2) 、コマンドの応答時間平均値 (統計情報 X 2 3) 、N A N D メモリの応答時間最大値 (統計情報 X 2 4) 、N A N D の応答時間平均値 (統計情報 X 2 5) 、現在温度 (統計情報 X 2 6) 、最高温度 (統計情報 X 2 7) 、最低温度 (統計情報 X 2 8) 、管理情報冗長度 (統計情報 X 2 9) 、R A M 4 0 への書き込みデータ量合計 (統計情報 X 3 0) 、統計情報増加率 (統計情報 X 3 1) 、N A N D 整理失敗フラグ (統計情報 X 3 2) などが含まれる。

30

40

【 0 0 6 6 】

バッドブロック数総計 (統計情報 X 0 1) について説明する。記憶部 2 内の N A N D メモリ 1 6 の物理ブロックが一つバッドブロックに追加されるごとに統計情報 X 0 1 が 1 インクリメントされる。統計情報 X 0 1 は記憶部 2 の製造時 (検査工程前) にゼロにリセッ

50

トされていることが望ましく、検査工程でエラーが発生したり、閾値分布の分布間マージンが少ないことが判明したブロックは、あらかじめバッドブロックに加えておくことがさらに望ましい。SSDC41は統計情報X01を統計情報65に格納せず、統計情報X01をBBT61から直接計算してもよい。統計情報X01が大きいほど信頼性が悪化していることを示す。

【0067】

バッド論理セクタ数総計（統計情報X02）について説明する。SSDC41は情報処理装置111から読み出し命令とLBAを受信してNAND型フラッシュメモリ16を読みだした時に読み出しデータをECC訂正できなかった場合、当該LBAをバッド論理セクタとして管理情報44内のバッド論理セクタテーブルに登録しても良い（図25参照）。SSDC41は、バッド論理セクタテーブルに登録されているLBAの数をバッド論理セクタ数総計（統計情報X02）として統計情報65に格納する。情報処理装置111から読み出し命令を受信した場合には、SSDC41はRAM40上のバッド論理セクタテーブルを読みだして受信したLBAをバッド論理セクタテーブルから検索し、バッド論理セクタテーブルから見つかった場合には、NAND型フラッシュメモリ16を読み出すことなく情報処理装置111に読み出しエラーを通知する。SSDC41はバッド論理セクタのLBAに対して情報処理装置111から書き込み命令を受信して書き込み処理を行った場合、SSDC41は書き込まれたLBAをバッド論理セクタテーブルから削除する。SSDC41はバッド論理セクタのLBAに対して情報処理装置111から削除通知を受信して削除通知処理を行った場合、SSDC41は削除通知処理されたLBAをバッド論理セクタテーブルから削除する。SSDC41は情報処理装置111から記憶部2の消去コマンド（Secure Eraseコマンド）を受信した場合にはバッド論理セクタテーブルを消去する。記憶部2の消去コマンドとして、たとえばACS-3のF4h Security Erase Unitコマンドや、NVM Express Revision 1.1の80h Format NVMコマンドを用いてもよい。なお、LBA単位（論理セクタ単位）でバッド論理セクタテーブルを管理する代わりに、図26に示すように、クラスタ単位でバッド論理セクタテーブルをバッドクラスタテーブルとして管理しても良い。SSDC41はバッド論理セクタテーブルに登録されているLBAの個数またはバッドクラスタテーブルに登録されているクラスタアドレスの個数を統計情報X02として管理する。SSDC41は統計情報X02を統計情報65に格納せず、統計情報X02をバッド論理セクタテーブルやバッドクラスタテーブルから直接計算してもよい。統計情報X02が大きいほど信頼性が悪化していることを示す。

【0068】

消去回数総計（統計情報X03）について説明する。統計情報X03は記憶部2内のNANDメモリ16の全ブロックの消去回数累計値を示す。SSDC41は、記憶部2内のNANDメモリ16の物理ブロックが一つ消去されるごとに統計情報X03を1インクリメントする。統計情報X03は記憶部2の製造時（検査工程時）にゼロにリセットされていることが望ましい。SSDC41は統計情報X03を統計情報65に格納せずFBT60、BBT61、ABT62から直接計算してもよい。統計情報X03が大きいほど信頼性が悪化していることを示す。

【0069】

消去回数平均値（統計情報X04）について説明する。SSDC41は、記憶部2内のNANDメモリ16の全ブロックについて1ブロックあたりの消去回数平均値を計算し、統計情報X04として統計情報65に格納する。SSDC41は、管理情報44を格納するブロックなど一部のブロックを統計情報X04の集計対象から除外してもよい。統計情報X04は記憶部2の製造時（検査工程前）にゼロにリセットされていることが望ましい。SSDC41は統計情報X04を統計情報65に格納せずFBT60、BBT61、ABT62から直接計算してもよい。SSDC41は、統計情報X04として、消去回数平均値のかわりに、消去回数最大値や消去回数最小値を使用してもよい。統計情報X04が大きいほど信頼性が悪化していることを示す。

【0070】

NANDメモリの書き込みエラー発生回数累積値（統計情報X05）について説明する。SSDC41は、記憶部2内のNANDメモリ16で書き込みエラーが1書き込み単位で発生するごとに、統計情報X05を1加算する（ブロック単位で加算してもよい）。統計情報X05は記憶部2の製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報X05が大きいほど信頼性が悪化していることを示す。

【0071】

NANDメモリの消去エラー発生回数累積値（統計情報X06）について説明する。統計情報X06は記憶部2製造時（検査工程前）にゼロにリセットされていることが望ましい。SSDC41は、記憶部2内のNANDメモリ16で消去エラーが1ブロックで発生するごとに統計情報X06を1加算する。SSDC41は、複数のブロックをまとめて消去単位とし、この消去単位1つで消去エラーが発生するごとに統計情報X06に1加算するようにしてもよい。統計情報X06が大きいほど信頼性が悪化していることを示す。

10

【0072】

読み出し論理セクタ数総計（統計情報X07）について説明する。SSDC41は、IFC42が読み出しデータとして情報処理装置111に送信したデータの論理セクタ数合計を、統計情報X07として統計情報65に格納する。統計情報X07は記憶部2の製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報X07が大きいほど信頼性が悪化していることを示す。

【0073】

書き込み論理セクタ数総計（統計情報X08）について説明する。SSDC41は、IFC42が書き込みデータとして情報処理装置111から受信したデータの論理セクタ数合計を、統計情報X08として統計情報65に格納する。統計情報X08は記憶部2の製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報X08が大きいほど信頼性が悪化していることを示す。

20

【0074】

ECC訂正不能回数総計（統計情報X09）について説明する。SSDC41は、ECC訂正によりエラービットが修復できなかった場合に、1読み出し単位ごとに統計情報X09を1インクリメントする。SSDC41は、エラー訂正できなかったエラービット数の推定値を加算するようにしてもよいし、エラー訂正できなかったブロックの数を加算するようにしてもよい。統計情報X09は記憶部2製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報X09が大きいほど信頼性が悪化していることを示す。

30

【0075】

リトライ読み出し発生回数総計（統計情報X10）について説明する。SSDC41は、データ読み出し時にエラービット数が多く誤り訂正が不可能（ECCエラー）であった場合、リトライ読み出しを行い、再度ECCを用いた誤り訂正を実行することが望ましい。特に、SSDC41は、リトライ読み出し時に図12の読み出しレベルVA1、VA2、VB2やVC2をデフォルトの値からシフトして読み出しを行うことで、誤り訂正不可能であったデータが誤り訂正可能になる場合がある。SSDC41はリトライ読み出し発生回数を統計情報X10として統計情報X09に格納し、寿命到達予測や寿命到達判定に用いてもよい。統計情報X10は記憶部2の製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報X10が大きいほど信頼性が悪化していることを示す。

40

【0076】

nビット～mビットECC訂正単位総計数（統計情報X11）について説明する。n, mは自然数で、0 < n < m 訂正可能最大ビット数である。SSDC41は、ECC訂正単位（たとえば物理ページ）に対してECC訂正を行った時、全エラービットが正常に修復され、かつ修復されたエラービット数がn以上m以下である場合、ECC訂正単位1つにつき「nビット～mビットECC訂正単位総計数」を1加算する。ECC訂正により1訂正単位につき最大64ビット訂正可能である場合、たとえば、SSDC41は、「1ビット～8ビットECC訂正単位総計数」「9ビット～16ビットECC訂正単位総計数」

50

「17ビット～24ビットECC訂正単位総計数」「25ビット～32ビットECC訂正単位総計数」「33ビット～40ビットECC訂正単位総計数」「41ビット～48ビットECC訂正単位総計数」「49ビット～56ビットECC訂正単位総計数」「57ビット～64ビットECC訂正単位総計数」の8つのパラメータを用意し、ECC訂正が正常に行われた場合、1ECC訂正単位のECC訂正につきこれら8つのパラメータのうちいずれか1つに1をインクリメントする。統計情報X11は記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。統計情報X11が大きいほど信頼性が悪化していることを示す。

【0077】

インタフェース19のデータ化けエラー発生回数(統計情報X12)について説明する。SSDC41は、インタフェース19上で信号のデータ化けを検出するごとに統計情報X12を1インクリメントする。インタフェース19上で送受信されるデータは、たとえばCyclic Redundancy Check(CRC)符号やBose-Chaudhuri-Hocquenghem(BCH)符号やReed-Solomon(RS)符号やLow-Density Parity-Check(LDPC)符号などを用いてSSDC41やIFC42やチップセット7によって誤り検出や誤り訂正が行われており、誤りが検出された場合や誤り訂正できなかった場合にSSDC41は統計情報X12を1インクリメントする。たとえば、インタフェース19がSATA規格である場合、SATA規格におけるRエラー(Reception Error, R_ERR)が1回発生するごとにSSDC41は統計情報X12を1インクリメントする。統計情報X12として、SATA規格のPhy Event Countersのカウンタのいずれかを採用してもよい。統計情報X12は記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。統計情報X12が大きいほど信頼性が悪化していることを示す。

【0078】

インタフェース19の通信速度ダウンシフト回数(統計情報X13)について説明する。SSDC41やIFC42やチップセット7がインタフェース19の通信速度が設計値よりも小さくなったことを検出した場合、SSDC41はX13を1インクリメントする。例えば、通信インタフェース19やIFC42やSSDC41がSATAの通信速度が最大6Gbpsで設計されているにも関わらず、実際に記憶部2や情報処理装置111との間で確立された通信速度が3Gbpsのようなより低速の通信速度であることが検出された場合、SSDC41はこれをSATA通信上のエラーとみなし、統計情報X13を1インクリメントする。たとえば、インタフェース19やIFC42やSSDC41がPCI Expressの通信速度が最大8Gbpsとして設計されているにも関わらず、実際に記憶部2や情報処理装置111との間で確立された通信速度が5Gbpsのようなより低速の通信規格であることが検出された場合、SSDC41はこれをPCI Express通信上のエラーとみなし、統計情報X13を1インクリメントする。統計情報X13は記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【0079】

インタフェース19のレーン数ダウンシフト回数(統計情報X14)について説明する。SSDC41やIFC42やチップセット7がインタフェース19のアクティブな伝送路数が設計値よりも小さくなったことを検出した場合、SSDC41はX14を1インクリメントする。たとえば、インタフェース19やIFC42やSSDC41がPCI Expressの伝送路数(Lane数)が最大8Laneとして設計されているにも関わらず、実際に記憶部2や情報処理装置111との間で確立された伝送路数が4Laneのようなより少数の伝送路数であることが検出された場合、SSDC41はこれをPCI Express通信上のエラーとみなし、統計情報X14を1インクリメントする。統計情報X14は記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【0080】

インタフェース19のエラー発生回数(統計情報X15)について説明する。SSDC

10

20

30

40

50

4 1 や I F C 4 2 やチップセット 7 がその他の (X 1 2 以外の) インタフェース 1 9 での異常を 1 回検出するごとに、S S D C 4 1 は統計情報 X 1 5 を 1 インクリメントする。統計情報 X 1 5 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【 0 0 8 1 】

R A M 4 0 のエラー発生回数 (統計情報 X 1 6) について説明する。S S D C 4 1 がデータを R A M 4 0 に書き込む場合、S S D C 4 1 または R A M 4 0 の E C C 符号化部または誤り検出符号作成部がデータを符号化して R A M 4 0 に書き込み、S S D C 4 1 がデータを R A M 4 0 から読み出す場合、S S D C 4 1 または R A M 4 0 の E C C 復号部または誤り検出部がデータを誤り訂正または誤り検出して R A M 4 0 からデータを読み出す。S S D C 4 1 が R A M 4 0 からデータを読み出す場合に誤り訂正できなかった時または誤り検出した時、S S D C 4 1 は統計情報 X 1 6 を 1 インクリメントする。統計情報 X 1 6 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【 0 0 8 2 】

記憶部 2 の使用時間総計 (統計情報 X 1 7) について説明する。記憶部 2 の電源が O N になっている間、S S D C 4 1 がクロックをカウントしたり、内部の時計回路から時刻情報を受信することで、S S D C 4 1 は経過時間として統計情報 X 1 7 をインクリメントする。あるいは、S S D C 4 1 が情報処理装置 1 1 1 から定期的に情報処理装置 1 1 1 の時刻情報を受信するようにし、その時刻情報の差分をインクリメントするようにしてもよい。統計情報 X 1 7 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【 0 0 8 3 】

起動回数 (統計情報 X 1 8) について説明する。記憶部 2 に電源が供給され起動されるたびに S S D C 4 1 は X 1 8 を 1 インクリメントする。電源起動時には N A N D 型フラッシュメモリ 1 6 に対して読み出し動作が発生し、書き込み動作が発生することがあるため、この値が大きいほど信頼性が悪化していることを示す。統計情報 X 1 8 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。

【 0 0 8 4 】

不正電源断発生回数 (統計情報 X 1 9) について説明する。通常、記憶部 2 の電源をオフにするときには、情報処理装置 1 1 1 は記憶部 2 に対して、たとえば Information technology ATA/ATAPI Command Set-3 (ACS-3) に記載されている E0h Standby Immediate コマンド等を発行したり、NVM Express Revision1.1 に記載されている Shutdown Notification (CC.SHN) を 01b にセットしたりするなどして、あらかじめ記憶部 2 を電源遮断可能な状態に遷移させておいてから記憶部 2 の電源を遮断する。一方、記憶部 2 が電源遮断可能ではない状態において意図せず電源遮断が発生する場合があります、これは不正電源断 (不正電源遮断、Ungraceful Power Down, Unsafe Shutdown, Unintentional Power Down) とよばれる。不正電源遮断後にはじめて記憶部 2 が起動する時、S S D C 4 1 は X 1 9 を 1 インクリメントする。不正電源遮断においてはユーザデータが破壊されたり、N A N D メモリ 1 6 に対して大量の読み出し・書き込み動作が発生して信頼性劣化要因にもなるため、X 1 9 が大きいほど信頼性が悪化していることを示す。統計情報 X 1 9 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。

【 0 0 8 5 】

推奨動作温度の最高値を上回った時間累計 (統計情報 X 2 0) について説明する。たとえば、記憶部 2 の基板上、S S D C 4 1 内、N A N D メモリ 1 6 内など、記憶部 2 内に温度計が実装されている場合、S S D C 4 1 は温度計から定期的に温度情報を受信している。受信した温度が推奨動作温度 (たとえば 1 0 0) を上回った場合、S S D C 4 1 はクロックや内部の時計や情報処理装置 1 1 1 から取得する時刻情報をもとに、推定動作温度以上で動作している時間数をインクリメントしていく。統計情報 X 2 0 は記憶部 2 の製造時 (検査工程前) にゼロにリセットされていることが望ましい。この値が大きいほど信頼

性が悪化していることを示す。

【 0 0 8 6 】

推奨動作温度の最低値を下回った時間累計（統計情報 X 2 1）について説明する。記憶部 2 内に温度計が実装されている場合、S S D C 4 1 は温度計から定期的に温度情報を受信している。受信した温度が推奨動作温度（たとえば - 4 0 ）を下回った場合、S S D C 4 1 はクロックや内部の時計や情報処理装置 1 1 1 から取得する時刻情報をもとに、推定動作温度以上で動作している時間数をインクリメントしていく。記憶部 2 の製造時（検査工程前）にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【 0 0 8 7 】

コマンドの応答時間最大値（統計情報 X 2 2）について説明する。S S D C 4 1 は、情報処理装置 1 1 1 からコマンドを受信してから、情報処理装置 1 1 1 へ応答するまで（またはコマンド実行完了するまで）に要した時間（またはクロック数）を計測し、その最大値を統計情報 X 2 2 として統計情報 6 5 に格納する。X 2 2 を上回る応答時間が発生した場合は、S S D C 4 1 はこの応答時間を X 2 2 に上書きする。S S D C 4 1 は、コマンドそれぞれに対して統計情報 X 2 2 を格納してもよい。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 2 がゼロにリセットされていることが望ましい。

【 0 0 8 8 】

コマンドの応答時間平均値（統計情報 X 2 3）について説明する。S S D C 4 1 は、情報処理装置 1 1 1 からコマンドを受信してから、情報処理装置 1 1 1 へ応答するまで（またはコマンド実行完了するまで）に要した時間（またはクロック数）を計測し、その平均値を統計情報 X 2 3 として統計情報 6 5 に格納する。たとえば、S S D C 4 1 は応答時間リストを一定数 R A M 4 0 に保持しておき、その応答時間リストの平均値を算出することにより統計情報 X 2 3 を計算する。S S D C 4 1 はコマンドそれぞれに対して統計情報 X 2 3 を格納してもよい。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 3 がゼロにリセットされていることが望ましい。

【 0 0 8 9 】

N A N D メモリの応答時間最大値（統計情報 X 2 4）について説明する。S S D C 4 1 は、S S D C 4 1 が N A N D メモリ 1 6 に命令してから応答を得る（またはコマンド実行完了通知を受信する）までに要した時間（またはクロック数）を計算し、その最大値を統計情報 X 2 4 として統計情報 6 5 に格納する。X 2 4 を上回る応答時間が発生した場合は、S S D C 4 1 はこの応答時間を X 2 4 に上書きする。S S D C 4 1 はコマンドそれぞれに対して統計情報 X 2 4 を保持してもよい。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 4 がゼロにリセットされていることが望ましい。

【 0 0 9 0 】

N A N D の応答時間平均値（統計情報 X 2 5）について説明する。S S D C 4 1 は、N A N D メモリ 1 6 に命令してから応答を得る（またはコマンド実行完了通知を受信する）までに要した時間（またはクロック数）を計測し、その平均値を統計情報 X 2 5 として統計情報 6 5 に格納する。たとえば S S D C 4 1 は応答時間リストを一定数 R A M 4 0 に格納しておき、その応答時間リストの平均値を算出することにより統計情報 X 2 5 を得る。S S D C 4 1 はコマンドそれぞれに対して統計情報 X 2 5 を保持してもよい。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 5 がゼロにリセットされていることが望ましい。

【 0 0 9 1 】

現在温度（統計情報 X 2 6）について説明する。記憶部 2 内に温度計が実装されている場合、S S D C 4 1 は温度計から定期的に温度情報を受信する。S S D C 4 1 は温度計から最後に受信した温度を現在温度として統計情報 X 2 6 に格納する。S S D C 4 1 は、この値が極端に大きいと（たとえば 8 5 以上）、記憶部 2 の信頼性に悪影響があり、また、この温度が極端に小さいと（たとえば - 1 0 以下）、記憶部 2 の信頼性に悪影響があると判定する。

10

20

30

40

50

【 0 0 9 2 】

最高温度（統計情報 X 2 7）について説明する。S S D C 4 1 は、現在温度 X 2 6 の最大値を最高温度として統計情報 X 2 7 に格納する。この値が極端に大きいと（たとえば 8 5 以上）、記憶部 2 の信頼性に悪影響がある。S S D C 4 1 は、X 2 7 よりも大きい現在温度を温度計から受信した時、X 2 7 を現在温度に書き換える。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 7 が記憶部 2 の動作温度にくらべて十分に小さい温度（たとえば - 4 0 ）にリセットされていることが望ましい。

【 0 0 9 3 】

最低温度（統計情報 X 2 8）について説明する。S S D C 4 1 は、現在温度 X 2 6 の最小値を最低温度として統計情報 X 2 8 に格納する。この値が極端に小さいと（たとえば - 4 0 以下）、記憶部 2 の信頼性に悪影響がある。S S D C 4 1 は、X 2 8 よりも小さい現在温度を温度計から受信した時、X 2 8 を現在温度に書き換える。記憶部 2 の製造時（検査工程前）または記憶部 2 の出荷時には X 2 8 が記憶部 2 の動作温度にくらべて十分に大きい温度（たとえば 1 2 0 ）にリセットされていることが望ましい。

【 0 0 9 4 】

管理情報冗長度（統計情報 X 2 9）について説明する。N A N D メモリ 1 6 の管理情報領域 4 0 M のデータや N A N D メモリ 1 6 に記憶されたシステムプログラム（ファームウェア）などのシステムデータが壊れて読みだし不能になると、記憶部 2 が正常な動作ができなくなる可能性がある。記憶部 2 の信頼性向上のため、S S D C 4 1 はシステムデータを R A I D 1 や R A I D 5 や R A I D 6 により複数物理ブロックや複数チャネルをまたいで冗長化して領域 4 0 M に格納することが望ましい。S S D C 4 1 は、システムデータの冗長度を数値化して管理情報冗長度（統計情報 X 2 9）として統計情報 6 5 に格納する。冗長度 X 2 9 = R の時、最大（R - 1）ブロックのデータ消失まではデータを復元することができる。たとえば、S S D C 4 1 が管理情報 4 5 を 4 ブロックにわたって R A I D 1 で管理する場合、管理情報 4 5 はたとえばブロック A とブロック B とブロック C とブロック D それぞれにクローンとして格納される。この場合、管理情報 4 5 は合計 4 つのクローンを保持するため、管理情報 4 5 の冗長度 X 2 9 は 4 である。たとえばブロック A のデータが壊れて読みだし不能になった場合、S S D C 4 1 はブロック B またはブロック C またはブロック D からデータ読み出しを行うことで管理情報 4 5 を読み出すことができる。この場合、管理情報 4 5 は合計 3 つのクローンを保持するため、管理情報 4 5 の冗長度 X 2 9 は 3 である。たとえば、S S D C 4 1 が管理情報 4 5 を 4 ブロックにわたって R A I D 5 で管理する場合、管理情報 4 5 はたとえばブロック A とブロック B とブロック C とブロック D それぞれにおいて 4 台の R A I D 5 で格納される。この場合、最大 1 ブロックのデータが消失したとしてもデータを復元可能であるため、管理情報 4 5 の冗長度 X 2 9 は 2 である。1 ブロック分のデータが消失している状態においては、冗長度 X 2 9 は 1 となる。冗長度 X 2 9 が下がると、システムデータが復元不可能になる確率が上がり記憶部 2 の故障率が増大するため、冗長度 X 2 9 が小さいほど信頼性が悪化していることを示す。冗長度 X 2 9 が減少している場合、S S D C 4 1 は復元したデータをデータ消失したブロックに再書き込みすることで冗長度を回復させることが望ましい。

【 0 0 9 5 】

R A M 4 0 への書き込みデータ量合計（統計情報 X 3 0）について説明する。S S D C 4 1 は、記憶部 2 内の R A M 4 0 へのデータ書き込み量の累計値を統計情報 X 3 0 として統計情報 6 5 に格納する。S S D C 4 1 は、たとえば R A M 4 0 に 1 ページのデータを書き込むごとに統計情報 X 3 0 を 1 インクリメントする。統計情報 X 3 0 は記憶部 2 の製造時（検査工程前）にゼロにリセットされていることが望ましい。統計情報 X 3 0 が大きいほど信頼性が悪化していることを示す。

【 0 0 9 6 】

統計情報増加率（統計情報 X 3 1）について説明する。S S D C 4 1 は、統計情報 X 0 1 ~ X 2 5 の最新でない情報（たとえば一定時刻前や、記憶部 2 をパワーオンした時の値や前回記憶部 2 をパワーダウンしたときの値など）を別途管理情報 4 4 に格納しておく。

SSDC41は、例えば、下記のいずれかの式で統計情報X31を計算する。

統計情報増加率 = (最新統計情報) (旧情報)

統計情報増加率 = ((最新統計情報) (旧情報)) / (旧情報を取得してからの経過時刻)

統計情報増加率 = ((最新統計情報) (旧情報)) / (旧情報を取得してからのNANDアクセス回数)

記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【0097】

NAND整理失敗フラグ(統計情報X32)について説明する。統計情報X32が1になっていると、NAND整理によっても動作に十分な数のフリーブロック数を確保できないことになる。記憶部2の製造時(検査工程前)にゼロにリセットされていることが望ましい。この値が大きいほど信頼性が悪化していることを示す。

【0098】

SSDC41は、統計情報65に上述したすべてのパラメータを格納してもよいし、これらの一部あるいはどれか一つのみを格納してもよい。SSDC41は、統計情報65を最新情報をRAM40上の領域40Aに保持し、定期的にNANDメモリ16上の領域40Mにバックアップデータとして退避保存することが望ましい。一方、SSDC41は、RAM40やNANDメモリ16のどちらか一方にのみ保存するようにしてもよいし、当該統計情報を情報処理装置111に送信して、情報処理装置111や情報処理装置111に接続された記憶装置に保存するようにしてもよい。

【0099】

(LBA正引き変換)

つぎに、図16を用いてSSDC41がLBAから物理アドレスを特定する手順(LBA正引き変換)について説明する。LBAが指定されたとき、SSDC41はLBAからトラックアドレスとクラスタアドレスとクラスタ内アドレスを計算する。

【0100】

SSDC41は、まずトラックテーブル63を検索し、計算したトラックアドレスに対応する物理ブロックIDを特定する(ステップS100、S101)。SSDC41は、特定した物理ブロックIDが有効であるか否かを判定し(ステップS102)、物理ブロックIDがヌルではなく有効な値である場合は(ステップS102: Yes)、この物理ブロックIDがABT62にエントリされているか否かを検索する(ステップS103)。ABT62に物理ブロックIDがエントリされている場合は(ステップS104: Yes)、この物理ブロックIDの指定する物理ブロックの先頭位置からトラック内アドレス分だけシフトした位置が指定されたLBAに対応するNANDメモリ16上の物理的な位置となる(ステップS105)。このような場合には、LBAに対応するNANDメモリ16上の物理的な位置の特定にはクラスタテーブル64を必要とせず、このようなLBAを「トラック管理されているLBA」と呼ぶ。ステップS104において、ABT62に物理ブロックIDがエントリされていない場合は(ステップS104: No)、指定されたLBAは対応する物理アドレスを持たないことになり、このような状態を「未書き込み状態」と呼ぶ(ステップS106)。

【0101】

ステップS102において、指定されたトラックアドレスに対応する物理アドレスがヌルであり無効な値の場合は(ステップS102: No)、SSDC41はLBAからクラスタアドレスを計算し、クラスタテーブル64を検索し、計算したクラスタアドレスに対応する物理ブロックID及び対応する物理ブロック内ページアドレスをクラスタテーブル64から取得する(ステップS107)。物理ブロックIDと物理ブロック内ページアドレスが指定する物理ページの先頭位置からクラスタ内アドレス分だけシフトした位置が指定されたLBAに対応するNANDメモリ16上の物理的な位置となる。このような場合は、LBAに対応するNANDメモリ16上の物理的な位置は、トラックテーブル63の

みからは特定できず、クラスタテーブル 6 4 の参照を必要とすることになり、このような L B A を「クラスタ管理されている L B A」という(ステップ S 1 0 8)。

【 0 1 0 2 】

(読み出し動作)

つぎに、図 1 7、図 1 8 を用いて情報処理装置 1 1 1 による記憶部 2 および記憶部 3 からの読み出し動作を説明する。本実施形態で説明する読み出し動作は、読み出しコマンドとして Information technology ATA/ATAPI Command Set-3 (ACS-3) に記載されている 60h READ FPDMA QUEUED を用いた場合であるが、25h READ DMA EXT などその他書き込みコマンドを採用してもよく、読み出しコマンドの種類の違いは発明の本質には影響しない。たとえば、読み出しコマンドとして NVM Express Revision 1.1 に記載されている 02h Read を用いても良い。記憶部 2 が情報処理装置 1 1 1 から読み出し命令を受信した場合は(ステップ S 1 1 0)、S S D C 4 1 がこの読み出し命令を R A M 4 0 上の読み出し命令待ち行列に追加し(ステップ S 1 1 1)、読み出し命令を受信した旨を情報処理装置 1 1 1 に返信する。

10

【 0 1 0 3 】

一方、S S D C 4 1 は、R A M 4 0 上の読み出し命令待ち行列に命令が存在している場合、読み出し処理を実行可能な状態になっているか否かを判定し(ステップ S 1 2 0)、読み出し処理を実行可能な状態になったと判定すると、先の図 1 6 に示した L B A 正引き変換手順にしたがって情報処理装置 1 1 1 から受信した L B A からデータの物理的な位置を特定する(ステップ S 1 2 1)。S S D C 4 1 は、特定した位置の物理ページからデータを読み出し(ステップ S 1 2 3)、読み出したデータのうち E C C 冗長ビットを用いて E C C 復号化し(ステップ S 1 2 4)、復号化したデータを I F C 4 2 を介して情報処理装置 1 1 1 に送信し(ステップ S 1 2 5)、統計情報 6 5 を更新する。なお、N A N D メモリ 1 6 から読み出したデータは、いったん R A M 4 0 に書き込み、R A M 4 0 に書き込んだデータを復号化して情報処理装置 1 1 1 に送信するようにしてもよいし、復号化したデータをいったん R A M 4 0 に書き込み、R A M 4 0 に書き込んだデータを情報処理装置 1 1 1 に送信するようにしてもよい。

20

【 0 1 0 4 】

ステップ S 1 2 4 において、S S D C 4 1 は E C C による復号化を試みるが、復号化できなかった場合、復号化できなかったページを含む物理ブロックを A B T 6 2 から削除して B B T 6 1 に登録し、統計情報 6 5 の E C C 訂正できなかった E C C 訂正単位数総計(統計情報 X 0 9)を加算する。その際、当該ブロックのデータを F B T 6 0 から割り当てたフリーブロックにコピーし、当該フリーブロックの物理ブロック I D を A B T 6 2 に登録してトラックテーブル 6 3 およびクラスタテーブル 6 4 の物理ブロックをコピー元物理ブロック I D からコピー先物理ブロック I D に書き換えることが望ましい。

30

【 0 1 0 5 】

(書き込み動作)

つぎに、図 1 9、図 2 0 を用いて情報処理装置 1 1 1 による記憶部 2 および記憶部 3 への書き込み動作を説明する。本実施形態で説明する書き込み動作は、書き込みコマンドとして Information technology ATA/ATAPI Command Set-3 (ACS-3) に記載されている 61h WRITE FPDMA QUEUED を用いた場合であるが、35h WRITE DMA EXT などその他書き込みコマンドを採用してもよく、書き込みコマンドの種類の違いは発明の本質には影響しない。たとえば、書き込みコマンドとして NVM Express Revision 1.1 に記載されている 01h Write を用いても良い。たとえば、記憶部 2 が情報処理装置 1 1 1 から書き込み命令を受信した場合は(ステップ S 1 3 0)、S S D C 4 1 がこの書き込み命令を R A M 4 0 上の読み出し命令待ち行列に追加し(ステップ S 1 3 1)、書き込み命令を受信した旨を情報処理装置 1 1 1 に返信する。

40

【 0 1 0 6 】

一方、S S D C 4 1 は、R A M 4 0 上の書き込み命令待ち行列に命令が存在している場合、書き込み処理を実行可能な状態になっているか否かを判定し(ステップ S 1 4 0)、

50

書き込み処理を実行可能な状態になったと判定すると、情報処理装置 1 1 1 に書き込み可能であることを通知し、情報処理装置 1 1 1 から書き込みデータを受信し、受信したデータを E C C 符号化し、符号化したデータを R A M 4 0 のキャッシュメモリ 4 6 に記憶する。なお、符号化しないデータをキャッシュメモリ 4 6 に記憶し、N A N D メモリ 1 6 に書き込む時に符号化するようにしてもよい。

【 0 1 0 7 】

つぎに、S S D C 4 1 は F B T 6 0 を読み出し (ステップ S 1 4 1)、F B T 6 0 からフリーブロックの物理ブロック I D を取得する。フリーブロックが存在しない場合は (ステップ S 1 4 2 : N o)、S S D C 4 1 は後述する N A N D メモリ 1 6 の整理 (N A N D 整理) を行い (ステップ S 1 4 3)、この整理の後、F B T 6 0 を読み出し (ステップ S 1 4 4)、F B T 6 0 からフリーブロックの物理ブロック I D を取得する。S S D C 4 1 は、物理ブロック I D を取得したフリーブロックに対し、消去動作を行う。消去エラーが発生した場合は、当該物理ブロック I D を B B T 6 1 に追加し、F B T 6 0 から削除し、S 1 4 1 からやり直してフリーブロックを再取得する。なお、一度消去エラーが発生した物理ブロックであっても、再度消去動作を行うと消去エラーが発生せずに正常に消去できることがあるため、バッドブロック数の不必要な増大を防止するという観点では、F B T 6 0 や A B T 6 2 に統計情報 X 0 6 としてのブロックごと消去エラー発生回数の項目を各ブロックごとに設け、ブロックの消去エラーが発生した場合にこれをインクリメントするようにし、ブロックごと消去エラー発生回数が所定値以上になった場合に当該ブロックを B B T 6 1 に登録するようにするのが望ましい。さらに望ましくは、連続して消去エラーが発生する物理ブロックのみをバッドブロック化するために、S S D C 4 1 は、前記「ブロックごと消去エラー発生回数」のかわりに「ブロックごと消去連続エラー回数」の項目をもうけ、ブロックの消去エラーが発生した場合にこれをインクリメントするようにし、消去をエラー無く行えた場合にこれをゼロにリセットするようにし、「ブロックごと消去連続エラー回数」が所定値以上になった場合に当該ブロックを B B T 6 1 に登録するようにするのが望ましい。

【 0 1 0 8 】

つぎに、S S D C 4 1 は、書き込み命令で指定されている L B A が未書き込み状態であるか否かを検索するために、先の図 1 6 に示した正引き変換手順にしたがって L B A に対応する有効なデータが N A N D メモリ 1 6 に記憶済みであるか否かを判定する (ステップ S 1 4 5、S 1 4 6)。

【 0 1 0 9 】

L B A が未書き込み状態である場合は (ステップ S 1 4 6 : Y e s)、S S D C 4 1 は、キャッシュメモリ 4 6 に記憶している受信データをフリーブロックに書き込み (ステップ S 1 4 7)、書き込みを行ったフリーブロック (新物理ブロック) の I D 及びその消去回数を A B T 6 2 に登録し、書き込みを行った物理ブロックの I D を F B T 6 0 から削除する (ステップ S 1 5 1)。この際、受信データの L B A をトラック単位の区画 (トラック区画) で区切り、トラック区画内がデータで埋め尽くされているか否かを判定することで、トラック管理するかクラスタ管理するかを判定する (ステップ S 1 5 2)。すなわち、トラック区画内がデータで埋め尽くされている場合は、トラック管理となり、トラック区画内がデータで埋め尽くされていない場合は、クラスタ管理となる。クラスタ管理の場合は、クラスタテーブル 6 4 を書き換えて、L B A に新物理ブロック I D を関連付け (ステップ S 1 5 3)、さらにトラックテーブル 6 3 を書き換えて、L B A に無効な物理ブロック I D (例えば、ヌル) を関連付ける。トラック管理の場合は、トラックテーブルを書き換えて、L B A に新物理ブロック I D を関連付ける (ステップ S 1 5 4)。

【 0 1 1 0 】

一方、ステップ S 1 4 6 において、L B A が未書き込み状態でない場合は、S S D C 4 1 は正引き変換により得られた物理ブロック I D をもとに、対応する物理ブロック内全データを N A N D メモリ 1 6 から読み出して、R A M 4 0 に書き込む (ステップ S 1 4 8)。そして、キャッシュメモリ 4 6 に記憶されているデータと N A N D メモリ 1 6 から読み

10

20

30

40

50

出してRAM 40に書き込んだデータとをRAM 40上で合成し(ステップS 149)、合成したデータをフリーブロックに書き込む(ステップS 150)。

【0111】

なお、ステップS 150で書き込みエラーが発生した場合は、当該物理ブロックIDをBBT 61に追加し、FBT 60から削除し、ステップS 141からやり直してフリーブロックを再取得する。なお、一度書き込みエラーが発生した物理ブロックであっても、再度書き込み動作を行うと書き込みエラーが発生せずに正常に書き込みできることがあるため、バッドブロック数の不必要な増大を防止するという観点では、FBT 60やABT 62に統計情報X 05としてのブロックごと書き込みエラー発生回数の項目を各ブロックごとに設け、ブロックの書き込みエラーが発生した場合にこれをインクリメントするようにし、「ブロックごと書き込みエラー発生回数」が所定値以上になった場合に当該ブロックをBBT 61に登録するようにするのが望ましい。さらに望ましくは、連続して書き込みエラーが発生する物理ブロックのみをバッドブロック化するために、SSDC 41は、前記「ブロックごと書き込みエラー発生回数」のかわりに「ブロックごと書き込み連続エラー回数」の項目をもうけ、ブロックの書き込みエラーが発生した場合にこれをインクリメントするようにし、書き込みをエラー無く行えた場合にこれをゼロにリセットするようにし、「ブロックごと書き込み連続エラー回数」が所定値以上になった場合に当該ブロックをBBT 61に登録するようにするのが望ましい。

10

【0112】

SSDC 41は、書き込みを行ったフリーブロック(新物理ブロック)のID及びその消去回数をABT 62に登録し、書き込みを行った物理ブロックのIDをFBT 60から削除する(ステップS 151)。LBAがクラスタ管理である場合は、クラスタテーブル64の旧物理ブロックIDを新物理ブロックIDに書き換える(ステップS 152、S 153)。トラック管理の場合は、トラックテーブルの旧物理ブロックIDを新物理ブロックIDに書き換える(ステップS 152、S 154)。さらに、SSDC 41は、旧物理ブロックID及びその消去回数をFBT 60に追加し、旧物理ブロックID及びその消去回数をABT 62から削除する(ステップS 155)。SSDC 41は以上の書き込み処理の内容を統計情報65に反映する。

20

【0113】

(NAND整理)

通常、記憶部2の全LBAの容量(全論理セクタ数)は、記憶部2のNANDメモリ16の全容量よりも小さく設計されているため(オーバープロビジョニング)、書き込み動作がトラック単位で書き込まれ続ける限りはフリーブロックが枯渇することはない。一方、未書き込みLBAに対してクラスタ単位の書き込みが多数発生した場合、クラスタ単位の書き込み一つに対してクラスタよりも容量の大きい物理ブロックが割り当てられることになるため、書き込まれるデータ容量よりも多くのNANDメモリ16物理ブロックを必要とすることになり、それによりフリーブロックが枯渇する可能性がある。フリーブロックが枯渇した場合は、以下に示すNANDメモリ16の整理によってフリーブロックを新たに確保することができる。

30

【0114】

図21を用いてSSDC 41によるNAND整理を説明する。物理ブロックに記憶されている全てのクラスタが有効クラスタであるとは限らず、有効クラスタに該当しない無効クラスタはLBAに対応付けられていない。有効クラスタとは最新のデータを記憶しているクラスタであり、無効クラスタとは同一LBAのデータが他の場所に書きこまれ、参照されることがなくなったクラスタである。物理ブロックは無効クラスタの分だけデータに空きがあることになり、有効クラスタのデータを集めて違うブロックに書き直すNAND整理を実行することでフリーブロックを確保することができる。

40

【0115】

まず、SSDC 41は、選択物理ブロックID = iを0にセットし、空き領域累積量Sを0にセットする(ステップS 160)。SSDC 41は、このIDがi = 0の物理プロ

50

ックがトラックテーブル 6 3 にエントリされているか否かを判定する (ステップ S 1 6 1)。トラックテーブルにエントリされている場合は i を + 1 し (ステップ S 1 6 2)、つぎの番号の ID を持つ物理ブロックについて同様の判定を行う (ステップ S 1 6 1)。すなわち、物理ブロック ID がトラックテーブル 6 3 に含まれている場合は、この物理ブロックのデータはトラック管理であるため、NAND 整理対象に含めない。

【0116】

S S D C 4 1 は、ID = i の物理ブロックがトラック管理でない場合は (ステップ S 1 6 1 : No)、つぎにクラスタテーブル 6 4 を参照し、ID = i の物理ブロックに含まれる有効クラスタのアドレスを全て取得する (ステップ S 1 6 3)。そして、S S D C 4 1 は、取得した有効クラスタの総容量分のサイズ v を求め (ステップ S 1 6 4)、 $v <$ 物理
10
ブロックサイズであるときは (ステップ S 1 6 5)、現在物理ブロックの ID を NAND 整理対象ブロックリストに加える (ステップ S 1 6 6)。さらに、S S D C 4 1 は、取得クラスタ累計量 S に現在物理ブロックの取得クラスタ容量 v を加算し、取得クラスタ累計量 S を更新する (ステップ S 1 6 7)。

【0117】

ステップ S 1 6 5 で、 $v <$ 物理ブロックサイズでないとき、あるいはステップ S 1 6 8 で取得クラスタ累計量 S が物理ブロックサイズに到達していない場合は、S S D C 4 1 は、 i を + 1 し (ステップ S 1 6 2)、つぎの番号の ID を持つ物理ブロックについて、ステップ S 1 6 1 ~ S 1 6 7 の手順を前記同様に実行する。そして、ステップ S 1 6 8 で、
20
取得クラスタ累計量 S が物理ブロックサイズに到達するまで、ステップ S 1 6 1 ~ S 1 6 7 の手順を繰り返す。

【0118】

そして、ステップ S 1 6 8 において、取得クラスタ累計量 S が物理ブロックサイズに到達した場合は、S S D C 4 1 は、NAND 整理対象ブロックリスト上の全物理ブロックについての全有効クラスタのデータを NAND メモリ 1 6 から読み出して R A M 4 0 に書き込み (ステップ S 1 6 9)、さらに、NAND 整理対象ブロックリスト上の全物理ブロックに対して消去処理を行い (ステップ S 1 7 0)、消去処理を行った全物理ブロックを A B T 6 2 から削除して F B T 6 0 に追加する (ステップ S 1 7 1)。その際、消去回数をインクリメントする。なお、ステップ S 1 7 0 で行う消去動作の対象は、ステップ S 1 7 2 でデータを書き込む対象のブロックに限定してもよく、ブロックの消去回数を抑制する
30
という観点ではそのように行うことが望ましい。

【0119】

消去エラーが発生した場合は、当該物理ブロック ID を B B T 6 1 に追加し、F B T 6 0 から削除する。なお、一度消去エラーが発生した物理ブロックであっても、再度消去動作を行うと消去エラーが発生せずに正常に消去できることがあるため、バッドブロック数の不必要な増大を防止するという観点では、F B T 6 0 や A B T 6 2 に「ブロックごと消去エラー発生回数」の項目を各ブロックごとに設け、ブロックの消去エラーが発生した場合にこれをインクリメントするようにし、ブロックごと消去エラー発生回数が所定値以上になった場合に当該ブロックを B B T 6 1 に登録するようにするのが望ましい。さらに望ましくは、連続して消去エラーが発生する物理ブロックのみをバッドブロック化するため
40
に、S S D C 4 1 は、前記「ブロックごと消去エラー発生回数」のかわりに「ブロックごと消去連続エラー回数」の項目をもうけ、ブロックの消去エラーが発生した場合にこれをインクリメントするようにし、消去をエラー無く行えた場合にこれをゼロにリセットするようにし、「ブロックごと消去連続エラー回数」が所定値以上になった場合に当該ブロックを B B T 6 1 に登録するようにするのが望ましい。

【0120】

そして、S S D C 4 1 は、F B T 6 0 から新たなフリーブロックを取得し、取得したフリーブロックに対し R A M 4 0 に書き込んだデータを書き込み (ステップ S 1 7 2)、データを書き込んだフリーブロックの物理ブロック ID 及び当該ブロックの消去回数を A B T 6 2 に追加し、さらにデータが書き込まれたブロックのブロック ID を F B T 6 0 から
50

削除する（ステップS 1 7 3）。さらに、S S D C 4 1は、今回のN A N D整理に対応するように、クラスタテーブル6 4におけるクラスタアドレス、物理ブロックID及び物理ブロック内ページアドレスを更新する（ステップS 1 7 4）。S S D C 4 1は上記N A N D整理の処理内容を統計情報6 5に反映する。

【0 1 2 1】

なお、ステップS 1 7 2で書き込みエラーが発生した場合は、当該物理ブロックIDをB B T 6 1に追加し、F B T 6 0から削除し、フリーブロックを再取得する。なお、一度書き込みエラーが発生した物理ブロックであっても、再度書き込み動作を行うと書き込みエラーが発生せずに正常に書き込みできることがあるため、パッドブロック数の不必要な増大を防止するという観点では、F B T 6 0やA B T 6 2に「ブロックごと書き込みエラー発生回数」の項目を各ブロックに設け、ブロックの書き込みエラーが発生した場合にこれをインクリメントするようにし、「ブロックごと書き込みエラー発生回数」が所定値以上になった場合に当該ブロックをB B T 6 1に登録するようにするのが望ましい。さらに望ましくは、連続して書き込みエラーが発生する物理ブロックのみをパッドブロック化するために、S S D C 4 1は、前記「ブロックごと書き込みエラー発生回数」のかわりに「ブロックごと書き込み連続エラー回数」の項目をもうけ、ブロックの書き込みエラーが発生した場合にこれをインクリメントするようにし、書き込みをエラー無く行えた場合にこれをゼロにリセットするようにし、「ブロックごと書き込み連続エラー回数」が所定値以上になった場合に当該ブロックをB B T 6 1に登録するようにするのが望ましい。

【0 1 2 2】

なお、図2 1の手順では、フリーブロックにデータを詰め込むことを優先するN A N D整理を行ったが、ステップS 1 6 4で、 v を物理ブロックサイズから取得したクラスタの容量を減算することによって求め、ステップS 1 6 5で $v > 0$ か否かを判定し、 $v > 0$ の場合はステップS 1 6 8に移行し、 $v > 0$ でない場合はステップS 1 6 2に移行させるようにすることで、フリーブロックを確保することを優先するN A N D整理を行うようにしてもよい。

【0 1 2 3】

（削除通知）

つぎに、図2 2を用いてS S D C 4 1による削除通知処理について説明する。削除通知は、情報処理装置1 1 1上のO S 1 0 0によってデータの削除が行われた場合に、情報処理装置1 1 1から記憶部2や記憶部3に対して送信される命令である。削除通知処理に使用されるコマンドは通称トリムコマンドと呼ばれる。トリムコマンドの例として、例えば、Information technology ATA/ATAPI Command Set-3 (ACS-3)に記載されている06h Data Set Managementコマンドや、NVM Express Revision 1.1に記載されている09h Dataset ManagementコマンドのDeallocateがあげられる。これは、O S 1 0 0上でデータが削除された場合、削除されたデータの存在する論理アドレス領域（L B A領域）を、L B A及びセクタ数の組みよりなるL B A Range Entryとして外部記憶装置に通知することにより、記憶部2や記憶部3上でもその領域を空き領域として扱うことができる方式である。削除通知により、S S D C 4 1はフリーブロックを新たに確保することができる。なお、トリムコマンドの機能は、前記コマンドだけでなく、たとえば、Information technology ATA/ATAPI Command Set-3 (ACS-3)で記述されているSCT Command Transportや、NVM Express Revision 1.1に記載されている08h Write Zeroesコマンドや、その他ベンダー独自のコマンドなどその他コマンドによって実現してもよい。

【0 1 2 4】

記憶部2および記憶部3が情報処理装置1 1 1から削除通知を受信した場合は（ステップS 1 8 0）、S S D C 4 1は、削除通知に含まれるL B Aを先の図1 6に示した手順に従ってL B A正引き変換する。S S D C 4 1は、削除通知に含まれるL B Aがトラック管理である場合は（ステップS 1 8 1：Y e s）、物理ブロックIDをF B T 6 0に追加してA B T 6 2から削除する（ステップS 1 8 4）。一方、S S D C 4 1は、削除通知に含まれるL B Aがクラスタ管理である場合は（ステップS 1 8 1：N o）、物理ブロックに

対応する全クラスタをクラスタテーブル 6 4 から削除し (ステップ S 1 8 2)、トラックテーブル 6 3 において、L B A に対応するトラックに対応する物理ブロック I D に適当な有効値 (例えば F F F F) を記入し (ステップ S 1 8 3)、物理ブロック I D を F B T 6 0 に追加して A B T 6 2 から削除する (ステップ S 1 8 4)。S S D C 4 1 は、N A N D 整理以外に、削除通知処理によってもフリーブロックを確保することができる。

【 0 1 2 5 】

このような N A N D 整理により、書き込みに対し十分な数のフリーブロックを確保できるのが通常である。N A N D 整理によっても、書き込みに対し十分な数のフリーブロックを確保できなかった場合は、統計情報 6 5 の N A N D 整理失敗フラグを 1 にして、情報処理装置 1 1 1 による統計情報 6 5 の取得を通じて S S D C 4 1 がフリーブロックを確保できなかったことを情報処理装置 1 1 1 に通知できるようにすることが望ましい。たとえば、N A N D 整理失敗フラグが 1 になってから実際に記憶部 2 が動作しなくなるまでの時間に猶予をもたせるという観点では、

(N A N D 整理をして確保できたフリーブロック数) < (書き込みに必要なフリーブロック数) + (マージン)

の条件を満たす場合に N A N D 整理失敗フラグを 1 にセットして信頼性情報として情報処理装置 1 1 1 に通知することが望ましい。

【 0 1 2 6 】

上記 N A N D 整理は、情報処理装置 1 1 1 からライト要求を受信した時だけでなく、情報処理装置から最後に命令を受信してから所定時間経過した時、または情報処理装置 1 1 1 からスタンバイやアイドルやスリープ状態に移行するコマンドを受信した時などに実行してもよいし、ACS-3に記載のSCT Command Transportやその他ベンダーコマンドなどを通じて、N A N D 整理を開始する命令を S S D C 4 1 が情報処理装置 1 1 1 から受信した時などに実行してもよい。

【 0 1 2 7 】

(エラー処理)

次に、図 2 3 を用いて S S D C 4 1 による N A N D メモリ 1 6 に関するエラー処理について説明する。情報処理装置 1 1 1 からのライト要求に対する処理や N A N D 整理処理など各種処理は通常は上記のように行われるが、N A N D メモリ 1 6 に対する書き込み動作 (プログラム動作) で書き込みエラーが発生する場合、N A N D メモリ 1 6 に対する消去動作 (イレーズ動作) で消去エラーが発生する場合、N A N D メモリ 1 6 に対する読み出し動作の際に E C C エラー (誤り訂正処理の失敗) が生じる場合などがあり、これらに対する例外処理が必要となる。

【 0 1 2 8 】

S S D C 4 1 は、上記の何れかのエラーが発生した場合 (ステップ S 1 9 0)、エラーが発生した物理ブロックを B B T 6 1 に追加し (ステップ S 1 9 1)、エラーが発生した物理ブロックを A B T 6 2 及び F B T 6 0 から削除することで (ステップ S 1 9 2)、以後は、エラーが発生した物理ブロックにアクセスできないようにする。この際、エラーが発生した物理ブロックのデータを別の物理ブロックにコピーしてもよい。S S D C 4 1 は上記エラー処理を統計情報 6 5 に反映させる。

【 0 1 2 9 】

上記では、エラー処理の一例を、読み出し処理、書き込み処理、N A N D 整理処理に関して示したが、エラー処理はこれらの例に限らず、N A N D メモリ 1 6 に対する全ての読み出し処理、書き込み処理、消去処理に対して適用可能である。

【 0 1 3 0 】

(寿命到達処理)

情報処理装置 1 1 1 が記憶部 2 を利用しているうちに、統計情報 6 5 内に格納された値は悪化していき、記憶部 2 は寿命に到達する。たとえば、情報処理装置 1 1 1 が記憶部 2 を利用しているうちに、N A N D メモリ 1 6 の各ブロックの信頼性が劣化していき、バッドブロックの数が増えていき、フリーブロック数とアクティブブロック数の和が減ってい

10

20

30

40

50

くことになる。さらに情報処理装置 1 1 1 が記憶部 2 を使用していると、SSDC 4 1 が NAND 整理を実行しても、書き込み処理を行うのに十分なフリースペース数を確保できなくなる。これが記憶部 2 の寿命到達の一例である。以下では、記憶部 2 の寿命が到達した場合の制御プログラム 2 0 0 の処理を示す。

【0131】

制御プログラム 2 0 0 は起動されると、主メモリ 6 に常駐して記憶部 2 の統計情報 6 5 などの信頼性情報を監視する。記憶部 2 の統計情報 6 5 を常に監視するためには、OS 1 0 0 が領域 1 6 B (または 2 0 B) から領域 6 A に読み出されるときまたはその直後に、制御プログラム 2 0 0 が領域 1 6 B (または領域 2 0 B) から領域 6 B に読み出されるのが望ましい(図 2、図 3 参照)。

10

【0132】

制御プログラム 2 0 0 は、例えば図 2 4 に示すように、一定時間おき(たとえば 1 分おき)または一定処理ごと(たとえば記憶部 2 へのアクセスが 1 0 0 0 アクセスごとや、記憶部 2 と授受したデータが 1 0 G B ごと、など)に記憶部 2 から統計情報 6 5 などの信頼性情報を取得する。なお、記憶部 2 が寿命に近づくほど信頼性情報を取得する頻度を上げることで、より堅牢なユーザデータの保護が可能になる。統計情報 6 5 を取得する方式として、例えば、INCITS ACS-3 に記述されている、メモリの自己診断機能である S.M.A.R.T (Self-Monitoring Analysis and Reporting Technology) のコマンドである、B0h/D0h SMART READ DATA コマンドや B0h/D5h SMART READ LOG コマンドを用いてもよい。又は、NVM Express Revision 1.1 に記述されている 02h Get Log Page コマンドを用いてもよい。又は ACS-3

20

【0133】

図 2 7 は、統計情報 6 5 をもとに生成されるか、あるいは統計情報 6 5 に直接格納されて、SSDC 4 1 から情報処理装置 1 1 1 に信頼性情報として送信されるテーブルデータ例を示すものである。SSDC 4 1 は、統計情報 6 5 として S.M.A.R.T を用いる場合、図 2 7 に示すように、統計情報 6 5 の構成要素それぞれに対し、属性 ID (attribute ID) を割り当てる。SSDC 4 1 は、統計情報 6 5 の構成要素の一部のみに属性 ID を割り当ててもよい。SSDC 4 1 は、これら統計情報 6 5 の構成要素において、信頼性が悪化すればするほど値が増加する要素に関しては、たとえば以下のようにして規格化を行った値である attribute value を計算する。

30

$$\text{attribute value} = \text{SMAL} + \text{SMAB} \times (1 - \text{AMALR}) \times (\text{RMAX} - \text{Raw Value}) / \text{RMAX}$$

【0134】

Raw Value は統計情報 6 5 に格納される値そのものである。RMAX は信頼性保証できる Raw Value の上限値である。SMAB は attribute value の初期値としてあらかじめ設定されているパラメータであり、たとえば 100 が設定される。SMAL (=attribute Threshold) は、Raw Value が上限値である RMAX に等しくなった時に attribute value が到達すべき値であり、あらかじめ設定されているパラメータである。SMAL とし、たとえば 30 が設定される。AMALR は、 $\text{AMALR} = \text{SMAL} / \text{SMAB}$ の関係から導かれるパラメータであり、 $0 < \text{AMALR} < 1$ となる。このようにして、SSDC 4 1 は、smart 情報の attribute value (図 2 7 の「Value」) を計算して制御プログラム 2 0 0 に送信する。attribute Threshold は、図 2 7 の「Threshold」であり、Raw Value は、図 2 7 の「Raw Data」である。

40

【0135】

たとえば、SSDC 4 1 が Raw Value として消去回数平均値(統計情報 X 0 4)を用いる場合、現在の消去回数平均値が 1000 回ならば Raw Data=1000 となり、NAND メモリ 1 6 の信頼性保証できる消去回数が 10000 回の時は RMAX=10000 となり、消去回数=0 の初期状態で attribute value の初期値が 100 となっているよう SSDC 4 1 が設計されている場合には SMAB=100 となり、消去回数が RMAX=10000 に到達すると attribute value は SMAL に到達する。

【0136】

50

S S D C 4 1 は、統計情報 6 5 の構成要素において、信頼性が悪化すればするほど値が減少する要素に関しては、たとえば以下のようにして規格化を行った値である attribute value を計算する。

$$\text{attribute value} = \text{SMAL} + \text{SMAB} \times (1 - \text{AMALR}) \times (\text{Raw Value} - \text{RMIN}) / (\text{RINIT} - \text{RMIN})$$

RMIN は信頼性保証できる Raw Value の下限値である。RINIT は Raw Value の初期値である。

【 0 1 3 7 】

RMAX、AMALR および SMAB は各 X 0 1 ~ X 3 2 に対して、それぞれ別の値を採用することができる。SMAB = 1 0 0 で、AMALR = 0 . 3 を採用した時、採用対象の統計情報に関して attribute value の最良値が 1 0 0 (たとえば出荷直後に 1 0 0) になっており、信頼性が劣化するにつれて徐々に減少していき、記憶部 2 が信頼性保証できなくなったとき (統計情報の Raw Value が RMAX に等しいかそれ以上になったとき) または記憶部 2 が寿命到着直前になったとき、attribute value が 3 0 かそれ以下の値に到達することになる。なお、制御プログラム 2 0 0 は、Attribute Value が Threshold を超過しているか否かを検知する手段として、ACS-3 に記載のコマンドである B0h/DAh SMART RETURN STATUS を用い、当該コマンドの Output から Attribute Value が Threshold を超過しているか否かで寿命到達を判定してもよい。

【 0 1 3 8 】

記憶部 2 の製造業者は、たとえば図 2 8 のように統計情報の Raw Value と記憶部 2 の不良率の関係を開発段階に導き出し、不良率が許容値を超える時の Raw Value を RMAX として採用することが望ましい。たとえば、記憶部 2 の開発段階で、多数個 (たとえば 1 0 0 個) の試験用記憶部 2 群に対して高温で書き込み動作を繰り返しながら、書き込んだデータが一定時間以上正しく記憶され続けるかを検証する摩耗試験を行い、同時に統計情報をモニタしつづけ、不良率が一定割合に到達する時点での統計情報の Raw Value を RMAX として採用すればよい。たとえば、摩耗した記憶部 2 を高温状態である時間以上放置して、その後記憶部 2 の温度を下げ、記憶部 2 に対して読み出し動作を行い、読み出したデータが ECC 訂正できない場合 (または ECC 訂正できないデータが一定数以上の場合)、これを記憶部 2 の不良として定義し、不良数を、同等の試験を行った記憶部 2 の個数で割った値を不良率として採用すればよい。この不良率が、前記許容できる不良率を統計的に有意に下回る Raw Value を RMAX として採用すればよい。RMAX にある程度マージンをもたせて、

$$\text{RMAX}' = \text{RMAX} - \text{マージン}$$

とした RMAX' を RMAX として採用してもよい。

【 0 1 3 9 】

図 2 7 の「Worst」を制御プログラム 2 0 0 が記憶部 2 の寿命を診断する指標として採用してもよい。「Worst」は attribute value の最悪値として S S D C 4 1 により計算される。たとえば、Worst は、たとえば記憶部 2 の出荷後の (または製造後の)、attribute value の最小値である。あるいは、Worst として、過去ある一定時間範囲内の attribute value の最小値を Worst Value として採用してもよいし、ある一定回数 (一定データ量) 通信または処理が行われた過去にさかのぼって、その過去から現在に至るまでの最小値を worst value として採用してもよい。

【 0 1 4 0 】

図 2 7 の「Raw Data」(Raw Value) を制御プログラム 2 0 0 が記憶部 2 の寿命を診断する仕様として採用してもよい。統計情報の Raw Value が Raw Data として記憶部 2 から制御プログラム 2 0 0 に送信される。この場合、制御プログラム 2 0 0 は、RMAX を制御プログラム 2 0 0 内にすでに保持しているか記憶部 2 から別途読み出すかその他記憶装置から読み出すかして RMAX を取得し、RMAX と Raw Data を比較し、Raw Data > RMAX または Raw Data

RMAX となった時、記憶部 2 が寿命に到達したと判定する。たとえば、NAND 整理失敗フラグの場合、これが 1 である場合に記憶部 2 が寿命に到達したと判定する。たとえば、バッドブロック数総計の場合、これが所定値を上回った場合に記憶部 2 が寿命に到達したと判定する。記憶部 2 から情報処理装置 1 1 1 に送信される Raw Data として、必ずしも統計情報の Raw Value を出力する必要はなく、たとえば、統計情報の Raw Value を四則演算し

た値をRaw DataとしてSSDC41が制御プログラム200に送信し、同じくRMAXを四則演算した値と比較することで判定を行ってもよい。また、統計情報のRaw Valueを暗号化するなど難読化したデータをRaw DataとしてSSDC41が制御プログラム200に送信し、SSDC41がこれを復号化して復号化後のデータをRMAXと比較することで判定を行ってもよい。

【0141】

上記のようにして、制御プログラム200は、記憶部2が寿命に到達したか否か（記憶部2が異常状態であるか否か）を判定し、記憶部2が寿命に到達したと判定された場合（記憶部2が異常状態であると判定された場合）、後述する寿命到達時処理（ステップS205）に移行する。統計情報65は、統計情報X01～X32以外にも種々の形態をとりうるが、本発明はこれらに対しても適用可能である。また、統計情報と不良率の関係に正の相関関係が存在する場合にだけでなく、統計情報と不良率の関係に負の相関関係が存在する場合に対しても本発明は適用可能である。たとえば、記憶部2が出荷後に体験した最低温度などである。その場合、RMAXの代わりに、信頼性保証できる下限値RMINを採用し、統計情報がRMINを下回った場合に記憶部2が寿命に到達したと判定すればよい。

10

【0142】

本実施形態では、制御プログラム200は図24に示すようにS.M.A.R.Tを用いて一定時間おき（たとえば1分おき）に統計情報の取得を行う（ステップS200：Yes）。制御プログラム200は、統計情報取得コマンドであるACS-3に記載のB0h/D0h SMART READ DATAを発行し（ステップS201）、記憶部2から統計情報を含むデータを受信し（ステップS202）、この受信したデータを診断する（ステップS203）。診断方法は前述したとおりである。ステップS204において、記憶部2が寿命に到達したと制御プログラム200が判定した時、または記憶部2が寿命到達目前であると制御プログラム200が判定した時（ステップS204：Yes）、制御プログラムは寿命到達時処理に移行する（ステップS205）。記憶部2が寿命に到達していなくても、統計情報があらかじめ定められたRMAXを上回ったり、あるいは通常動作ではありえない異常値を示した場合などにおいても、ステップS205の処理に移行することが望ましい。

20

【0143】

特許文献22や特許文献23の寿命予測技術を用いて記憶部2の寿命予測を行い、記憶部2の寿命がある一定期間後に迫っていると判定された場合に、寿命到達時処理に移行してもよい。

30

【0144】

統計情報65以外の信頼性情報を用いて、寿命到達時処理に移行してもよい。たとえば、図29に示すように、制御プログラム200は、OS100が記憶部2から受信する応答情報（図7参照）をOS100から取得（監視）して信頼性情報として用い（ステップS210）、エラー応答であった場合（ステップS211）、記憶部2が異常状態に至ったと判定し、寿命到達時処理に移行する（ステップS205）。監視する応答はどのようなコマンドに対する応答でもよいが、たとえばACS-3記載の61h WRITE FPDMA QUEUEDや35h WRITE DMA EXTや、NVM Express Revision 1.1記載の01h Writeコマンドなど、記憶部2に対する書き込みコマンドに対する応答のみを監視することが、CPU5への負荷低減の観点から望ましい。特に、記憶部2が特許文献3の発明を採用したSSDである場合、記憶部2が寿命に至ると、記憶部2への書き込みコマンドに対する応答がエラーで返ってくるようになるため、統計情報の取得無しに寿命到達を判定することができる。なお、記憶部2が特許文献3の発明を採用したSSDではない場合であっても本発明が適用できることはもちろんである。

40

【0145】

なお、特許文献3の発明を採用した記憶部2である場合、記憶部2が書き込みコマンドに対してエラーを返す状態である時には、後述する記憶部2のBoot Loader領域の書き換えを行う場合は、特許文献3のRead Onlyモード状態において特殊な書き込みコマンド（例えばACS-3に記述されているSCT Command Transportやその他ベンダー独自コマンドなど

50

）に対してはエラーを返さないよう S S D C 4 1 を構成しておき、前述した特殊な書き込みコマンドを用いて記憶部 2 に書き込みを行うことが望ましい。なお、この特殊な書き込みコマンドは、記憶部 2 以外の記憶装置への書き込みには用いなくてもよい。あるいは、O S 1 0 0 が、書き込みコマンドとしてある書き込みコマンド（たとえば 6 1 h WRITE FPDMA QUEUED）のみを使うような O S である場合には、S S D C 4 1 が特許文献 3 の Read Only に至ると、当該書き込みコマンド（たとえば 6 1 h WRITE FPDMA QUEUED）に対してエラーを返すように S S D C 4 1 を構成し、別の書き込みコマンド（たとえば 3 0 h WRITE SECTOR(S)）に対してはエラーを返さないように S S D C 4 1 を構成し、別の書き込みコマンド（たとえば 3 0 h WRITE SECTOR(S)）を用いて記憶部 2 に対して Boot Loader 領域の書き込みを行うようにしてもよい。

10

【 0 1 4 6 】

監視対象のコマンドは、書き込みコマンド以外のコマンドであってもよいのはもちろんである。たとえば、コマンド応答として ACS-3 記載の B0h/D4h SMART EXECUTE OFF-LINE IMMEDIATE の応答 (Outputs) やレポートを監視してもよいし、90h EXECUTE DEVICE DIAGNOSTIC の応答を監視してもよい。たとえば、制御プログラム 2 0 0 は、記憶部 2 から取得される SMART の self-test の結果を用いて、寿命到達時処理に移行してもよい。制御プログラム 2 0 0 は、ACS-3 に記載の B0h/D4h SMART EXECUTE OFF-LINE IMMEDIATE コマンドを記憶部 2 に送信することで S S D C 4 1 は self-test を実行し、ACS-3 に記載の B0h/D0h SMART READ DATA や B0h/D5h SMART Read Log コマンドを記憶部 2 に送信することで self-test の結果を信頼性情報として取得する。たとえば、制御プログラム 2 0 0 は取得した self-test の結果にエラーが含まれていることを検知したとき、記憶部 2 が寿命に到達したと判定する。

20

【 0 1 4 7 】

あるコマンド応答がエラーであっても、同コマンドを再度送信するとエラーでない可能性があり、この場合は記憶部 2 が寿命に到達していない可能性があるため、再現性のあるコマンドエラーが発生した場合にのみ寿命到達時処理を行う観点では、コマンドエラーが複数回発生した場合に寿命到達時処理を行うことが望ましい。さらに、エラー再現性を厳密に判定する観点では、コマンドエラーが複数回連続で発生した場合に寿命到達時処理を行うことが望ましい。あるいは、図 3 0 に示すように、記憶部 2 へのコマンド監視中にコマンド応答がエラーで返ってきた場合（ステップ S 2 2 0、S 2 2 1:Yes）、制御プログラム 2 0 0 または O S 1 0 0 が同コマンドを記憶部 2 に再度送信（コマンドリトライ）し（ステップ S 2 2 2）、リトライしたコマンドがエラーした場合に（ステップ S 2 2 3:Yes）、寿命到達時処理を行うようにしてもよい（ステップ S 2 0 5）。

30

【 0 1 4 8 】

また、制御プログラム 2 0 0 は、情報処理装置 1 1 1 から取得される信頼性情報を用いて、寿命到達時処理に移行してもよい。たとえば、情報処理装置 1 1 1 内に温度計が設置されているとき、制御プログラム 2 0 0 は温度計から出力される温度を監視し、温度が上限値を上回った場合や下限値を下回った場合に、信頼例劣化時処理としての寿命到達時処理を行うようにしてもよい。

【 0 1 4 9 】

40

（通常状態のデータ構成）

図 3 1 は寿命到達時処理 S 2 0 5 を行うより以前の、情報処理装置 1 1 1 が管理するデータの構成例である。前述のとおり、情報処理装置 1 1 1 は N A N D メモリ 1 6 の物理アドレスを直接指定してデータの読み出し、書き込みを記憶部 2 に要求するのではなく、論理アドレスである L B A を指定してデータの読み出し、書き込み命令を記憶部 2 に送信する。記憶部 2 内の S S D C 4 1 は、L B A と、N A N D メモリ 1 6 の物理アドレスの対応を管理情報 4 4 によって動的に割り当てる。このように、情報処理装置 1 1 1 が直接管理できるデータは L B A によって管理されており、情報処理装置 1 1 1 によって管理可能なアドレス空間として、L B A 領域 2 0 0 1 が記憶部 2 に割り当てられている。L B A 領域 2 0 0 1 は、Boot Loader 領域 2 0 0 2 と、メタデータ領域 2 0 0 3 と、ユーザデータ領

50

域 2 0 0 4 からなる。なお、これらの領域のうち一部の領域を ACS-3 の SMART Read Log コマンドや Read Log コマンドでアクセス可能なログページアドレスに割り当ててもよい。記憶部 2 は、論理ドライブ 4 に割り当てられている。

【 0 1 5 0 】

Boot Loader 領域 2 0 0 2 は、情報処理装置 1 1 1 の起動時に読み出される領域である。本実施形態では Boot Loader 領域 2 0 0 2 は L B A の固定領域に割り当てられているとするが、情報処理装置 1 1 1 は Boot Loader 領域 2 0 0 2 を動的に割り当て可能としてもよい。なお、Boot Loader 領域の例として、たとえば Master Boot Record (M B R) があり、M B R ではたとえば L B A = 0 x 0 0 0 0 の合計 1 論理セクタ (合計 5 1 2 Byte) の領域が固定 Boot Loader 領域として割り当てられている。また、Boot Loader 領域の例として、たとえば G U I D パーティションテーブル (G P T) が存在する。本実施形態では、Boot Loader 領域 2 0 0 2 は、図 3 2 に示すように、メタデータ領域 2 0 0 3 の先頭 L B A を格納したメタデータへのポインタ領域 2 0 0 5 と、ディスクステータスを格納したステータス格納領域 2 0 0 6 と、データ引き継ぎ対象記憶装置アドレスを格納したアドレス領域 2 0 0 7 よりなる。例えば、記憶部 2 がデータ引き継ぎ元になり、記憶部 3 がデータ引き継ぎ先になる場合、記憶部 2 のアドレス領域 2 0 0 7 には、記憶部 3 のディスク識別 I D が格納され、記憶部 3 のアドレス領域 2 0 0 7 には、記憶装置のディスク識別 I D が格納される。本実施形態では、ステータス格納領域 2 0 0 6 に格納されるデータは、0 ~ 5 の値を取りうる。それぞれの値は、格納先の記憶部がそれぞれ

- 0 : 初期ディスク状態
- 1 : 通常状態
- 2 : 信頼性劣化状態
- 3 : データ引き継ぎ元状態 (保護状態)
- 4 : データ引き継ぎ先状態
- 5 : 廃棄対象状態

の状態であることを示す。

【 0 1 5 1 】

情報処理装置 1 1 1 は起動時にポインタ 2 0 0 5 を読み出してメタデータ領域 2 0 0 3 の L B A を特定し、メタデータ 3 0 0 を L B A 領域 2 0 0 1 のメタデータ領域 2 0 0 3 から主メモリ 6 のメタデータ領域 6 C に読み出す。O S 1 0 0 はファイルの書き換えが発生すると主メモリ 6 のメタデータ領域 6 C のメタデータ 3 0 0 を書き換え、さらにメタデータ領域 6 C のメタデータ 3 0 0 を定期的に記憶部 2 のメタデータ領域 2 0 0 3 に退避保存したり、メタデータ 3 0 0 のジャーナルを記憶部 2 のメタデータ領域 2 0 0 3 に逐次記録したりする。

【 0 1 5 2 】

図 3 3 はメタデータ 3 0 0 の構成例である。ファイル識別 I D はアプリケーションプログラム 4 0 0 がデータを識別するために用いられるデータのアドレスまたはファイル名であり、論理ドライブアドレスは論理ドライブ 4 (図 3 1 、図 3 8 参照) を特定するために用いられる各論理ドライブに割り当てられたアドレスである。また、ディスク識別 I D は、記憶部 2 や記憶部 3 などの物理的な記憶装置を特定するために用いられる、物理的記憶装置に割り当てられたアドレスである。本実施形態においては、ディスク識別 I D として、W W N (World Wide Name) を用いる。W W N は各記憶装置それぞれに異なる値が割り当てられているため、物理的な記憶装置を区別するのに用いることができる。W W N は、たとえば、ACS-3 に記述されている ECh Identify Device コマンドで読み出されるデータの Word 108-111 に格納される。あるいは、ディスク識別 I D として、たとえば ACS-3 に記述されている ECh Identify Device コマンドで読み出されるデータの Word 10-19 に割り当てられたシリアル番号を用いてもよいし、NVM Express Revision 1.1 に記述されている 06h Identify コマンドで読み出されるデータの Byte 23:04 に割り当てられた Serial Number (S N) を用いてもよいし、ネットワーク・プロトコルにおける M A C (Media Access Control) アドレスや I P (Internet Protocol) アドレスを用いてもよい。ディスク識別 I D は記憶装置

の製造時にあらかじめ記憶装置に書き込まれていてもよいし、情報処理装置 111 との接続時に情報処理装置 111 によって新規に割り当てられたりしてもよい。メタデータ 300 の LBA にはメタデータ 300 によって紐付けされているユーザデータ領域 2004 の LBA アドレスが格納されている。メタデータ 300 の LBA にメタデータ領域 2003 や Boot Loader 領域 2002 の LBA など、ユーザデータ領域 2004 以外の LBA が格納されてもよい。セクタカウントは、データ長を示している。

【0153】

メタデータ 300 は、OS 100 がファイル識別 ID および論理ドライブアドレスからディスク識別 ID と LBA とセクタカウントに正引きしたり、OS 100 がディスク識別 ID と LBA とセクタカウントから論理ドライブアドレスとファイル識別 ID を逆引きしたりするのに用いられる。通常、アプリケーションプログラム 400 は記憶部 2 や LBA を直接指定して記憶装置への読み出し、書き込みを行わない。OS 100 は記憶部 2 を論理ドライブ 4 として認識（マウント）し、アプリケーションプログラム 400 が論理ドライブアドレスおよびファイル ID を OS 100 に送信すると、OS 100 がメタデータ 300 を読みだして論理ドライブおよびファイル ID に対応する記憶装置と LBA を特定し、当該記憶装置に命令と LBA を送信する。

【0154】

論理ドライブは、図 7 における OS 100 および制御プログラム 200 によって、下位の階層である 1 乃至複数台の物理ドライブやあるいはそれらの一部の LBA 領域に割り当てられ、それにより上位の階層であるアプリケーションプログラム 400 は論理ドライブを仮想的に 1 台のドライブとして認識する。本実施形態においては、寿命到達時処理 S205 以前の状態において、論理ドライブ 4 が 1 つの物理記憶装置である記憶部 2 に割り当てられる場合について述べる。なお、寿命到達時処理 S205 以前の状態であっても、複数の物理記憶装置を用いて Redundant Arrays of Inexpensive Disks (RAID) のディスクアレイ、たとえば RAID0 や RAID5 のディスクアレイを構成し、これを一つの論理ドライブ 4 として認識してもよく、その場合であっても本発明は適用可能である。

【0155】

図 34 はアプリケーションプログラム 400 が OS 100 に対して論理ドライブ 4 へのアクセス要求を行った場合のフローチャートである。アプリケーションプログラム 400 がファイル読み出し命令と論理ドライブアドレスとファイル識別 ID を OS 100 に送信すると（ステップ S300）、OS 100 はメタデータ領域 6C からメタデータ 300 を読みだす（ステップ S301）。OS 100 は、論理ドライブアドレスとファイル識別 ID をディスク識別 ID と LBA に正引き変換し（ステップ S302）、ディスク識別 ID に対応する記憶装置に対して命令および LBA を送信する（ステップ S303）。記憶装置は、命令に従ってユーザデータ領域 2004 への書き込みまたはユーザデータ領域 2004 からの読み出し動作を行い、OS 100 に応答またはデータを送信する（ステップ S304）。OS 100 は、記憶装置から応答およびデータを受信し、さらにアプリケーションプログラム 400 に応答およびデータを送信して処理が終了する（ステップ S305）。

【0156】

本実施形態においては、記憶部 2 に関する寿命到達時処理 S205 を行う以前の状態においては、論理ドライブ 4 は記憶部 2 のみから構成されているため、メタデータ 300 における論理ドライブ 4 に対応するディスク識別 ID は全て記憶部 2 の WWN となっている。一方、論理ドライブにおいて RAID アレイを構築する場合のように、寿命到達時処理 S205 を行う以前であっても、メタデータ 300 において一つの論理ドライブアドレスに対して複数のディスク識別 ID が割り当てられていてもよい。

【0157】

ディスクステータス領域 2006 には、OS 100 に対する記憶装置の情報が格納される。図 35 は情報処理装置 111 が起動した時、およびインタフェース 19 に記憶装置が接続された時の制御プログラム 200 の処理手順を示すものである。制御プログラム 20

10

20

30

40

50

0 はインタフェース 19 を経由して記憶装置の Boot Loader 領域 2002 のディスクステータス領域 2006 を読み出し、読み出された値に応じて OS 100 に通知する記憶装置の状態を変更する。ディスクステータス = 0 であるとき (ステップ S311)、制御プログラム 200 は OS 100 に対して記憶装置が初期ディスクであると通知し、OS 100 は記憶装置を初期ディスク状態として認識する (ステップ S312)。記憶装置の工場出荷時や、ACS-3 の F4h Security Erase Unit コマンドや、NVM Express Revision 1.1 の 80h Format NVM コマンドなどをもちいることによって情報処理装置 111 が記憶装置を消去したとき、当該記憶装置のディスクステータス領域 2006 の値はディスクステータス = 0 に変更される。また、情報処理装置 111 が記憶装置をフォーマットした時、当該記憶装置のディスクステータス領域 2006 の値はディスクステータス = 0 に変更される。

10

【0158】

ディスクステータス = 1 であるとき (ステップ S313)、制御プログラム 200 は OS 100 に対して記憶装置が通常状態であると通知し、OS 100 は記憶装置を通常状態として認識する (ステップ S314)。寿命到達時処理 S205 を行う前の記憶部 2 のディスクステータスはディスクステータス = 1 となっている。

【0159】

ディスクステータス = 2 であるとき (ステップ S315)、制御プログラム 200 は OS 100 に対して記憶装置が信頼性劣化状態であると通知し、制御プログラム 200 は記憶装置が信頼性劣化状態として認識し (ステップ S316)、寿命到達時処理 S205 を行う。

20

【0160】

ディスクステータス = 3 であるとき (ステップ S317)、制御プログラム 200 は OS 100 に対して記憶装置がデータ引き継ぎ元としてデータ引き継ぎ作業中であると通知し、OS 100 は記憶装置を保護状態として認識する (ステップ S318)。

【0161】

ディスクステータス = 4 であるとき (ステップ S319)、制御プログラム 200 は OS 100 に対して記憶装置がデータ引き継ぎ先としてデータ引き継ぎ作業中であると通知し、OS 100 は記憶装置をデータ引き継ぎ先状態として認識する (ステップ S320)。

【0162】

ディスクステータス = 5 であるとき (ステップ S321)、制御プログラム 200 は OS 100 に対して記憶装置が廃棄対象状態であると通知し、OS 100 は記憶装置を廃棄対象状態として認識する (ステップ S322)。ディスクステータスが 0 ~ 5 以外であるときは、不正なディスクとみなし、マウントしないことが望ましい (ステップ S323)。

30

【0163】

図 36 および図 37 は記憶装置のライフサイクルを示す状態遷移図である。図 36 に示すように、記憶装置製造直後、記憶装置購入直後、記憶装置消去直後、記憶装置フォーマット直後の記憶装置のディスクステータス領域 2006 の値は初期ディスク状態を示すディスクステータス = 0 となっており (ステップ S330)、その後通常状態を示すディスクステータス = 1 となることによって、OS 100 に通常の記憶装置として利用され (ステップ S331)、記憶装置が信頼性劣化するにつれて信頼性劣化状態を示すディスクステータス = 2 へと遷移し (ステップ S332)、さらに保護状態を示すディスクステータス = 3 を経由した (ステップ S333) 後、廃棄対象状態を示すディスクステータス = 5 となった (ステップ S334) 後、最終的に廃棄される。

40

【0164】

図 37 は、情報処理装置 111 にデータ引き継ぎ先状態としてマウントされた場合の記憶装置のライフサイクルを示すもので、この場合は、ステップ S330 のディスクステータス = 0 の初期ディスク状態の後、データ引き継ぎ先状態 (ディスクステータス = 4) を経由した (ステップ S330b) 後、データ引き継ぎ元状態の記憶装置が脱着されること

50

で、ディスクステータス = 1 の通常状態として利用されている (ステップ S 3 3 1)。

【 0 1 6 5 】

図 3 8 は、記憶部 2 が信頼性劣化状態として認識されたときに、記憶部 2 とは異なる別の記憶部 3 が接続された状態を示すものである。新たに接続された記憶部 3 にも、情報処理装置 1 1 1 によって管理可能なアドレス空間として、L B A 領域 3 0 0 1 が割り当てられている。L B A 領域 3 0 0 1 は、Boot Loader 領域 3 0 0 2 と、メタデータ領域 3 0 0 3 と、ユーザデータ領域 3 0 0 4 からなる。

【 0 1 6 6 】

なお、記憶部 2 が信頼性劣化状態として認識されたとき、初期ディスク状態の記憶装置が情報処理装置 1 1 1 に接続されると該初期ディスク状態の記憶装置をデータ引き継ぎ先状態として認識するのではなく、図 3 9 の破線部のように、すでに情報処理装置 1 1 1 に接続されている通常状態の記憶部 3 をデータ引き継ぎ先状態として認識するようにしてもよい。

【 0 1 6 7 】

(寿命到達時処理)

図 4 0 に寿命到達時処理 S 2 0 5 の際に、制御プログラム 2 0 0 が行う処理のフローチャートを示す。制御プログラム 2 0 0 は、接続されている記憶部 2 が寿命到達するかあるいは寿命到達直前になって寿命到達時処理が開始されると、記憶部 2 の Boot Loader 領域 2 0 0 2 のディスクステータス領域 2 0 0 6 がディスクステータス = 2 (信頼性劣化状態) であるか否かを判定する (ステップ S 3 4 0)。この判定結果が否である場合、制御プログラム 2 0 0 は、ディスクステータス領域 2 0 0 6 をディスクステータス = 2 に書き換え (ステップ S 3 4 1)、手順をステップ S 3 4 2 に移行する。ステップ S 3 4 0 の判定結果が Y e s の場合は、手順をステップ S 3 4 2 に移行する。

【 0 1 6 8 】

ステップ S 3 4 2 のように、制御プログラム 2 0 0 は、「新しいディスクをインタフェース 1 9 に接続してください」などの新しいディスクの接続を促すメッセージをディスプレイ 9 に表示することが望ましい。制御プログラム 2 0 0 は、ディスクステータス = 0 (初期ディスク状態) の記憶装置が接続されたか否かを判定する (ステップ S 3 4 3)。

【 0 1 6 9 】

新たな記憶装置が接続されると、制御プログラム 2 0 0 は、記憶部 2 の Boot Loader 領域 2 0 0 2 を、新たに接続されたディスクステータス = 0 (初期ディスク状態) の記憶部 3 の Boot Loader 領域 3 0 0 2 にコピーする (ステップ S 3 4 4、図 3 2 参照)。制御プログラム 2 0 0 は、記憶部 3 の Boot Loader 領域 3 0 0 2 のディスクステータス領域 3 0 0 6 をディスクステータス = 4 (データ引き継ぎ先状態) に書き換える (ステップ S 3 4 5)。

【 0 1 7 0 】

なお、ステップ S 3 4 0 : Y e s や S 3 4 1 の時点で、すでにディスクステータス = 0 (初期ディスク状態) の記憶装置がすでに接続されている場合や、データ引き継ぎ先として割り当て可能なディスクステータス = 1 (通常状態) の記憶装置がすでに接続されている場合には、制御プログラム 2 0 0 は、この記憶装置を記憶部 3 とし、記憶部 2 の Boot Loader 領域 2 0 0 2 を、記憶部 3 の Boot Loader 領域 3 0 0 2 にコピーし、記憶部 3 の Boot Loader 領域 3 0 0 2 のディスクステータス領域 3 0 0 6 をディスクステータス = 4 に書き換えてもよい。

【 0 1 7 1 】

制御プログラム 2 0 0 は、記憶部 3 の Boot Loader 領域 3 0 0 2 のアドレス領域 3 0 0 7 に記憶部 2 のディスク識別 I D を書き込む (ステップ S 3 4 6)。制御プログラム 2 0 0 は、記憶部 2 の Boot Loader 領域 2 0 0 2 のディスクステータス領域 2 0 0 6 をディスクステータス = 3 (保護状態) に書き換える (ステップ S 3 4 7)。制御プログラム 2 0 0 は、記憶部 2 の Boot Loader 領域 2 0 0 2 のアドレス領域 2 0 0 7 に記憶部 3 のディスク識別 I D を書き込む (ステップ S 3 4 6)。記憶部 2 のアドレス領域 2 0 0 7 および記

10

20

30

40

50

憶部 3 のアドレス領域を見ることで、記憶部 2 と記憶部 3 がデータ引き継ぎ処理のペアであることを認識することができる。

【 0 1 7 2 】

制御プログラム 2 0 0 は、最新のメタデータ 3 0 0 を主メモリ 6 からまたは記憶部 2 のメタデータ格納領域 2 0 0 3 から読み出して、読み出した最新のメタデータを記憶部 3 のメタデータ格納領域 3 0 0 3 に書き込んでコピーする(ステップ S 3 4 9)。制御プログラム 2 0 0 は、記憶部 2 + 記憶部 3 を 1 つの論理ドライブ 4 として O S 1 0 0 に認識させる(ステップ S 3 5 0)。制御プログラム 2 0 0 は、主メモリ 6 上の領域 6 D の論理ドライブステータステーブル 4 5 0 を、ステータスが「通常状態」から「データ引き継ぎ状態」になるように書き換える(ステップ S 3 5 1)。

10

【 0 1 7 3 】

図 4 1 は、主メモリ 6 の領域 6 D に格納されている論理ドライブステータステーブル 4 5 0 を示すものである。テーブル 4 5 0 では、複数の論理ドライブとステータスの対応が管理されている。制御プログラム 2 0 0 は、論理ドライブの状態(通常状態またはデータ引き継ぎ状態)に応じて、テーブル 4 5 0 を随時書き換える。

【 0 1 7 4 】

本実施形態では、寿命到達時処理 S 2 0 5 によって、記憶部 2 と記憶部 3 は、図 3 8 に示したように、単一の論理ドライブである論理ドライブ 4 として認識される。新しい記憶部 3 が接続されてから、記憶部 3 を用いた論理ドライブ 4 が構築されるまでのデータの読み出し、書き込みは、高々、メタデータ 3 0 0 のデータ量しか発生しないため、R A I D ディスクアレイでディスクを交換して論理ドライブにマウントするまでの時間と比べ、記憶部 3 のマウントが非常に高速に行われる。

20

【 0 1 7 5 】

寿命到達時処理 S 2 0 5 以前においては、主メモリ 6 のメタデータ領域 6 C のメタデータ 3 0 0 やメタデータ 3 0 0 のジャーナルは O S 1 0 0 によって定期的にデータ引き継ぎ元の記憶部 2 のメタデータ領域 2 0 0 3 に退避保存されていたが、寿命到達時処理 S 2 0 5 以後においては、主メモリのメタデータ領域 6 C のメタデータ 3 0 0 やメタデータ 3 0 0 のジャーナルは O S 1 0 0 によって定期的にデータ引き継ぎ先の記憶部 3 のメタデータ領域 3 0 0 3 に退避保存される。これにより、メタデータ領域 3 0 0 3 には最新に近いメタデータが格納され、メタデータ領域 2 0 0 3 には寿命到達時処理 S 2 0 5 以前の古いメタデータが格納されることになる。

30

【 0 1 7 6 】

なお、上記の説明では、ディスクステータス = 2 の信頼性劣化状態を定義することにしたが、制御プログラム 2 0 0 は、信頼性情報と閾値との比較の結果、記憶部 2 が寿命に到達したと判定した場合、ディスクステータス = 2 の信頼性劣化状態を経ることなく、ディスクステータスをディスクステータス = 3 のデータ引き継ぎ元状態(保護状態)に即座に移行させるようにしてもよい。すなわち、この場合は、制御プログラム 2 0 0 は、寿命到達時処理が開始されると、図 4 0 のステップ S 3 4 0 において、記憶部 2 の Boot Loader 領域 2 0 0 2 のディスクステータス領域 2 0 0 6 がディスクステータス = 3 (データ引き継ぎ元状態)であるか否かを判定する。この判定結果が否である場合、制御プログラム 2 0 0 は、図 4 0 のステップ S 3 4 1 において、ディスクステータス領域 2 0 0 6 をディスクステータス = 3 に書き換え、手順をステップ S 3 4 2 に移行する。それ以降の手順は、ステップ S 3 4 7 が削除される以外は、図 4 0 と同様である。

40

【 0 1 7 7 】

(論理ドライブへの書き込み)

図 4 2 はアプリケーションプログラム 4 0 0 から O S 1 0 0 にファイルデータ書き込み要求が送信された時の、O S 1 0 0 の処理手順を示すものである。O S 1 0 0 は、アプリケーションプログラム 4 0 0 から書き込み命令と、論理ドライブアドレスと、ファイル I D と、データを受信する(ステップ S 3 6 0)。O S 1 0 0 は、主メモリ 6 から論理ドライブステータステーブル 4 5 0 を読み出し(ステップ S 3 6 1)、主メモリ 6 からメタデ

50

ータ300を読み出し(ステップS362)、メタデータ300を参照してデータ書き込み用のLBAを割り当てる(ステップS363)。

【0178】

OS100は、論理ドライブステータステーブル450に基づいて書き込み命令で指定された論理ドライブが通常状態かデータ引き継ぎ状態であるかを判定する(ステップS364)。論理ドライブが通常状態である場合、OS100は、記憶部2に対し、書き込み命令、LBAおよび書き込みデータを送信する(ステップS365)。OS100は、記憶部2から応答を受信する(ステップS366)。OS100は、主メモリ6上のメタデータを書き換えて、書き込みファイルIDと記憶部2とLBAとセクタカウントを紐付けする(ステップS367)。OS100は、アプリケーションプログラム400に応答を送信する(ステップS371)。

10

【0179】

論理ドライブがデータ引き継ぎ状態である場合、OS100は、データ引き継ぎ先の記憶部3に対し、書き込み命令、LBAおよび書き込みデータを送信する(ステップS368)。OS100は、記憶部3から応答を受信する(ステップS369)。OS100は、主メモリ6上のメタデータを書き換えて、書き込みファイルIDと記憶部3とLBAとセクタカウントを紐付けする(ステップS370)。OS100は、アプリケーションプログラム400に応答を送信する(ステップS371)。すなわち、OS100は、論理ドライブがデータ引き継ぎ状態である場合、データ引き継ぎ先の記憶部3への書き込みに伴って、主メモリ6上のメタデータを書き換えて、記憶部2および記憶部3の記憶データのアドレスを更新する。書き込みに伴う記憶データのアドレスの更新とは、記憶部3への書き込み処理と同時並行的に行うこと、記憶部3への書き込み処理の中で行うこと、記憶部3への書き込み処理を行う前に行うこと、記憶部3への書き込み処理を行った後に行うことを全て含む。

20

【0180】

(論理ドライブへのファイル削除命令)

図43はアプリケーションプログラム400からOS100にファイル削除要求が送信された時の、OS100の処理手順を示すものである。OS100は、アプリケーションプログラム400から削除命令と、論理ドライブアドレスと、ファイルIDを受信する(ステップS900)。OS100は、主メモリ6から論理ドライブステータステーブル450を読み出し(ステップS901)、主メモリ6からメタデータ300を読み出し(ステップS902)、メタデータ300を参照して論理ドライブアドレスとファイルIDをディスク識別IDとLBAに正引き変換する(ステップS903)。OS100は、主メモリ6内のメタデータから削除対象ファイルのファイルIDが含まれる行を削除するか、主メモリ6内のメタデータ上の削除対象ファイルのファイルIDを無効なIDに書き換えるかすることで、削除対象ファイルIDをメタデータ300から削除する(ステップS904)。

30

【0181】

OS100は、論理ドライブステータステーブル450に基づいて削除命令で指定された論理ドライブが通常状態かデータ引き継ぎ状態であるかを判定する(ステップS905)。論理ドライブが通常状態である場合、OS100は、記憶部2に対し、削除通知およびLBAを送信する(ステップS906)。OS100は、記憶部2から応答を受信する。OS100は、アプリケーションプログラム400に応答を送信する(ステップS910)。

40

【0182】

論理ドライブがデータ引き継ぎ状態である場合、OS100は、正引き変換後のディスク識別IDがデータ引き継ぎ元である記憶部2であるかデータ引き継ぎ先の記憶部3であるかを判定する(ステップS907)。正引き変換後のディスク識別IDが記憶部2である場合、OS100は、記憶部2に対し、削除通知およびLBAを送信し(ステップS908)、記憶部2から応答を受信し、アプリケーションプログラム400に応答を送信す

50

る（ステップS 9 1 0）。正引き変換後のディスク識別IDが記憶部3である場合、OS 1 0 0は、記憶部3に対し、削除通知およびLBAを送信し（ステップS 9 0 9）、記憶部3から応答を受信し、アプリケーションプログラム4 0 0に応答を送信する（ステップS 9 1 0）。

【0 1 8 3】

（論理ドライブからの読み出し）

図4 4はアプリケーションプログラム4 0 0からOS 1 0 0にファイルデータ読み出し要求が送信された時の、OS 1 0 0の処理手順を示すものである。OS 1 0 0は、アプリケーションプログラム4 0 0から読み出し命令と、論理ドライブアドレスと、ファイルIDとを受信する（ステップS 3 8 0）。OS 1 0 0は、主メモリ6から論理ドライブステータステーブル4 5 0を読み出し（ステップS 3 8 1）、主メモリ6からメタデータ3 0 0を読み出し（ステップS 3 8 2）、メタデータ3 0 0を参照して、論理ドライブアドレスとファイルIDをデータ読み出し用のディスク識別IDとLBAとセクタカウントに正引き変換する（ステップS 3 8 3）。

【0 1 8 4】

正引き変換後のディスク識別IDが記憶部2を指定している場合（ステップS 3 8 4）、OS 1 0 0は、記憶部2に対し、読み出し命令、LBAおよびセクタカウントを送信する（ステップS 3 8 5）。OS 1 0 0は、記憶部2から応答および読み出しデータを受信する（ステップS 3 8 6）。OS 1 0 0は、アプリケーションプログラム4 0 0に読み出しデータおよび応答を送信する（ステップS 3 8 9）。

【0 1 8 5】

正引き変換後のディスク識別IDが記憶部3を指定している場合（ステップS 3 8 4）、OS 1 0 0は、記憶部3に対し、読み出し命令、LBAおよびセクタカウントを送信する（ステップS 3 8 7）。OS 1 0 0は、記憶部3から応答および読み出しデータを受信する（ステップS 3 8 8）。OS 1 0 0は、アプリケーションプログラム4 0 0に読み出しデータおよび応答を送信する（ステップS 3 8 9）。例えば、LBA = 0がデータ引き継ぎ済み、LBA = 1が未データ引き継ぎの状態、LBA = 0、セクタカウント = 1で読み出しを行うと、記憶部3から読み出しが行われ、LBA = 1、セクタカウント = 1で読み出しを行うと、記憶部2から読み出しが行われ、LBA = 0、セクタカウント = 2で読み出しを行うと、記憶部2および記憶部3から読み出しが行われる。

【0 1 8 6】

このように、データ引き継ぎ元記憶装置への書き込みを禁止し、データ引き継ぎ先の記憶部3への書き込みを利用してデータ引き継ぎ元記憶装置からデータ引き継ぎ先記憶装置へのデータ引き継ぎを実現しているため、ユーザ自身によるバックアップ作業は不要となる。また、データ引き継ぎ中において、ユーザデータ2 0 0 4のコピーそのものは行われず、データ引き継ぎはユーザデータ2 0 0 4の新規書き込み処理を利用して行われるため、データ引き継ぎ中であってもアプリケーションプログラム4 0 0の書き込み処理性能は低下しない。また、寿命到達時処理S 2 0 5後において、記憶部2に発生する書き込み処理はたかだかディスクステータス領域2 0 0 6に対する書き込みのみであるため、記憶部2への書き込み処理はほとんど発生しない。このようにして、記憶部2に対する寿命到達時処理S 2 0 5以後であっても、アプリケーションプログラム4 0 0にとって論理ドライブ4そのものは読み出しおよび書き込み可能なドライブとして認識されるが、情報処理装置1 1 1にとって実際には記憶部2はリードオンリーデバイスのように扱われることになる。

【0 1 8 7】

（書き戻しバックアップ）

なお、データ引き継ぎ元である記憶部2のデータを主メモリ6内のキャッシュメモリ領域に読みだした時には、キャッシュメモリ領域に読みだしたデータをデータ引き継ぎ先である記憶部3に書き込み（書き戻し）、書き込み先LBAとファイルIDを関連付けるようにメタデータ3 0 0を書き換えてもよい。以下、図4 5を用いて説明する。

【 0 1 8 8 】

OS 100は、アプリケーションプログラム400から読み出し命令と、論理ドライブアドレスと、ファイルIDとを受信する(ステップS400)。OS 100は、主メモリ6から論理ドライブステータステーブル450を読み出し(ステップS401)、主メモリ6からメタデータ300を読み出し(ステップS402)、メタデータ300を参照して、論理ドライブアドレスとファイルIDをデータ読み出し用のディスク識別IDとLBAとセクタカウントに正引き変換する(ステップS403)。

【 0 1 8 9 】

正引き変換後のディスク識別IDが記憶部3を指定している場合(ステップS404)、OS 100は、記憶部3に対し、読み出し命令、LBAおよびセクタカウントを送信する(ステップS409)。OS 100は、記憶部3から応答および読み出しデータを受信する(ステップS410)。OS 100は、アプリケーションプログラム400に記憶部3から読み出したデータおよび応答を送信する(ステップS411)。

10

【 0 1 9 0 】

正引き変換後のディスク識別IDが記憶部2を指定している場合(ステップS404)、OS 100は、記憶部2に対し、読み出し命令、LBAおよびセクタカウントを送信する(ステップS405)。さらに、OS 100は、記憶部3に対し、書き込み命令、LBAおよびセクタカウントを送信する(ステップS406)。OS 100は、記憶部2から応答および読み出しデータを受信する(ステップS407)。OS 100は、記憶部3に対し、書き込み命令、LBAおよび記憶部2から読み出したデータを送信し、これにより記憶部2から読み出したデータを記憶部3に書き込むバックグラウンド書き込みを行う(ステップS408)。OS 100は、アプリケーションプログラム400に記憶部2から受信したデータおよび応答を送信する(ステップS412)。OS 100は、主メモリ6上のメタデータを書き換えて、書き込みファイルIDと記憶部3とLBAとセクタカウントを紐付けする(ステップS413)。OS 100は、アプリケーションプログラム400に応答を送信する(ステップS372)。

20

【 0 1 9 1 】

このようにすることで、論理ドライブ4から情報処理装置111へのデータ読み出しのバックグラウンドで、記憶部3へのデータ移行が可能になり、後述するバックグラウンドバックアップするべきLBA領域量を削減し、データ引き継ぎ状態開始から完了までの期間がより短縮化される。特に、論理ドライブ4の読み出し動作において、記憶部2のデータ読み出しと記憶部3へのデータ書き戻しを並列して行うことで、より高速にデータ引き継ぎを行うことができる。

30

【 0 1 9 2 】

(バックグラウンドバックアップ)

論理ドライブ4が論理ドライブステータステーブル450においてデータ引き継ぎ状態である場合、アプリケーションプログラム400やOS 100による論理ドライブ4へのアクセスがほとんど発生しない時に(アイドル時に)、データ引き継ぎ元の記憶部2からデータ引き継ぎ先の記憶部3にバックグラウンドでバックアップを行う(バックグラウンドバックアップ)。制御プログラム200は主メモリ6からメタデータ300を読み出し、記憶部2に紐付けられているファイルIDを検索し、記憶部2に紐付けられているファイルが存在すれば、当該ファイルのLBAに対して記憶部2に読み出し命令を送信し、記憶部3の当該LBAに対して書き込み命令と読み出されたデータを送信して、書き込みを行い、主メモリ6上のメタデータ300を書き換えて当該ファイルIDを記憶部3に紐付ける。

40

【 0 1 9 3 】

(データ引き継ぎ完了時)

図46は、データ引き継ぎ完了時の制御プログラムの動作手順を示すものである。論理ドライブステータステーブル450において、論理ドライブ4のステータスが「データ引き継ぎ状態」の時(ステップS420)、制御プログラム200は、主メモリ6上のメタ

50

データ 300 を定期的に読み出し（ステップ S 4 2 1）、記憶部 2 に紐付けられている引き継ぎ対象ファイル ID が存在するか否かを定期的にチェックする（ステップ S 4 2 2）。たとえば、制御プログラム 200 は、論理ドライブ 4 に格納されている全ファイルのファイル ID のうち、記憶部 2 に紐付けられている引き継ぎ対象ファイル ID が存在するか否かを定期的にチェックする。存在する場合には、まだデータ引き継ぎが完了していないため、データ引き継ぎ状態のステータスを続行する。

【0194】

一方、存在しない場合には、制御プログラム 200 はデータ引き継ぎ先である記憶部 3 のディスクステータス領域 3006 をディスクステータス = 1（通常状態）に書き換え（ステップ S 4 2 3）、データ引き継ぎ元である記憶部 2 の領域 2006 をディスクステータス = 5（廃棄対象状態）に書き換える（ステップ S 4 2 4）。制御プログラム 200 は、記憶部 2 を論理ドライブ 4 から切り離し、記憶部 3 を論理ドライブ 4 として認識（マウント）し（ステップ S 4 2 5）、論理ドライブステータステーブル 450 で論理ドライブ 4 のステータスを「データ引き継ぎ状態」から「通常状態」に書き換える（ステップ S 4 2 6）。

【0195】

これにより、記憶部 2 はいつでも物理的（機械的）に取り外し可能な状態で破棄可能な状態になり、かつ記憶部 3 が寿命到達時処理 S 205 以前の記憶部 2 の役割を担うことになる。以後、記憶部 3 を記憶部 2 とみなしてよく、情報処理装置 111 のデータ構成は寿命到達時処理 S 205 以前のデータ構成である図 31 の状態に戻る。

【0196】

なお、記憶部 2 を情報処理装置 111 から安全に取り外すためには、記憶部 2 を論理ドライブ 4 から切り離し後に、Information technology ATA/ATAPI Command Set-3 (ACS-3) に記載されている E0h Standby Immediate コマンド等を発行したり、NVM Express Revision 1.1 に記載されている Shutdown Notification (CC.SHN) を 01b にセットしたりするなどして、あらかじめ記憶部 2 を電源遮断可能な状態に遷移させておくことが望ましい。

【0197】

また、廃棄可能な状態の記憶部 2 の消費電力を低減するためには、記憶部 2 を論理ドライブ 4 から切り離し後に、Information technology ATA/ATAPI Command Set-3 (ACS-3) に記載されている E0h Standby Immediate コマンドや E6h SLEEP コマンド等を発行したり、記憶部 2 への供給電源を遮断したり、記憶部 2 を Serial ATA Revision 3.1 Gold Revision に記載の Partial 状態や Slumber 状態に遷移させたり、“Serial ATA Technical Prop O S a l: SATA31_TPR_C108 Title: Device Sleep” に記載の DEVSLP 信号をアクティブにして記憶部 2 を DevSleep 状態に遷移させたり、記憶部 2 を PCI Express Base Specification Revision 3.0 に記載の D1 状態や D2 状態や D3 状態に遷移させたり、PCI Express Base Specification Revision 3.0 に記載の L1 状態や L2 状態や L3 状態に遷移させたりしてもよい。

【0198】

図 47 は、記憶部 2 がディスクステータス = 3（保護状態）で、記憶部 3 がディスクステータス = 4（データ引き継ぎ先状態）である状態における情報処理装置 111 による論理ドライブ 4 からの読み出し状態を示すものである。記憶部 2、3 で、LBA は重複していない。この状態の時には、指定される LBA に応じて、記憶部 2、3 の少なくとも一方からデータが読み出される。

【0199】

図 48 は、記憶部 2 がディスクステータス = 3（保護状態）で、記憶部 3 がディスクステータス = 4（データ引き継ぎ先状態）である状態における情報処理装置 111 による論理ドライブ 4 に対する書き込み状態を示すものである。この状態の時には、書き込みは、記憶部 3 に対してのみ行われる。すなわち、記憶部 2 は、リードオンリーデバイスとし機能する。記憶部 3 に書き込みが行われた LBA に関しては、記憶部 2 から割り当てが解除される。

【0200】

以上のようにして、本実施形態により、記憶部 2 が寿命到達した場合または寿命到達直前になった場合、記憶部 2 の書き換えは、ディスクステータス領域 2 0 0 6 の書き換えしか発生しないため、書き込み処理がほとんど行われなくなり、あたかもリードオンリーデバイスであるかのように取り扱われる。一方、論理ドライブは読み出しと書き込み両方可能なドライブとして振る舞うため、アプリケーションプログラム 4 0 0 にとっての論理ドライブ 4 は寿命到達以前と同様の振る舞いになる。記憶部 2 から記憶部 3 へのデータのデータ引き継ぎは、アプリケーションプログラム 4 0 0 や OS 1 0 0 から論理ドライブ 4 への書き込みが要求された場合に発生し、記憶部 2 から記憶部 3 へのデータのコピーではなく、アプリケーションプログラム 4 0 0 や OS 1 0 0 から記憶部 3 へのデータ書き込み処理とメタデータ書き換えによる論理的なデータ移行という形で行われる。それにより、記憶部 2 から記憶部 3 へのデータのデータ引き継ぎは、アプリケーションプログラム 4 0 0 や OS 1 0 0 から記憶部 2 への通常のデータ書き込みのバックグラウンドで実行可能であり、記憶部 2 からデータを読み出し記憶部 3 に当該データを書き込む処理をアプリケーションプログラム 4 0 0 や OS 1 0 0 からの論理ドライブ 4 へのアクセスとは独立して行う比較例のバックアップ処理に比べて飛躍的に高速に行われる。つまり、アプリケーションプログラム 4 0 0 や OS 1 0 0 からの書き込みが発生する L B A に対しては実質的にデータ引き継ぎ時間がゼロとなる。

【 0 2 0 1 】

アプリケーションプログラム 4 0 0 や OS 1 0 0 からの書き込みが発生しないような L B A に対しては、別途バックアップ処理が必要となるが、新しい記憶装置のマウント前にデータをコピーしなければならない比較例のバックアップ処理や R A I D アレイの再構築と異なり、記憶部 2 と記憶部 3 をマウントした後に、アイドル時にバックグラウンドで行うことができるため、アプリケーションプログラム 4 0 0 の性能劣化を抑制することができる。新しい記憶装置のマウント前にユーザデータのコピーが必要となる比較例のバックアップ処理や、新しい記憶装置のマウント前にユーザデータおよびパリティデータの再構築が必要となる R A I D などによる論理ドライブの再構築と異なり、本実施形態によるデータ引き継ぎ先記憶装置接続に伴う論理ドライブ再構築は、図 4 0 に示したように、ディスクステータス領域およびディスク識別 ID 領域の書き換えとメタデータ領域のコピーのみ必要となるため、非常に高速に行うことが可能になる。

【 0 2 0 2 】

(第 2 の実施形態)

第 1 の実施形態では、記憶部 3 に引き継がれたデータを検索するための情報として、主メモリ 6 に格納されるメタデータ 3 0 0 が使用される例について説明した。これにより、たとえば OS 1 0 0 がアプリケーションプログラム 4 0 0 からファイル ID の指定をうけて論理ドライブ 4 のデータ読み出しを命令されたとき、メタデータ 3 0 0 を読み出すことで、記憶部 2 と記憶部 3 のどちらからデータを読み出すべきかの情報と、どの L B A からデータを読み出すべきかの情報を取得することができる。第 2 の実施形態では、記憶部 3 に引き継がれたデータを検索するための情報として、記憶部 3 に格納される引き継ぎ履歴 5 5 0 が使用される例について説明する。たとえば OS 1 0 0 がアプリケーションプログラム 4 0 0 から論理ドライブ 4 のデータ読み出しを命令されたとき、引き継ぎ履歴 5 5 0 を読み出すことで、記憶部 2 と記憶部 3 のどちらからデータを読み出すべきかの情報を取得することができる。なお、本実施形態においては、アプリケーションプログラム 4 0 0 は L B A を直接指定して OS 1 0 0 に読み出し命令、書き込み命令を送信する場合について説明する。アプリケーションプログラム 4 0 0 が第 1 の実施形態のようにファイル ID を指定して OS 1 0 0 に読み出し命令・書き込み命令を送信する場合についても本実施形態の発明は適用可能であり、その場合、制御プログラム 2 0 0 または OS 1 0 0 がメタデータ 3 0 0 を読みだすことでファイル ID を L B A に変換することができ、変換した L B A についての読み出し、書き込み処理は本実施形態と同様に行われる。

【 0 2 0 3 】

図 4 9 に第 2 の実施形態の情報処理システム 1 の構成を示す。情報処理システム 1 の基

10

20

30

40

50

本構成は第1の実施形態1と同様である。論理ドライブは、OS100が認識しうる論理的なドライブであり、論理ドライブID（ドライブ名やボリューム番号や論理ユニット番号など）が割り当てられ、OS100は1乃至複数の物理デバイスとしての記憶部を論理ドライブとして認識する。論理ドライブは、論理セクタ（論理ブロック）に分割され、各論理セクタには、LBAが割り当てられる。論理ドライブは、図7におけるOS100および制御プログラム200によって、下位の階層である1乃至複数台の物理ドライブやあるいはそれらの一部のLBA領域に割り当てられ、OS100は論理ドライブのLBAと物理ドライブのLBAを相互変換する。上位の階層であるアプリケーションプログラム400は論理ドライブを仮想的に1台のドライブとして認識する。本実施形態においては、寿命到達時処理S205以前の状態において、論理ドライブ4が1つの物理記憶装置である記憶部2に割り当てられる場合について述べる。この場合、論理ドライブのLBAと物理ドライブのLBAは同じ値となる。なお、寿命到達時処理S205以前の状態であっても、複数の物理記憶装置を用いてRedundant Arrays of Inexpensive Disks (RAID)のディスクアレイ、たとえばRAID0やRAID5のディスクアレイを構成し、これを一つの論理ドライブ4として認識してもよく、その場合であっても本実施形態は適用可能である。アプリケーションプログラム400は、論理ドライブID及びLBAからなる論理アドレスを含むコマンドをOS100に与えることにより、特定の論理ドライブにおける特定の論理セクタに対してアクセスすることができる。なお、論理ドライブIDは記憶部の全LBA領域ではなく一部のLBA領域に割り当てられても良いし、そうすることで、記憶部2および記憶部3は複数の論理ドライブに分割して管理されそれぞれの論理ドライブに別々の論理ドライブIDが割り当てられてもよい。

10

20

【0204】

本実施形態は、例として、記憶部2として第1の実施形態に記載の記憶部2であるSSDを用い、記憶部3として第1の実施形態に記載の記憶部3であるSSDを用いる場合について説明する。信頼性劣化後に廃棄して設置スペースを削減しかつシステム1全体の消費電力を削減するという観点では、記憶部2は情報処理装置111に対して物理的に着脱可能であることが望ましい。

【0205】

データ引き継ぎ先である記憶部3の記憶可能容量は、データ引き継ぎ元である記憶部2の記憶可能容量と同じか、記憶部2の記憶可能容量よりも大きいことが望ましいが、記憶部2の記憶可能容量よりも小さい場合であっても本発明は適用可能である。

30

【0206】

本実施形態では、記憶部3は記憶部2が寿命到達または寿命到達目前であると判定された後に、新たに情報処理装置111に接続される記憶部である。記憶部2が寿命到達または寿命到達目前であると判定された後に、すでに情報処理装置111に接続されている通常状態の記憶部3を新たな接続無しに引き継ぎ先として利用する場合についても本発明は適用可能である。記憶部3の接続前の設置スペースを削減しかつシステム1全体の消費電力を削減するという観点、および記憶部3の信頼性劣化後に記憶部3を廃棄して設置スペースを削減しかつシステム1全体の消費電力を削減するという観点では、記憶部3は情報処理装置111に対して物理的に着脱可能であることが望ましい。

40

【0207】

主メモリ6に格納される制御プログラム200は、記憶部2や記憶部3の統計情報、それぞれのステータス記憶領域510、論理ドライブID記憶領域520、および引き継ぎ履歴記憶領域550の制御および管理を行い、統計情報に基づく寿命到達時処理やデータ引き継ぎ処理等を行う。

【0208】

記憶部2および記憶部3はそれぞれステータス格納領域510と論理ドライブID格納領域520を具備し、記憶部3は引き継ぎ履歴格納領域550を具備する。

【0209】

本実施形態では、ステータス記憶領域510に格納されるデータは、0乃至5の値を取

50

りうる。それぞれの値は、対応する記憶部がそれぞれ

- 0：初期ディスク状態
- 1：通常状態
- 2：信頼性劣化状態
- 3：データ引き継ぎ元状態（保護状態）
- 4：データ引き継ぎ先状態
- 5：廃棄対象状態

の状態であることを示す。なお、ステータス記憶領域 510 や論理ドライブ ID 記憶領域 520 や引き継ぎ履歴記憶領域 550 は各記憶部内に格納されるのではなく、主メモリ 6 に格納されるようにしてもよい。データ引き継ぎ元の記憶部 2 とデータ引き継ぎ先の記憶部 3 の論理ドライブ ID 記憶領域 520 には、同じ論理ドライブ ID が格納される。

【0210】

引き継ぎ履歴格納領域 550 は、データ引き継ぎ先である記憶部 3 に具備され、データ引き継ぎ状態の記憶部 3 の LBA に対して書き込みが行われたり、情報処理装置 1 から論理ドライブ 4 に削除通知が発行されて記憶部 2 の LBA のデータが無効化されたりすると、その更新履歴（データ引き継ぎ履歴）として、図 50 に示すように、書き込み LBA と書き込みセクタサイズを格納する。引き継ぎ履歴格納領域 550 は書き込み時や削除通知時や記憶部 3 のアイドル時に随時整理・最適化される。たとえば、 $LBA = cLBA \sim cLBA + X - 1$ のセクタサイズ X の領域（以後 $(cLBA, X)$ と表記する）に対し、オーバーラップする LBA や連続する LBA が引き継ぎ履歴 550 にすでに記録されている場合には、それらを合併（統合）したログを引き継ぎ履歴 550 に記録し、引き継ぎ履歴 550 から合併されたログを削除することが望ましい。たとえば、 $A = (cLBA, X)$ の LBA 領域 A を新たに引き継ぎ履歴格納領域 550 に記入する場合であって、 $B = (cLBA - a, a)$ の LBA 領域 B がすでに引き継ぎ履歴 550 に格納されている時、LBA 領域 A と LBA 領域 B は連続している LBA 領域であるため、 $A + B$ の LBA 領域である $(cLBA - a, a + X)$ のデータを引き継ぎ履歴 550 における $B = (cLBA - a, a)$ の履歴格納部に上書きすることで、ログデータ量を増大させることなく引き継ぎ履歴 550 を更新することができる。引き継ぎ履歴 550 は後述するデータ引き継ぎが完了すると消去されることにより、ユーザデータ格納など他の用途に使用されてもよい。

【0211】

図 51 は引き継ぎ履歴 550 に対する履歴の書き込み例である。寿命到達前の論理ドライブ 4 の $LBA = cLBA$ のデータ $D(cLBA)$ は記憶部 2 の $LBA = cLBA$ に格納されている。記憶部 2 が寿命到達すると、たとえば論理ドライブ 4 の $LBA = cLBA$ への 1 セクタサイズのデータ $D(cLBA)$ の書き込みは記憶部 3 の $LBA = cLBA$ に対して行われ、引き継ぎ履歴 550 に $LBA = cLBA$ とセクタカウント = 1 が格納される。寿命到達後の論理ドライブ 4 の $LBA = cLBA$ へのセクタサイズ = X のデータ $D(cLBA), D(cLBA + 1), \dots, D(cLBA + X - 1)$ の書き込みは、記憶部 3 の $LBA = cLBA, cLBA + 1, \dots, cLBA + X - 1$ に対して行われ、引き継ぎ履歴 550 に $LBA = cLBA$ とセクタカウント = X が格納される。

【0212】

図 51 では、 $LBA = 1$ 、 $LBA = 3$ 、 $LBA = 4$ に対する書き込みは記憶部 3 の $LBA = 1$ 、 $LBA = 3$ 、 $LBA = 4$ に行われ、引き継ぎ履歴として、 $LBA = 1$ およびセクタカウント = 1 の履歴が記録され、 $LBA = 3$ およびセクタカウント = 1 と、 $LBA = 4$ およびセクタカウント = 1 の履歴は $LBA = 3$ およびセクタカウント = 2 の履歴に統合されて記録される。

【0213】

ステータス格納領域 510 と論理ドライブ ID 格納領域 520 と引き継ぎ履歴格納領域 550 には LBA の領域が割り当てられてもよい。あるいは、ステータス格納領域 510 と論理ドライブ ID 格納領域 520 と引き継ぎ履歴格納領域 550 には LBA 領域ではなくログページアドレスが割り当てられたログ領域が割り当てられてもよく、この場合には

、たとえば、非特許文献 1 の ACS-3 に記載されている 2Fh Read Log Ext により該ログページ領域の読み出しが行われ、非特許文献 1 の ACS-3 に記載されている 3Fh Write Log Ext により該領域への書き込みが行われる。

【 0 2 1 4 】

制御プログラム 200 は、第 1 の実施形態の図 24 や図 29 と同様にして、CPU5 に接続されたそれぞれの記憶部が寿命に到達しているか否か、寿命到達目前であるか否か、あるいは故障目前であるか否かを判定し、寿命に到達しているか場合、寿命到達目前である場合、あるいは故障目前である場合に、該記憶部の寿命到達時処理を行う。寿命到達判定は、第 1 の実施形態と同様に、図 24 のような一定時間ごと、または一定処理数おき、または一定データ送受信おき、または図 29 や図 30 のように記憶部からのコマンド応答がエラーであった場合に実施される。

10

【 0 2 1 5 】

(寿命到達処理)

図 52 は、たとえば、記憶部 2 が寿命に到達していると判定された場合の、本実施例における記憶部 2 の寿命到達時処理を示す。記憶部 2 が寿命に到達していると判定された場合 (ステップ S430)、制御プログラム 200 は記憶部 2 のステータスを 1 の通常状態から 2 の信頼性劣化状態に書き換える (ステップ S431)。制御プログラム 200 はインタフェース 19 の空きポートに新しい記憶部である記憶部 3 を接続するようディスプレイ装置やポート近傍に設置された LED 等を経由してユーザや管理者に通知することが望ましい (ステップ S432)。あるいは、ディスクロード・アンロード装置 (図示せず) として記憶部 2 や記憶部 3 のインタフェース 19 に対する物理的な着脱を自動的に行う機械装置がシステム 1 に搭載されている場合、新しい記憶部である記憶部 3 をインタフェース 19 に接続するよう該ディスクロード・アンロード装置に対して命令してもよい。

20

【 0 2 1 6 】

新たな記憶部として記憶部 3 が接続された時 (ステップ S433)、制御プログラム 200 は記憶部 3 のステータス 510 を 4 のデータ引き継ぎ先状態に書き換え (ステップ S434)、記憶部 3 の論理ドライブ ID 格納領域 520 に記憶部 2 の論理ドライブ ID 格納領域 520 の値をコピーすることで、記憶部 2、3 の論理ドライブ ID を一致させる (ステップ S435)。今回の例では、記憶部 2 は図 49 に示すように論理ドライブ 4 として割り当てられているため、記憶部 3 の論理ドライブ ID 格納領域 520 には論理ドライブ 4 の ID が書き込まれる。制御プログラム 200 は記憶部 2 のステータス 510 を 3 のデータ保護状態 (データ引き継ぎ元状態) に書き換え (ステップ S436)、OS100 に記憶部 2 と記憶部 3 を同一の論理ドライブである論理ドライブ 4 として認識させる (ステップ S437)。寿命到達処理後、記憶部 2 のステータスは 3 に、記憶部 3 のステータスは 4 になっており、論理ドライブ 4 および記憶部 2 および記憶部 3 はデータ引き継ぎ状態に遷移していることになる。

30

【 0 2 1 7 】

本実施形態では、各記憶部のステータスは、各記憶部のステータス格納領域 510 に不揮発に格納されており、図 35 に示すように、OS100 の起動ごとに制御プログラム 200 はステータス格納領域 510 を読み出すことで各記憶部のステータスを認識する。制御プログラム 200 は、各記憶部のステータスを認識し、かつ各記憶部の論理ドライブ ID を論理ドライブ ID 格納領域 520 から読み出すことで、論理ドライブ 4 がデータ引き継ぎ状態か否かを認識する。

40

【 0 2 1 8 】

(論理ドライブからの読み出し)

制御プログラム 200 は、アプリケーションプログラム 400 からの読み出し命令に対し、図 53 に示すようにしてデータを読み出す。制御プログラム 200 は、アプリケーションプログラム 400 から読み出し命令と、読み出し対象論理ドライブ ID と、読み出し LBA と、セクタカウントを受信する (ステップ S440)。制御プログラム 200 は、論理ドライブ ID 格納部 520 の値が、読み出し対象論理ドライブ ID に等しい記憶部を

50

全て検索して、記憶部 2 や記憶部 3 を特定する（ステップ S 4 4 1）。制御プログラム 2 0 0 は、検索した記憶部のステータス格納部 5 1 0 の値を読み出し、ステータスを判定することで、検索した記憶部それぞれが記憶部 2 と記憶部 3 のいずれであるかを特定する（ステップ S 4 4 2）。ステータス格納部 5 1 0 の読み出しにともなう性能劣化を抑制するためには、記憶部 2 や記憶部 3 のステータス格納部 5 1 0 の各データは情報処理装置 1 1 1 起動時に主メモリ 6 に読み出され、以後ステータス格納部 5 1 0 のデータは主メモリ 6 から読み出されるようにすることが望ましい。

【 0 2 1 9 】

対象論理ドライブに属する記憶部のステータスが 1 の場合、論理ドライブのステータスは通常状態であるということになる（ステップ S 4 4 3：Y e s）。制御プログラム 2 0 0 は、記憶部 2 に、読み出し命令と、読み出し L B A と、セクタカウントを送信する（ステップ S 4 4 4）。制御プログラム 2 0 0 は、記憶部 2 から応答および読み出しデータを受信する（ステップ S 4 4 5）。制御プログラム 2 0 0 は、アプリケーションプログラム 4 0 0 に読み出しデータおよび応答を送信する（ステップ S 4 4 6）。

【 0 2 2 0 】

対象論理ドライブに属する記憶部のステータスが 1 ではない場合、論理ドライブのステータスはデータ引き継ぎ状態ということになる（ステップ S 4 4 3：N o）。制御プログラム 2 0 0 は、記憶部 3 の引き継ぎ履歴 5 5 0 を読み出し（ステップ S 4 4 7）、読み出し L B A が引き継ぎ履歴に含まれるか否かを判定する（ステップ S 4 4 8）。制御プログラム 2 0 0 は、読み出し L B A が引き継ぎ履歴に含まれる場合は（ステップ S 4 4 8：Y e s）、記憶部 3 に読み出し命令と、読み出し L B A と、セクタカウントを送信する（ステップ S 4 5 2）。制御プログラム 2 0 0 は、記憶部 3 から応答および読み出しデータを受信する（ステップ S 4 5 3）。制御プログラム 2 0 0 は、アプリケーションプログラム 4 0 0 に読み出しデータおよび応答を送信する（ステップ S 4 5 5）。

【 0 2 2 1 】

制御プログラム 2 0 0 は、読み出し L B A に引き継ぎ履歴が含まれない場合は（ステップ S 4 4 8：N o）、記憶部 2 に読み出し命令と、読み出し L B A と、セクタカウントを送信する（ステップ S 4 4 9）。制御プログラム 2 0 0 は、記憶部 2 から応答および読み出しデータを受信する（ステップ S 4 5 0）。制御プログラム 2 0 0 は、アプリケーションプログラム 4 0 0 に読み出しデータおよび応答を送信する（ステップ S 4 5 5）。

【 0 2 2 2 】

（書き戻しバックアップ）

なお、例えば図 5 3 において、ステップ S 4 5 1 の書き戻しバックアップを行うようにしてもよいし、行わないようにしてもよい。ステップ S 4 5 1 では、データ引き継ぎ中のデータ引き継ぎ元である記憶部 2 のデータを主メモリ 6 内のキャッシュメモリ領域に読みだした時には、キャッシュメモリ領域に読みだしたデータをデータ引き継ぎ中のデータ引き継ぎ先である記憶部 3 に書き込み、書き込み先 L B A とセクタサイズを引き継ぎ履歴 5 5 0 に書きこむ。そうすることで、論理ドライブ 4 から情報処理装置 1 1 1 へのデータ読み出しのバックグラウンドで記憶部 3 へのデータ移行が可能になり、バックグラウンドバックアップすべき L B A 領域量を削減し、データ引き継ぎ状態開始から完了までの期間がより短縮化される。特に、論理ドライブ 4 の読み出し動作において、記憶部 2 のデータ読み出しと記憶部 3 へのデータ書き戻しを並列して行うことで、より高速にデータ引き継ぎを行うことができる。

【 0 2 2 3 】

（論理ドライブへの削除命令）

図 5 4 はアプリケーションプログラム 4 0 0 から O S 1 0 0 に L B A データ削除要求が送信された時の、O S 1 0 0 の処理手順を示すものである。O S 1 0 0 は、アプリケーションプログラム 4 0 0 から削除命令と、論理ドライブアドレスと、削除対象 L B A を受信する（ステップ S 9 2 0）。制御プログラム 2 0 0 は、論理ドライブ I D 格納部 5 2 0 の値が、L B A データ削除対象の論理ドライブ I D に等しい記憶部を全て検索して、記憶部

2や記憶部3を特定する(ステップS921)。制御プログラム200は、検索した記憶部のステータス格納部510の値を読み出し(ステップS922)、ステータスを判定することで、検索した記憶部それぞれが記憶部2と記憶部3のいずれであるかを特定する。

【0224】

対象論理ドライブに属する記憶部のステータスが1の場合(ステップS923: Yes)、論理ドライブのステータスは通常状態であるということになる。OS100は、記憶部2に対し、削除通知およびLBAを送信する(ステップS924)。OS100は、記憶部2から応答を受信する。OS100は、アプリケーションプログラム400に応答を送信する(ステップS932)。

【0225】

対象論理ドライブに属する記憶部のステータスが1ではない場合(ステップS923: No)、論理ドライブのステータスはデータ引き継ぎ状態ということになる。制御プログラム200は、記憶部3の引き継ぎ履歴550を読み出し(ステップS925)、データ削除対象LBAが引き継ぎ履歴に含まれるか否かを判定する(ステップS926)。削除LBAが引き継ぎ履歴に含まれる場合は、削除対象のデータは記憶部3に格納されていることになる。制御プログラム200は、記憶部3に対し、削除通知およびLBAを送信し(ステップS927)、記憶部3は削除通知対象LBAのデータを無効化し、記憶部3から応答を受信し、アプリケーションプログラム400に応答を送信する(ステップS930)。

【0226】

読み出しLBAに引き継ぎ履歴が含まれない場合は(ステップS926)、削除対象のデータは記憶部2に格納されていることになる。制御プログラム200は記憶部2に対し削除通知およびLBAを送信し(ステップS928)、記憶部2から応答を受信する。記憶部2に対する削除通知は送信しなくても良い。アプリケーションプログラム400からの削除命令の対象となるLBAはアプリケーションプログラム400やOS100にとって今後読み出す必要の無いデータであり、記憶部3に引き継ぐ必要のないデータであるため、制御プログラム200は、引き継ぎ履歴550に削除対象LBAとセクタカウントを記録することにより、削除対象LBAと記憶部2の紐付けを無効にする(ステップS929)。制御プログラム200はアプリケーションプログラム400に応答を送信する(ステップS930)。

【0227】

このようにして、データ削除命令によって、引き継ぎ履歴550の更新により、記憶部2から記憶部3にデータが論理的に引き継がれていくことになり、データ削除がデータ引き継ぎ動作を兼用することになる。

【0228】

(論理ドライブへの書き込み)

制御プログラム200は、アプリケーションプログラム400からの書き込み命令に対し、図55に示すようにしてデータを書き込む。制御プログラム200は、アプリケーションプログラム400から書き込み命令と、書き込み対象論理ドライブIDと、書き込みLBAと、セクタカウントを受信する(ステップS460)。制御プログラム200は、論理ドライブID格納部520の値が、書き込み対象論理ドライブIDに等しい記憶部を全て検索して、記憶部2や記憶部3を特定する(ステップS461)。制御プログラム200は、検索した記憶部のステータス格納部510の値を読み出し、ステータスを判定することで、検索した記憶部それぞれが記憶部2と記憶部3のいずれであるかを特定する(ステップS462)。

【0229】

対象論理ドライブに属する記憶部のステータスが通常状態である場合、論理ドライブのステータスは通常状態であるということになる(ステップS463: Yes)。制御プログラム200は、記憶部2に、書き込み命令と、書き込みLBAと、セクタカウントを送信する(ステップS464)。また、制御プログラム200は、アプリケーションから受

10

20

30

40

50

信した書き込みデータを記憶部 2 に送信する（ステップ S 4 6 5）。

【 0 2 3 0 】

対象論理ドライブに属する記憶部のステータスがデータ引き継ぎ状態である場合、論理ドライブのステータスはデータ引き継ぎ状態ということになる（ステップ S 4 4 3 : N o）。制御プログラム 2 0 0 は、データ引き継ぎ先の記憶部 3 に書き込み命令と、書き込み L B A と、セクタカウントを送信する（ステップ S 4 6 6）。また、制御プログラム 2 0 0 は、アプリケーションから受信した書き込みデータを記憶部 3 に送信する（ステップ S 4 6 7）。制御プログラム 2 0 0 は、記憶部 3 の引き継ぎ履歴 5 5 0 を読み出し（ステップ S 4 6 8）、書き込み L B A が引き継ぎ履歴に含まれるか否かを判定する（ステップ S 4 6 9）。書き込み L B A が引き継ぎ履歴に含まれる場合は、書き込み L B A はすでに引き継ぎ済みということになるので、制御プログラム 2 0 0 は引き継ぎ履歴 5 5 0 を特に更新しない。書き込み L B A が引き継ぎ履歴に含まれない場合は、書き込み L B A は新たに引き継ぎ完了した L B A ということになるので、制御プログラム 2 0 0 は引き継ぎ履歴 5 5 0 に書き込み L B A とセクタカウントを記録する（ステップ S 4 7 0）。記憶部 3 では、書き込み L B A に対し、書き込みデータを書き込む。

10

【 0 2 3 1 】

このようにして、データ引き継ぎ中には、データ引き継ぎ元の記憶部 2 に対して書き込み命令を送信せず、データ引き継ぎ先の記憶部 3 に対して書き込み命令を送信し、記憶部 3 の引き継ぎ履歴 5 5 0 に引き継ぎ履歴を記録する。論理ドライブ 4 において、アプリケーションプログラム 4 0 0 から書き込み要求を受信することにより、記憶部 2 に格納される有効なデータが記憶部 3 に徐々に移行されていくことになり、新規データ書き込みがデータ引き継ぎ動作を兼用することになる。

20

【 0 2 3 2 】

アプリケーションプログラム 4 0 0 からのデータ書き込み命令が記憶部 2 の全 L B A に対して均等な確率分布で送信される場合、十分な量のデータが書き込まれると、記憶部 2 のほぼ全ての有効データが記憶部 3 に移行され、記憶部 2 内の有効データはほとんどなくなる。たとえば記憶部 2 の記憶容量である全論理セクタ数が C 2 で、記憶部 3 の全論理セクタ数が C 3 とし、 $C 2 = C 3 = C$ である場合を考える。モデルケースとして、全 L B A に対する書き込み分布が均一な確率分布である場合を考えると、ある書き込みコマンドである L B A = c L B A が書き込まれる確率は $1 / C$ である。n 回の書き込み命令が処理された場合、L B A = c L B A が一度も書き込まれない確率は $(1 - (1 / C))^n$ となる。 n は n 乗である。よって、n 回の書き込み命令が処理されたあとに書き込みが完了した論理セクタ数の期待値は $C - C \times (1 - (1 / C))^n$ となる。

30

【 0 2 3 3 】

1 回の書き込み命令で 1 論理セクタの書き込みが行われる場合には、記憶部 2 と記憶部 3 の記憶容量の N 倍の容量のデータが書き込まれたとき、処理された書き込みコマンド数は $n = N C$ であるため、書き込みが行われていない論理セクタ数の期待値 E は $E = C \times (1 - (1 / C)^{(N C)})$ となる。たとえば Gbyte 単位での記憶部 2 の記憶容量 G が IDEMA (International Disk Drive Equipment and Materials Association) 規格に基づき $G = 5 1 2 \text{ GB} (= 4 7 6 . 9 \text{ GiByte})$ であるとき、 $C = 9 7 , 6 9 6 , 3 6 8 + 1 , 9 5 3 , 5 0 4 \times (G = 5 1 2 - 5 0) = 1 , 0 0 0 , 2 1 5 , 2 1 6$ であり、一般に C は十分に大きい整数であるため、 $E = C \times e^{(- N)}$ と近似可能である（e は自然対数の底）。よって、期待値 E は N の増加に対して指数関数的に減少する。たとえば、 $G = 5 1 2 \text{ GByte} (= 4 7 6 . 9 \text{ GiByte})$ の容量の論理ドライブ 4 に対して、1 周分である 4 7 6 . 9 GiByte のデータ書き込みが発生した場合、論理ドライブ 4 の約 6 3 . 2 % の L B A が書き込み完了したことになり、論理ドライブ 4 の半数以上のデータが記憶部 2 から記憶部 3 に移行完了したといえる。たとえば、 $G = 5 1 2 \text{ GByte} (= 4 7 6 . 9 \text{ GiByte})$ の容量の論理ドライブ 4 に対して、4 . 6 周分である 1 3 TiByte のデータ書き込みが発生した場合、論理ドライブ 4 の約 9 9 % の L B A が書き込み完了したことになり、論理ドライブ 4 のほぼ全てのデータが記憶部 2 から記憶部 3 に移行完了したといえる。 $K_i = 1 0 2 4 , M_i = 1 0 2 4 \times K_i , G_i = 1 0 2$

40

50

4 × Mi, Ti = 1 0 2 4 × Giである。

【 0 2 3 4 】

(データ引き継ぎ状況の監視)

制御プログラム 2 0 0 は引き継ぎ履歴 5 5 0 を読み出すことでデータ引き継ぎ状態の論理ドライブ 4 のデータ引き継ぎ状況を監視する。図 5 6 は、引き継ぎ履歴を用いたデータ引き継ぎ状況の監視手順を示すものである。たとえば、制御プログラム 2 0 0 は所定時間経過毎に引き継ぎ履歴 5 5 0 を読み出すことでデータ引き継ぎ状況を監視する (ステップ S 4 8 0、S 4 8 1)。制御プログラム 2 0 0 は、全ての引き継ぎ対象 L B A が引き継ぎ履歴 5 5 0 に含まれる場合、データ引き継ぎが完了したと判定する。たとえば、データ引き継ぎ元である記憶部 2 の全 L B A が引き継ぎ履歴 5 5 0 に含まれる場合、データ引き継ぎが完了したと判定される (ステップ S 4 8 2)。あるいは、データ引き継ぎ完了の判定として、例えば、データ引き継ぎ先である記憶部 3 の全 L B A が引き継ぎ履歴 5 5 0 に含まれるか否かを判定するようにしてもよい。

10

【 0 2 3 5 】

データ引き継ぎが完了したと判定された場合、制御プログラム 2 0 0 はデータ引き継ぎ元である記憶部 2 のステータスを 5 の廃棄対象状態に変更し、データ引き継ぎ先である記憶部 3 のステータスを 1 の通常状態に変更することで (ステップ S 4 8 3)、論理ドライブ 4 のデータ引き継ぎ状態を終了し、論理ドライブ 4 のデータ引き継ぎ状況監視を終了する。記憶部 2 が消費する電力を低減するという観点では、制御プログラム 2 0 0 は記憶部 2 に低消費電力モードへの遷移要求を送信することが望ましい。データ引き継ぎ状態終了後、制御プログラム 2 0 0 はディスプレイ 9 やポート近傍に設置された L E D を通じて、記憶部 2 をインタフェース 1 9 から取り外すようユーザや管理者に通知したり、記憶部 2 をインタフェース 1 9 から取り外すよう前記ディスクロード・アンロード装置に命令を送信することが望ましい。

20

【 0 2 3 6 】

(全体のステータスの遷移)

図 5 7 は、記憶部 2 が寿命到達する際の論理ドライブ 4 のステータスの遷移を示したものである。記憶部 2 が通常状態である場合には記憶部 2 のステータスは 1 になっており、記憶部 3 は C P U 5 に未接続状態である (ステップ 1)。制御プログラム 2 0 0 は記憶部 2 が寿命到達したと判定すると、制御プログラム 2 0 0 は記憶部 2 のステータスを 2 に変更する (ステップ 2)。制御プログラム 2 0 0 からの通知または命令にもとづき、新しい記憶部として記憶部 3 がインタフェース 1 9 に接続される (ステップ 3)。記憶部 3 が接続されると、制御プログラム 2 0 0 はデータ引き継ぎ元として記憶部 2 のステータスを 3 に変更し、データ引き継ぎ先として記憶部 3 のステータスを 4 に変更することで、論理ドライブ 4 はデータ引き継ぎ状態に移行完了する (ステップ 4)。引き継ぎ履歴 5 5 0 の情報をもとに、記憶部 2 のすべての有効データが記憶部 3 にデータ引き継ぎされたと判定されると、制御プログラム 2 0 0 は、記憶部 2 のステータスを 5 に変更し、記憶部 3 のステータスを 1 に変更する。以後、記憶部 3 が、元の記憶部 2 であるかのように振る舞う (ステップ 1 に戻る)。その後さらに記憶部 3 = 新記憶部が寿命到達したと判定された場合、同様のステップ 2 ~ 5 が繰り返される。このようにして、情報処理装置 1 1 1 の記憶部のいずれか一つが寿命到達した場合、寿命到達目前である場合、故障寸前である場合となっても、本実施形態により記憶部 2 のデータを新しい記憶部に容易にデータ引き継ぎすることが可能になる。

30

40

【 0 2 3 7 】

(バックグラウンドバックアップ)

制御プログラム 2 0 0 は、たとえば、アプリケーションプログラム 4 0 0 から命令を一定時間以上受信しない場合などのアイドル時、またはスタンバイモード遷移要求をアプリケーションプログラム 4 0 0 から受信した時、または情報処理装置 1 1 1 や O S 1 0 0 のシャットダウン時などに、記憶部 3 へのデータ引き継ぎ未完了の L B A に対して記憶部 2 からデータを自動的に読み出し記憶部 3 へデータを自動的に書き込むバックアップ動作を

50

行う、バックグラウンドバックアップを行うことが望ましい。たとえばバックグラウンドバックアップ動作は、制御プログラム200が、記憶部3の引き継ぎ履歴550を読み出し、引き継ぎ履歴550に含まれないcLBAに対して、記憶部2からデータ読み出しを行い、記憶部3に該データを書き込み、引き継ぎ履歴550に該cLBAと書き込みデータのセクタサイズを格納することで実行される。引き継ぎ履歴550へのcLBAの格納の際、連続したLBA領域が存在する場合、またはオーバーラップしたLBA領域が存在する場合には、これらを統合したLBA領域が引き継ぎ履歴550に格納され、統合前のLBA領域は引き継ぎ履歴550から削除されることが望ましい。

【0238】

引き継ぎ履歴550のデータサイズを低減するという観点、および後述するデータ引き継ぎ終了判定を高速に行うという観点では、引き継ぎ履歴550に登録されたcLBA領域以外のcLBA領域のうち小さく断片化されたcLBA領域（断片化した未データ引き継ぎcLBA領域）に対してバックグラウンドバックアップが優先的に行われることが望ましい。断片化した未データ引き継ぎcLBA領域をバックアップすることにより、断片化未データ引き継ぎ領域の前記断片化未データ引き継ぎ領域の前後の既データ引き継ぎcLBA領域と、新たにデータ引き継ぎされたcLBA領域が引き継ぎ履歴550に統合されて格納される。例えば、領域LBA=0（（LBA=0，セクタサイズ=1））および領域LBA=2（（LBA=2，セクタサイズ=1））のデータがバックアップ済みで引き継ぎ履歴に登録されている場合、LBA1のデータをバックグラウンドバックアップすることで、連続したLBA0からLBA2のLBA領域（LBA=0，セクタサイズ=3）が引き継ぎ完了したことになり、引き継ぎ履歴550の情報量が削減される。

【0239】

たとえば、前述のように、512GByte（=476.9GiByte）のSSDに対してアプリケーションプログラム400から合計476.9GiByteのデータが書き込まれた時に、制御プログラム200が記憶容量の36.8%の容量である175.5GiBの領域に対しバックグラウンドバックアップ動作を行うことで、記憶部2の全LBA領域がデータ引き継ぎ完了したことになる。SSDである記憶部2と、SSDである記憶部3の典型的な読み出し速度と書き込み速度はたとえば、約400MiB/秒であり、前述の476.9GiBのLBA領域の記憶部2からの読み出しは約449秒で完了し、記憶部3への書き込みは約449秒で完了する。よって、このような状況下では、前記バックグラウンドバックアップは高々約15分で完了し、記憶部2からの読み出しと記憶部3への書き込みが並列に行われる場合には約8分で完了する。

【0240】

さらには、たとえば、前述のように、512GByte（=476.9GiByte）のSSDに対してアプリケーションプログラム400から合計13TiByteのデータが書き込まれた時に、制御プログラム200が記憶容量の1%の容量である4.8GiBの領域に対しバックグラウンドバックアップ動作を行うことで、記憶部2の全LBA領域がデータ引き継ぎ完了したことになる。SSDである記憶部2と、SSDである記憶部3の典型的な読み出し速度と書き込み速度はたとえば、約400MiB/秒であり、前述の4.8GiBのLBA領域の記憶部2からの読み出しは約12秒で完了し、記憶部3への書き込みは約12秒で完了する。よって、このような状況下では、前記バックグラウンドバックアップは高々24秒で完了し、記憶部2からの読み出しと記憶部3への書き込みが並列に行われる場合には12秒で完了する。

【0241】

一方、本実施形態を適用せずに、512GBの容量の記憶部2から全データを読みだして記憶部3に書き込むことによるデータのバックアップには、20分～41分の時間を要する。すなわち、本実施形態の適用によって、実質的なバックアップに要する時間は、記憶容量×1に相当するデータ書き込みの後には63%が減少し、記憶容量×4.6に相当するデータ書き込みの後には99%が減少することになる。

【0242】

このように、本実施形態の適用により、ユーザ自身によるデータバックアップ作業は不要となる。また、バックグラウンドバックアップによる情報処理装置 1 1 1 への処理の負荷は大幅に低減され、アプリケーションプログラム 4 0 0 は記憶部 2 から記憶部 3 へのデータバックアップの影響をほとんど受けることなく、論理ドライブ 4 を利用可能である。信頼性劣化した記憶部 2 のデータ書き換えは、ディスクステータス格納部の書き換えしか発生しないため、記憶部 2 へのさらなるデータ書き込みによる記憶部 2 の故障の可能性を低減することができる。論理ドライブ 4 への新規データ書き込みは記憶部 2 ではなく信頼性のよい記憶部 3 に書き込まれるため、書き込みデータの損失を防止することができる。記憶部 2 が寿命到達してさらなるデータ書き込みが阻止される場合であっても、記憶部の上位レイヤである論理ドライブ 4 は読み出しと書き込みの両方が可能なドライブとして振る舞うため、アプリケーションプログラムなど上位のソフトウェアレイヤは論理ドライブ 4 が寿命到達モードであっても通常状態であっても同等の論理ドライブとして取り扱うことができるため、本実施形態導入のためのアプリケーションプログラム修正を必要とせず、本実施形態採用のシステムへの移行が容易である。

【 0 2 4 3 】

(第 3 の実施形態)

第 3 の実施形態では、本発明をディスクアレイからなる情報処理システム 1 に適用した場合について説明する。図 5 8 は、第 3 の実施形態の情報処理システム 1 を示すものである。情報処理システム 1 はストレージアレイ装置 1 0 0 3、記憶部 2 A ~ 2 D、記憶部 3、ストレージアレイ装置 1 0 0 3 と記憶部 2 A ~ 2 D および記憶部 3 とを接続するインタフェース 1 9、クライアント 1 0 0 2、クライアント 1 0 0 2 とストレージアレイ装置 1 0 0 3 とを接続するストレージネットワーク 1 0 0 0 とを具備する。情報処理システム 1 において、記憶部 2 A ~ 2 D はストレージアレイ装置 1 0 0 3 に接続されてそれぞれが論理単位としての論理スロットとして認識され、それら論理スロットを用いて R A I D (Redundant Arrays of Inexpensive Disks)アレイが構築される。ストレージアレイ装置 1 0 0 3 には、さらにデータ引き継ぎ先としての記憶部 3 が接続可能である。本実施形態においては寿命到達時処理前の R A I D アレイを構築する記憶部は 4 つである場合について説明するが、2 乃至複数の任意の数の記憶部を用いて R A I D アレイを構築してもよい。また、本実施形態では R A I D アレイとして R A I D 5 を用いるが、R A I D 0、R A I D 2、R A I D 3、R A I D 4、R A I D 6 など他の R A I D 技術や他のディスクアレイ実装形態を用いてディスクアレイを構築した場合についても本実施形態は適用可能である。

【 0 2 4 4 】

ネットワーク 1 0 0 0 は、ストレージアクセスのためのストレージネットワークであり、たとえばファイバチャネル(Fibre Channel)やイーサネット(登録商標)が用いられる。特に、ストレージネットワーク 1 0 0 0 として、例えば S A N (Storage Area Network)や N A S (Network Attached Storage)が用いられる。S A N には、たとえば、F C - S A N (Fibre Channel Storage Area Network)や I P - S A N (Internet Protocol Storage Area Network)が用いられ、その上位プロトコルとして、たとえば、S C S I (Small Computer System Interface)が用いられる。本実施形態では、ストレージネットワーク 1 0 0 0 として I P - S A N を採用した例を示し、その上位プロトコルとしては i S C S I (Internet Small Computer System Interface)が利用されるものとする。ストレージネットワーク 1 0 0 0 にはネットワークスイッチ 1 0 0 0 1 やハブ(図示せず)が含まれる。

【 0 2 4 5 】

クライアント 1 0 0 2 は、ストレージネットワーク 1 0 0 0 に接続され、所望の処理を遂行するコンピュータである。クライアント 1 0 0 2 は、典型的には、プロセッサと、メインメモリと、通信インタフェースと、ローカル入出力装置等のハードウェア資源を備え、また、デバイスドライバやオペレーティングシステム(O S)、アプリケーションプログラム等のソフトウェア資源を備える(図示せず)。これによって、クライアント 1 0 0 2 は、プロセッサの制御の下、各種のプログラムを実行して、ハードウェア資源との協働作用により処理を実現する。例えば、クライアント 1 0 0 2 は、プロセッサの制御の下、

業務アプリケーションプログラムを実行することにより、ストレージネットワーク 1000 を経由してストレージレイ装置 1003 に I/O アクセスし、所望の業務システムを実現する。または、クライアント 1002 は、データベース管理システム(Database management system, DBMS)が稼働しているサーバであるデータベースサーバ(DBサーバ)であってもよく、その場合、ストレージネットワーク 1000 またはその他のネットワーク(図示せず)を経由して DBサーバに接続されたクライアント(図示せず)からのデータ読み出し要求を受信すると、ストレージレイ装置 1003 からデータを読みだしてそのクライアントに読み出しデータを送信し、そのクライアントからのデータ書き込み要求を受信するとそのクライアントから書き込みデータを受信してストレージレイ装置 1003 にデータを書き込む。

10

【0246】

ストレージレイ装置 1003 は論理スロット 0 ~ 3 を RAID の構成単位として使用する。論理スロットは、第 2 の実施形態における論理ドライブに対応する。記憶部 2A ~ 2D のいずれか一つが寿命到達する前である通常状態において、インタフェース 19 経由で記憶部 2A ~ 2D がストレージレイ装置 1003 に接続されており、論理スロット 0 には記憶部 2A が、論理スロット 1 には記憶部 2B が、論理スロット 2 には記憶部 2C が、論理スロット 3 には記憶部 2D が割り当てられている。これにより、ストレージレイ装置 1003 は、4 台の記憶部 2A ~ 2D に対応する 4 つの論理スロットを RAID 5 により仮想的な 1 台の論理デバイスとしてクライアント 1002 に通知する。クライアント 1002 はストレージレイ装置にアクセスするための LBA (以下、「アレイ LBA」「ALBA」と呼ぶ)を送信する。RAID コントローラ 1005 内の制御部 200 はアレイ LBA を論理スロット番号と記憶部 2A ~ 2D にアクセスするための LBA (以下、「記憶部 LBA」「SLBA」と呼ぶ)に変換し、制御部 200 は論理スロット番号で特定された記憶部 2A ~ 2D のうち少なくとも 1 台の記憶部の SLBA に対してアクセス命令を送信する。

20

【0247】

ストレージレイ装置 1003 は、それ単体で、クライアント 1002 に対するデータストレージサービスを提供してもよいし、ストレージレイ装置 1003 と図示しない別のストレージレイ装置から仮想的に構成される 1 つのストレージ装置として、クライアント 1002 に対するデータストレージサービスを提供してもよい。ストレージレイ装置 1003 には、クライアント 1002 に提供するための 1 つ以上の論理デバイス(LDEV)が形成される。

30

【0248】

論理デバイスは、クライアント 1002 が認識しうる論理的な記憶装置であり、論理ユニット(LU)が割り当てられ、クライアント 1002 は、物理デバイス上に形成された論理デバイスを論理ユニットとして認識する。各論理ユニットには、論理ユニット番号(LUN)が付与される。また、論理ユニットは、論理セクタ(論理ブロック)に分割され、各論理セクタには、アレイ LBA が割り当てられる。クライアント 1002 は、論理ユニット番号及びアレイ LBA からなる論理アドレスを含むコマンドをストレージレイ装置 1003 に与えることにより、特定の論理ユニットにおける特定の論理セクタに対してアクセスすることができる。iSCSI を利用した本実施形態では、クライアント 1002、ストレージレイ装置 1003 はそれぞれ、iSCSI ネームが割り当てられた iSCSI ノードとしてのイニシエータ及びターゲットとして機能し、従って、クライアント 1002 及びストレージレイ装置 1003 は、IP アドレス及び TCP ポート番号の組み合わせで特定されるネットワークポータルを介して、iSCSI PDU を送受する。従って、クライアント 1002 は、iSCSI ネーム、IP アドレス及び TCP ポート番号を指定することにより、ネットワーク 1000 上のストレージレイ装置 1003 を認識し、その論理ユニット内の論理セクタに対してアクセスする。

40

【0249】

記憶部 2A ~ 2D は、インタフェース 19 を経由してストレージレイ装置 1003 に

50

対して接続される記憶部である。記憶部 2 A ~ 2 D は、たとえば、第 1 の実施形態に記載の記憶部 2 と同等の記憶部をそれぞれ用いることができる。本実施形態においては、例として、記憶部 2 A ~ 2 D は第 1 の実施形態に記載の S S D を用いた場合について説明する。信頼性劣化後に廃棄して設置スペースを削減しかつ情報処理システム 1 全体の消費電力を削減するという観点では、記憶部 2 A ~ 2 D はストレージレイ装置 1 0 0 3 に対して物理的に着脱可能であることが望ましい。

【 0 2 5 0 】

記憶部 3 は記憶部 2 A ~ 2 D のいずれかが一つが寿命到達または寿命到達目前であると判定された後に、新たにストレージレイ装置 1 0 0 3 に接続される記憶部であり、たとえば第 1 の実施形態に記載の記憶部 3 と同等の記憶部を用いることができる。本実施形態においては、記憶部 3 は第 1 の実施形態に記載の S S D を用いた場合について説明する。記憶部 3 の接続前の設置スペースを削減しかつシステム 1 全体の消費電力を削減するという観点、および記憶部 3 の信頼性劣化後に記憶部 3 を廃棄して設置スペースを削減しかつ情報処理システム 1 全体の消費電力を削減するという観点では、記憶部 3 はストレージレイ装置 1 0 0 3 に対して物理的に着脱可能であることが望ましい。

【 0 2 5 1 】

R A I D コントローラ 1 0 0 5 はディスクインタフェース 1 0 0 7 に接続された記憶部の R A I D アレイの構築・管理をつかさどり、制御部 2 0 0 を具備する。制御部 2 0 0 は、R A I D コントローラ内のメモリに格納されたファームウェアやソフトウェア、あるいは R A I D コントローラ 1 0 0 5 内のハードウェアなどの実装形態をとる。制御部 2 0 0 は、ネットワークインタフェース 1 0 0 4 経由でクライアント 1 0 0 2 やネットワークスイッチ 6 0 0 9 等から命令を受信すると、ディスクインタフェース 1 0 0 7 経由で各記憶部に読み出し命令や書き込み命令やその他命令やデータを送信したり、各記憶部から応答やデータを受信したり、ネットワークインタフェース 1 0 0 4 経由でクライアント 1 0 0 2 に応答やデータを送信したりする。制御部 2 0 0 は記憶部 2 A ~ 2 D や記憶部 3 の統計情報やステータス記憶領域 5 1 0 やスロット番号記憶領域 5 3 0 や引き継ぎ履歴記憶領域 5 5 0 の制御および管理を行い、統計情報に基づく寿命到達時処理やデータ引き継ぎ処理等を行う。なお、ステータス記憶領域 5 1 0、スロット番号記憶領域 5 3 0 および引き継ぎ履歴記憶領域 5 5 0 は各記憶部内に格納されるのではなく、R A I D コントローラ 1 0 0 5 内のメモリ領域（図示せず）など、情報処理システム 1 内の記憶領域に格納されるようにしてもよい。本実施形態では、ステータス記憶領域 5 1 0 に格納されるデータは、0 乃至 5 の値を取りうる。それぞれの値は、対応する記憶部がそれぞれ

- 0 : 初期ディスク状態
- 1 : 通常状態
- 2 : 信頼性劣化状態
- 3 : データ引き継ぎ元状態（保護状態）
- 4 : データ引き継ぎ先状態
- 5 : 廃棄対象状態

の状態であることを示す。

【 0 2 5 2 】

図 5 9 に本実施形態における通常状態のアレイ L B A (A L B A) と記憶部 L B A (S L B A) の変換方法を示す。R A I D 5 を採用する制御部 2 0 0 は、3 つの連続する論理セクタである $A L B A = 3 q, 3 q + 1, 3 q + 2$ (q は 0 以上の任意の整数) を組として、それぞれのデータであるアレイデータ $D(A L B A = 3 q), D(A L B A = 3 q + 1), D(A L B A = 3 q + 2)$ に対して 1 論理セクタ分のパリティデータである $P(3 q, 3 q + 2)$ を生成する。

【 0 2 5 3 】

パリティデータ $P(3 q, 3 q + 2)$ は、

$$P(3 q, 3 q + 2) = (D(A L B A = 3 q) \text{ XOR } D(A L B A = 3 q + 1) \text{ XOR } D(A L B A = 3 q + 2))$$

のように、 $D(ALBA = 3q)$, $D(ALBA = 3q+1)$, $D(ALBA = 3q+2)$ 内の論理セクタ内オフセットが同じである各ビットに対して排他的論理和をとって得られたビットをパリティデータ $P(3q, 3q+2)$ の論理セクタ内の各オフセットのビットデータとすることにより得られる。例えば、 $ALBA = 0$ のデータである $D(0)$ および $ALBA = 1$ のデータである $D(1)$ および $ALBA = 2$ のデータである $D(2)$ からパリティデータ $P(0, 2)$ が制御部 200 によって計算される。通常状態においては、アレイデータ $D(ALBA = 3q)$, $D(ALBA = 3q+1)$, $D(ALBA = 3q+2)$ およびパリティデータ $P(3q, 3q+2)$ は、図 59 に示すように、記憶部 2A ~ 記憶部 2D に分散管理される。たとえば $ALBA = 1$ の場合、対応するデータ $D(1)$ は論理スロット 1 に割り当てられている記憶部 2B の $LB A = SLBA = 1$ に格納されており、クライアント 1002 から $ALBA = 1$ を受信すると制御部 200 は論理スロット番号 = 1 および $SLBA = 1$ を特定する。

10

【0254】

たとえば、制御部 200 がクライアント 1002 から $ALBA = 1$ の読み出し命令を受信すると、制御部 200 は $D(1)$ の格納先である論理スロット番号 = 1 および $SLBA = 1$ を特定し、論理スロット番号 = 1 に接続されている記憶部 2B に対して $SLBA = 1$ の読み出し命令を送信し、受信した読み出しデータをクライアント 1002 に送信する。もし記憶部 2B から応答が無かった場合または記憶部 2B からエラーが返信されてきた場合、制御部 200 は論理スロット 1 以外のスロットである論理スロット 0 および論理スロット 2 ~ 3 に接続された記憶部 2A と記憶部 2C と記憶部 2D に $SLBA = 1$ の読み出し命令を送信し、受信したデータである $D(0)$, $D(2)$, $P(0, 2)$ から排他的論理和によって $D(1)$ を復元し、復元した $D(1)$ をクライアント 1002 に送信する。なお、 $D(1)$ の読み出し時に $D(0)$, $D(2)$, $P(0, 2)$ のデータ読み出しを並列して行なっても良い。

20

【0255】

たとえば、制御部 200 がクライアント 1002 から $ALBA = 1$ の書き込み命令および書き込みデータを受信すると、制御部 200 は $D(1)$ の格納先である論理スロット番号 = 1 および $SLBA = 1$ を特定し、論理スロット番号 = 1 の記憶部 2B に $SLBA = 1$ の書き込み命令を送信して書き込みデータを書き込むとともに、論理スロット番号 1 以外のスロットに接続されておりかつ $SLBA = 1$ にパリティデータが格納されていない記憶部である記憶部 2A および記憶部 2C の $SLBA = 1$ からデータ $D(0)$ と $D(2)$ を読み出し、 $D(0)$, $D(1)$, $D(2)$ からパリティデータ $P(0, 2)$ を生成して記憶部 2D の $SLBA = 1$ に書き込む。

30

【0256】

図 60 に、論理スロット 1 の引き継ぎ状態遷移直後のアレイ $LB A$ と記憶部 $LB A$ の変換方法を示す。引き継ぎ状態遷移直後においては、アレイ $LB A$ と記憶部 $LB A$ の変換方法は、図 59 に示す通常状態と同様である。

【0257】

図 58 に示したように、記憶部 2A ~ 2D および記憶部 3 はそれぞれステータス格納領域 510 と論理スロット番号格納領域 530 を具備し、記憶部 3 は引き継ぎ履歴格納領域 550 を具備する。引き継ぎ履歴格納領域 550 は後述するデータ引き継ぎが完了すると消去されることにより、ユーザデータ格納など他の用途に使用されてもよい。ステータス格納領域 510 と論理スロット番号格納領域 530 と引き継ぎ履歴格納領域 550 には、 $LB A$ の領域が割り当てられても良い。あるいは、ステータス格納領域 510 と論理スロット番号格納領域 530 と引き継ぎ履歴格納領域 550 には $LB A$ ではなくログページアドレスが割り当てられたログ領域が割り当てられてもよく、この場合には、たとえば、非特許文献 1 の ACS-3 に記載されている 2Fh Read Log Ext により該領域の読み出しが行われ、非特許文献 1 の ACS-3 に記載されている 3Fh Write Log Ext により該領域への書き込みが行われる。引き継ぎ履歴格納領域 550 には後述するデータ引き継ぎ状態のログが格納される。

40

【0258】

引き継ぎ履歴格納領域 550 に記録される履歴データとして、本実施形態では図 61 に

50

示すようなテーブルデータを使用する。データ引き継ぎ状態に記憶部3のLBA(=SLBA)にセクタ長Xのデータが書き込まれると、制御部200は引き継ぎ履歴格納領域550に(SLBA,X)を追記する。LBA=SLBA~SLBA+X-1の領域に対し、オーバーラップするLBAや連続するLBAが引き継ぎ履歴格納領域550にすでに記録されている場合には、それらを合併したログを引き継ぎ履歴格納領域550に記録し、引き継ぎ履歴格納領域550から合併されたログを削除することが望ましい。たとえば、(SLBA,X)のLBA領域Aを新たに引き継ぎ履歴格納領域550に記入する場合であって、(SLBA-a,a)のLBA領域Bがすでに引き継ぎ履歴格納領域550に記録されている時、LBA領域AとLBA領域Bは連続しているため、A+Bの領域を示す(SLBA-a,a+X)のデータを(SLBA-a,a)の履歴部に上書きすることで、ログデータ量を増大させることなく引き継ぎ履歴格納領域550を更新することができる。

10

【0259】

制御部200は、第1の実施形態の図24や図29と同様にして、ディスクインタフェース1007に接続されたそれぞれの記憶部が寿命に到達しているか否か、寿命到達目前であるか否か、あるいは故障目前であるか否かを判定し、寿命に到達しているか場合、寿命到達目前である場合、あるいは故障目前である場合に寿命到達時処理を行う。寿命到達判定は、第1の実施形態と同様に、図24のような一定時間ごと、または一定処理数おき、または一定データ送受信おき、または図29や図30のように記憶部からのコマンド応答がエラーであった場合などに実施される。

【0260】

20

図62は、たとえば、記憶部2Bが寿命に到達していると判定された場合の、記憶部2Bの寿命到達時処理を示す。記憶部2Bが寿命に到達していると判定された場合(ステップS500)、制御部200は記憶部2Bのステータスを1から2に書き換える(ステップS501)。制御部200は空きスロットに新しい記憶部をディスクインタフェース1007に接続するようディスプレイ装置やLED等を経由してネットワーク管理者に通知することが望ましい(ステップS502)。あるいは、ディスクロード・アンロード装置(図示せず)として記憶部2A~2Dや記憶部3のディスクインタフェース1007に対する物理的な着脱を自動的に行う機械装置がシステム1に搭載されている場合、新たな記憶部をインタフェース1007に接続するよう該ディスクロード・アンロード装置に対して命令してもよい(ステップS502)。

30

【0261】

新たな記憶部として記憶部3が接続された時(ステップS503)、制御部200は記憶部3のステータス510を4に書き換え(ステップS504)、記憶部3の論理スロット番号格納領域530に記憶部2Bの論理スロット番号格納領域530のデータをコピーする(ステップS505)。今回の例では、記憶部2Bは図57に示すように論理スロット1として割り当てられているため、記憶部3の論理スロット番号格納領域530には1が書き込まれる。制御部200は記憶部2Bのステータス510を3に書き換え(ステップS506)、RAIDコントローラ1005に記憶部2Bと記憶部3を同一の論理スロットである論理スロット1として認識させる(ステップS507)。該寿命到達処理後、記憶部2Bのステータスは3に、記憶部3のステータスは4になっており、論理スロット1はデータ引き継ぎ状態に遷移していることになる。

40

【0262】

本実施形態では、各記憶部のステータスは、各記憶部のステータス格納領域510に不揮発に格納されている。図63はRAIDコントローラ1005が起動した時、およびディスクインタフェース1007に記憶部が接続された時の制御部200の処理手順を示すものである。RAIDコントローラ1005の起動ごとに制御部200はステータス格納領域510を読み出すことで各記憶部のステータスを認識する。制御部200は、各記憶部のステータスを認識し、かつ各記憶部の論理スロット番号を論理スロット番号格納領域530から読み出すことで、各論理スロット0~3それぞれがデータ引き継ぎ状態か否かを判定する。

50

【 0 2 6 3 】

すなわち、制御部 2 0 0 は、ディスクステータス = 0 であるとき（ステップ S 5 1 1）、記憶部が初期ディスク状態であると認識し（ステップ S 5 1 2）、ディスクステータス = 1 であるとき（ステップ S 5 1 3）、記憶部が通常状態であると認識し（ステップ S 5 1 4）。ディスクステータス = 2 であるとき（ステップ S 5 1 5）、記憶部が信頼性劣化状態であると認識し（ステップ S 5 1 6）、ディスクステータス = 3 であるとき（ステップ S 5 1 7）、記憶部がデータ引き継ぎ作業中のデータ引き継ぎ元状態（保護状態）として認識し（ステップ S 5 1 8）、ディスクステータス = 4 であるとき（ステップ S 5 1 9）、記憶部がデータ引き継ぎ作業中のデータ引き継ぎ先状態として認識し（ステップ S 5 2 0）、ディスクステータス = 5 であるとき（ステップ S 5 2 1）、記憶部が廃棄対象状態として認識し（ステップ S 5 2 2）、ディスクステータスが 0 ~ 5 以外であるときは、不正なディスクとみなす（ステップ S 5 2 3）。

10

【 0 2 6 4 】

（論理ドライブからの読み出し 1）

図 6 4 はクライアント 1 0 0 2 からストレージアレイ装置 1 0 0 3 に読み出し要求が送信された時の、制御部 2 0 0 の処理手順を示すものである。制御部 2 0 0 は、クライアント 1 0 0 2 から読み出し対象アレイ L B A である A L B A = c A L B A の読み出し命令を受信する（ステップ S 5 3 0）。制御部 2 0 0 は、c A L B A から論理スロット番号 c S L O T と、読み出し対象記憶部 L B A である S L B A = c S L B A を計算する（ステップ S 5 3 1）。制御部 2 0 0 は、c S L O T の論理スロットの記憶部が故障中であるか否かを判定し（ステップ S 5 3 2）、c S L O T の論理スロットの記憶部が故障中でない場合、つぎに c S L O T の論理スロットの記憶部がデータ引き継ぎ状態であるか否かを判定する（ステップ S 5 3 3）。

20

【 0 2 6 5 】

c S L O T の論理スロットの記憶部がデータ引き継ぎ状態である場合（ステップ S 5 3 3）、制御部 2 0 0 は、c S L O T のスロット以外からパリティを含むデータを読み出して、それらのデータを使用して c S L O T の c S L B A のデータを復元し、復元したデータをクライアント 1 0 0 2 に送信する（ステップ S 5 3 4）。そして、制御部 2 0 0 は、c S L O T の c S L B A の復元データをデータ引き継ぎ先の記憶部 3 に書き戻し、引き継ぎ履歴格納領域 5 5 0 に引き継ぎ履歴を記録する（ステップ S 5 3 5）。c S L O T の論理スロットの記憶部がデータ引き継ぎ中でない場合（ステップ S 5 3 3）、制御部 2 0 0 は、c S L O T の記憶部からデータ D（c S L B A）を読み出し、読み出したデータをクライアント 1 0 0 2 に送信する。このように、読み出し対象の論理スロットがデータ引き継ぎ状態である場合は、読み出し対象以外のスロットから読み出し対象データを復元することで、データ引き継ぎ状態の記憶部のデータを損失した場合であっても救済可能であるとともに、データ引き継ぎ状態の記憶部からの読み出しを減らし、リードディスタurb（データを読み出したメモリセルと同一のブロックに含まれる非選択メモリセルのフローティングゲートに微妙な電荷が蓄えられることによって、記憶したデータに誤りが生じる現象）を抑制する。

30

【 0 2 6 6 】

c S L O T の論理スロットの記憶部が故障中である場合、制御部 2 0 0 は、データ引き継ぎ状態のスロットが存在するか否かを判定し（ステップ S 5 3 7）、データ引き継ぎ状態のスロットが存在しない場合、c S L O T のスロット以外からデータを読み出して、それらのデータを使用して c S L O T の c S L B A のデータを復元し、復元したデータをクライアント 1 0 0 2 に送信する（ステップ S 5 3 8）。このように読み出し対象のスロットが故障中であり、かつ他のスロットにデータ引き継ぎ状態のスロットが存在しない場合は、読み出し対象以外のスロットから読み出し対象データを復元する。

40

【 0 2 6 7 】

c S L O T の論理スロットの記憶部が故障中であって、データ引き継ぎ状態のスロットが存在する場合（ステップ S 5 3 7）、制御部 2 0 0 は、引き継ぎ履歴格納領域 5 5 0 を

50

読み出して読み出し対象 S L B A のデータが引き継ぎ元と引き継ぎ先のどちらに存在するかを判定し（ステップ S 5 3 9）、データ引き継ぎ先記憶部、データ引き継ぎ元記憶部および通常状態記憶部から読み出したデータから c S L O T の c S L B A のデータを復元し、復元したデータをクライアント 1 0 0 2 に送信する（ステップ S 5 4 0）。データ引き継ぎ元のデータを使用した場合は、使用したデータ引き継ぎ元データをデータ引き継ぎ先の記憶部 3 に書き戻し、データ引き継ぎ先の引き継ぎ履歴格納領域 5 5 0 に引き継ぎ履歴を記録する（ステップ S 5 4 1）。

【 0 2 6 8 】

（論理ドライブからの読み出し 2）

図 6 5 はクライアント 1 0 0 2 からストレージアレイ装置 1 0 0 3 に読み出し要求が送信された時の、制御部 2 0 0 の他の処理手順を示すものである。図 6 5 に示す処理手順では、図 6 4 のステップ 5 3 3 の判定が Y e s のときの処理を、ステップ S 5 3 4、S 5 3 5 からステップ S 5 5 0 ~ S 5 5 4 に置換している。また、図 6 5 では、図 6 4 のステップ S 5 4 1 の処理を削除している。c S L O T の論理スロットの記憶部がデータ引き継ぎ状態である場合（ステップ S 5 3 3）、制御部 2 0 0 は、データ引き継ぎ先の記憶部 3 から引き継ぎ履歴を読み出し、c S L B A のデータが引き継ぎ履歴に含まれるか否かを判定する（ステップ S 5 5 1）。制御部 2 0 0 は、c S L B A のデータが引き継ぎ履歴に含まれる場合は、データ引き継ぎ先の記憶部 3 から c S L B A のデータを読み出し、読み出したデータをクライアント 1 0 0 2 に送信する（ステップ S 5 5 4）。制御部 2 0 0 は、c S L B A のデータが引き継ぎ履歴に含まれない場合は、データ引き継ぎ元の記憶部から c S L B A のデータを読み出し、読み出したデータをクライアント 1 0 0 2 に送信する（ステップ S 5 5 2）。そして、データ引き継ぎ元から読み出したデータをデータ引き継ぎ先の記憶部 3 に書き戻し、データ引き継ぎ先の記憶部 3 の更新履歴 5 5 0 を更新する（ステップ S 5 5 3）。

【 0 2 6 9 】

（論理ドライブへの書き込み）

図 6 6 は、本実施形態におけるクライアント 1 0 0 2 からの書き込み命令に対する処理を示すフローチャートである。制御部 2 0 0 は、クライアント 1 0 0 2 からアレイ L B A である c L B A とセクタ長を含む書き込み命令を受信すると（ステップ S 5 6 0）、クライアント 1 0 0 2 から書き込みデータを受信する（ステップ S 5 6 1）。制御部 2 0 0 は、データを書き込むべき論理スロット番号 c S L O T を c A L B A から計算し、データを書き込むべき記憶部 L B A である c S L B A を c A L B A から検索し、c S L B A のパリティデータ格納先の論理スロット番号 c P S L O T を c A L B A から計算する（ステップ S 5 6 2）。制御部 2 0 0 は、全スロットの c S L B A からデータを並列に読み出す。故障スロットがある場合は、故障スロットのデータは、故障スロット以外からデータを読み出して、故障スロットのデータを復元する（ステップ S 5 6 3）。以後、制御部 2 0 0 は、本体データ書き込みタスクとパリティデータ書き込みタスクを並列に処理する。

【 0 2 7 0 】

本体データ書き込みタスクは次のように実行される。制御部 2 0 0 は、c S L O T が故障中かどうかを判定し（ステップ S 5 6 4）、故障中の場合にはデータを書き込まずタスクを終了する。故障中でない場合には、c S L O T がデータ引き継ぎ状態かどうかを判定し（ステップ S 5 6 5）、データ引き継ぎ状態でない場合には、スロット c S L O T の c S L B A に対してクライアント 1 0 0 2 からの受信データを書き込む（ステップ S 5 6 6）。データ引き継ぎ状態の場合には、スロット c S L O T に割り当てられたデータ引き継ぎ元記憶部とデータ引き継ぎ先記憶部のうちデータ引き継ぎ先記憶部の c S L B A に受信データを書き込み（ステップ S 5 6 7）、データ引き継ぎ先記憶部の引き継ぎ履歴格納部 5 5 0 に c S L B A とセクタサイズをデータ引き継ぎ履歴として記録する（ステップ S 5 6 8）。引き継ぎ履歴格納部 5 5 0 に、c S L B A の前後の連続した L B A や重複した L B A の引き継ぎ履歴が存在する場合には、それらのログを合併した引き継ぎ履歴を引き継ぎ履歴 5 5 0 に記入し、合併元の引き継ぎ履歴を削除する（無効化する）（ステップ S 5

10

20

30

40

50

68)。

【0271】

なお、引き継ぎ履歴550からの引き継ぎ履歴削除は論理的な削除でよく、NANDメモリ16のブロック消去によって物理的に消去しなくてもよい。たとえば、ログを無効化するフラグを引き継ぎ履歴550に記入したり、更新後の引き継ぎ履歴550を更新前の引き継ぎ履歴550とは別の領域に格納し引き継ぎ履歴550の格納位置を示すポイントを更新したりすることで、引き継ぎ履歴550から引き継ぎ履歴が論理的に削除される。

【0272】

パリティデータ書き込みタスクは次のように実行される。制御部200は、全論理スロットからRAIDコントローラ1005内のメモリに読み出されたデータに対してクライアント1002から受信した書き込みデータを上書きし、パリティデータを再計算する(ステップS570)。制御部200は、cPSLOTが故障中かどうかを判定し(ステップS571)、故障中の場合にはパリティデータを書き込まずタスクを終了する。故障中でない場合には、cPSLOTがデータ引き継ぎ状態かどうかを判定し(ステップS572)、データ引き継ぎ状態でない場合には、スロットcPSLOTのcSLBAに対してパリティデータを書き込む(ステップS573)。データ引き継ぎ状態の場合には、スロットcPSLOTに割り当てられたデータ引き継ぎ元記憶部とデータ引き継ぎ先記憶部のうちデータ引き継ぎ先記憶部のcSLBAにパリティデータを書き込み(ステップS574)、データ引き継ぎ先記憶部の引き継ぎ履歴格納部550にcSLBAとセクタサイズを引き継ぎ履歴として記録する。引き継ぎ履歴550に、cSLBAの前後の連続したLBAや重複したLBAの引き継ぎ履歴が存在する場合には、それらの引き継ぎ履歴を合併した引き継ぎ履歴を引き継ぎ履歴格納領域550に記入し、合併元の引き継ぎ履歴を引き継ぎ履歴550から削除する(ステップS575)。

【0273】

このようにして、制御部200は、クライアントからの書き込み命令に対し、データ引き継ぎ状態ではない各論理スロット0,2,3に割り当てられた各記憶部2A,2C,2Dに対してはSLBAの書き込み命令を送信する一方で、データ引き継ぎ中である論理スロット1に対しては、データ引き継ぎ元である記憶部2Bに対して書き込み命令を送信せず、データ引き継ぎ先である記憶部3に対して書き込み命令を送信し、記憶部3の引き継ぎ履歴格納領域550に引き継ぎ履歴を記録する。

【0274】

図67に更に詳細な例を示す。図67は、論理スロット1がデータ引き継ぎ状態で、データ引き継ぎ元として記憶部2Bが割り当てられ、データ引き継ぎ先として記憶部3が割り当てられており、論理スロット0,2,3は通常状態である場合に、ストレージレイ装置1003のALBA=1にD(1)を、ALBA=7にD(7)を、ALBA=16にD(16)を書き込む命令がクライアント1002から送信された場合の例を示すものである。

【0275】

(例1)

制御部200はALBA=1に対する新たなデータD(1)newの書き込み命令をクライアント1002から受信すると、ALBA=1からSLBA=0および論理スロット1を特定する。論理スロット1がデータ引き継ぎ状態ではない場合には、D(1)newは記憶部2BのSLBA=0に書き込まれるが、今回の例では論理スロット1はデータ引き継ぎ状態であるため、制御部200はD(1)newを記憶部3に書き込む。制御部200は前述した読み出し動作に基づき論理スロット0の記憶部2AからD(0)を読み出し、論理スロット2の記憶部2CからD(2)を読み出し、D(0)とD(1)newとD(2)から排他的論理和で新しいパリティデータP(0,2)newを計算し、P(0,2)newを論理スロット3の記憶部2DのSLBA=0に格納する。なお、処理速度向上のためには、記憶部3に対するD(1)newの書き込み命令と、記憶部2Aに対するD(0)の読み出し命令と、記憶部2Cに対するD(2)の読み出し命令は、制御部200から同時に送信される

ことが望ましく、さらには $D(0)$ と $D(2)$ の読み出しが完了した時点で、 $D(1)$ の書き込み完了を待たずして $P(0,2)_{new}$ を計算し記憶部2Dに書き込み命令を送信することが望ましい。制御部200は記憶部3の引き継ぎ履歴格納領域550に、記憶部3に書き込みを行ったSLBAとセクタカウントである(SLBA,セクタカウント)を記録する。たとえば、この例では記憶部3のSLBA = 0から1セクタのSLBA領域に書き込みを行ったので、引き継ぎ履歴格納領域550に(0,1)を追記する。前述のとおり、連続するLBA領域やオーバーラップするLBA領域がすでに引き継ぎ履歴格納領域550に記録されている場合には、これらを統合したSLBA領域を引き継ぎ履歴格納領域550に記録する。データ引き継ぎ状態ではない論理スロットに書き込みを行った場合には、引き継ぎ履歴550の更新は行わない。

10

【0276】

(例2)

制御部200はALBA = 7に対する新たなデータ $D(7)_{new}$ の書き込み命令をクライアント1002から受信すると、ALBA = 7からSLBA = 2および論理スロット2を特定する。論理スロット2はデータ引き継ぎ状態ではないため $D(7)_{new}$ は記憶部2BのSLBA = 0に書き込まれる。制御部200は、前述の読み出し動作に基づき論理スロット0、論理スロット3から $D(6)$ と $D(8)$ の読み出しを行い、パリティデータ $P(6,8)_{new}$ を計算し、論理スロット1に格納する。論理スロット1はデータ引き継ぎ状態であるため、制御部200は $P(6,8)_{new}$ を記憶部2Bではなく記憶部3に書き込み、記憶部3の引き継ぎ履歴格納領域550に、記憶部3に書き込みを行ったSLBAとセクタカウントである(SLBA,セクタカウント) = (2,1)を記録する。

20

【0277】

(例3)

制御部200はALBA = 16に対する新たなデータ $D(16)_{new}$ の書き込み命令をクライアント1002から受信すると、ALBA = 16からSLBA = 5および論理スロット1を特定する。今回の例では論理スロット1はデータ引き継ぎ状態であるため、制御部200は $D(16)_{new}$ を記憶部3に書き込む。制御部200は前述した読み出し動作に基づき論理スロット0の記憶部2Aから $D(15)$ を読み出し、論理スロット2の記憶部2Dから $D(17)$ を読み出し、 $D(15)$ と $D(16)_{new}$ と $D(17)$ から排他的論理和で新しいパリティデータ $P(15,17)_{new}$ を計算し、論理スロット2の記憶部2CのSLBA = 5に格納する。制御部200は記憶部3の引き継ぎ履歴格納領域550に、記憶部3に書き込みを行ったSLBAとセクタカウントである(SLBA,セクタカウント)を記録する。たとえば、この例では記憶部3のSLBA = 5から1セクタのLBA領域に書き込みを行ったので、引き継ぎ履歴領域550に(5,1)を追記する。前述のとおり、連続するSLBA領域やオーバーラップするSLBA領域がすでに引き継ぎ履歴550に記録されている場合には、これらを統合したSLBA領域を引き継ぎ履歴格納領域550に記録する。データ引き継ぎ状態ではないスロットに書き込みを行った場合には、引き継ぎ履歴格納領域550の更新は行わない。このようにして、第2の実施形態と同じように、クライアント1002から書き込み要求を受信するごとに、対応するSLBAのデータをデータ引き継ぎ元の記憶部2Bではなくデータ引き継ぎ先の記憶部3に書き込みデータ引き継ぎ履歴を記録することで、新規データ書き込みとデータ引き継ぎ動作を同時に行うことが可能になる。

30

40

【0278】

(バックグラウンドバックアップ)

制御部200はクライアント1002から命令を一定時間以上受信しない場合などのアイドル時、またはスタンバイモード遷移要求をクライアント1002から受信した場合などに、記憶部2Bから記憶部3へのバックグラウンドバックアップを行う。制御部200は、記憶部3の引き継ぎ履歴格納領域550を読み出し、引き継ぎ履歴格納領域550に記録されていないSLBAに対して、記憶部2Bからデータ読み出しを行い、記憶部3に該データを書き込み、記憶部3の引き継ぎ履歴格納領域550に該SLBAと書き込みデ

50

ータのセクタサイズを格納することでバックグラウンドバックアップが行われる。引き継ぎ履歴 550 への S L B A の格納の際、連続した S L B A 領域が存在する場合、またはオーバーラップした S L B A 領域が存在する場合には、この S L B A 領域と統合した S L B A 領域が引き継ぎ履歴 550 に格納され、統合前の S L B A 領域は引き継ぎ履歴格納領域 550 から削除されることが望ましい。

【0279】

引き継ぎ履歴格納領域 550 に格納される引き継ぎ履歴のデータサイズを低減するという観点、および後述するデータ引き継ぎ終了判定を高速に行うという観点では、引き継ぎ履歴 550 で示される S L B A 領域以外の S L B A 領域のうち小さく断片化された S L B A 領域（断片化した未データ引き継ぎ S L B A 領域）に対してバックグラウンドバックアップが優先的にされることが望ましい。断片化した未データ引き継ぎ S L B A 領域をバックアップすることにより、前記断片化未データ引き継ぎ領域の前記断片化未データ引き継ぎ領域の前後の既データ引き継ぎ S L B A 領域と、新たにデータ引き継ぎされた S L B A 領域が引き継ぎ履歴 550 に統合されて格納される。

【0280】

（データ引き継ぎ状況の監視）

制御部 200 は引き継ぎ履歴格納領域 550 を読み出すことでデータ引き継ぎ中の論理スロットのデータ引き継ぎ状況を監視する。図 68 は、引き継ぎ履歴を用いたデータ引き継ぎ状況の監視手順を示すものである。たとえば、制御部 200 は所定時間経過毎に引き継ぎ履歴格納領域 550 からデータ引き継ぎ履歴を読み出すことでデータ引き継ぎ状況を監視する（ステップ S600、S601）。制御部 200 は、全ての引き継ぎ対象 S L B A が引き継ぎ履歴 550 に含まれる場合、データ引き継ぎが完了したと判定する。たとえば、データ引き継ぎ元である記憶部 2B の全 S L B A が引き継ぎ履歴 550 に含まれる場合、データ引き継ぎが完了したと判定される（ステップ S602）。なお、データ引き継ぎ完了の判定として、例えば、データ引き継ぎ先である記憶部 3 の全 S L B A が引き継ぎ履歴 550 に含まれるか否かを判定するようにしてもよい。

【0281】

データ引き継ぎが完了したと判定された場合、制御部 200 はデータ引き継ぎ元である記憶部 2 のステータスを 5 の廃棄対象状態に変更し、データ引き継ぎ先である記憶部 3 のステータスを 1 の通常状態に変更することで（ステップ S603）、論理スロット 1 のデータ引き継ぎ状態を終了し、論理スロット 1 のデータ引き継ぎ状況監視を終了する。記憶部 2B が消費する電力を低減するという観点では、制御部 200 は記憶部 2 に低消費電力モードへの遷移要求を送信することが望ましい。データ引き継ぎ状態終了後、制御部 200 はディスプレイ 9 やポート近傍に設置された L E D を通じて、記憶部 2B をインタフェース 19 から取り外すようユーザや管理者に通知したり、記憶部 2B をインタフェース 19 から取り外すよう前記ディスクロード・アンロード装置に命令を送信することが望ましい。

【0282】

（全体のステータスの遷移）

図 69 は記憶部 2B が寿命到達する際の各記憶部のステータスの遷移を示したものである。記憶部 2A ~ 2D すべてが通常状態である場合には、記憶部 2A ~ 2D すべてのステータスは 1 になっており、記憶部 3 はディスクインタフェース 1007 に未接続状態である（ステップ 1）。記憶部 2B が寿命到達したと制御部 200 が判定すると、制御部 200 は記憶部 2B のステータスを 2 に変更する（ステップ 2）。制御部 200 からの通知または命令にもとづき、新しい記憶部として記憶部 3 がインタフェース 1007 に接続される（ステップ 3）。記憶部 3 が接続されると、制御部 200 はデータ引き継ぎ元としての記憶部 2B のステータスを 3 に変更し、データ引き継ぎ先としての記憶部 3 のステータスを 4 に変更することで、論理スロット 1 の状態をデータ引き継ぎ状態に移行させる（ステップ 4）。引き継ぎ履歴格納領域 550 から読み出される情報をもとに、記憶部 2B のすべての有効データが記憶部 3 にデータ引き継ぎされたと判定されると、制御部 200 は、

記憶部 2 B のステータスを 5 に変更し、記憶部 3 のステータスを 1 に変更する。以後、記憶部 2 A , 記憶部 3 , 記憶部 2 C , 記憶部 2 D が、元の記憶部 2 A , 記憶部 2 B , 記憶部 2 C , 記憶部 2 D であるかのように振る舞う (ステップ 1 に戻る)。その後さらに記憶部 2 A , 記憶部 3 , 記憶部 2 C , 記憶部 2 D のいずれかが寿命到達したと判定された場合、同様のステップ 2 ~ 5 が繰り返される。このようにして、ストレージアレイ装置 1 0 0 3 の記憶部のいずれか一つが寿命到達した場合、寿命到達目前である場合、故障寸前である場合となっても、本実施形態により記憶部のデータを新しい記憶部に容易にデータ引き継ぎすることが可能になる。本実施形態において、データ引き継ぎ動作はバックグラウンドで行われるため、バックアップ動作による性能劣化を抑制することが可能になる。

【 0 2 8 3 】

10

(1 台のデータ引き継ぎ状態と 1 台の故障が同時に発生した場合)

図 7 0 は、1 台のデータ引き継ぎ状態と他の 1 台の故障が同時に発生した場合の例を示すものである。図 7 0 では、論理スロット 1 がデータ引き継ぎ状態で、データ引き継ぎ元として記憶部 2 B が割り当てられ、データ引き継ぎ先として記憶部 3 が割り当てられており、論理スロット 0 , 2 , 3 は通常状態である場合に、論理スロット 3 の記憶部 2 D に故障が発生した場合の例を示すものである。

【 0 2 8 4 】

S L B A = 0 からの読み出しの場合、論理スロット 0 の D (0)、論理スロット 2 の D (2) に関しては復元不要であり、論理スロット 1 の D (1) に関しては、データ引き継ぎ先の記憶部 3 から最新のデータである D (1) n e w を読み出し可能である。S L B A = 2 からの読み出しの場合、論理スロット 0 の D (6)、論理スロット 2 の D (7) に関しては復元不要であり、論理スロット 3 の D (8) に関しては、D (6) と D (7) と、データ引き継ぎ先の記憶部 3 から P (6 , 8) n e w との排他的論理和をとることで復元可能である。S L B A = 5 からの読み出しの場合、論理スロット 0 の D (1 5) に関しては復元不要であり、論理スロット 1 の D (1 6) に関しては、データ引き継ぎ先の記憶部 3 から最新のデータである D (1 6) n e w を読み出し可能であり、論理スロット 3 の D (1 7) に関しては、D (1 5) とデータ引き継ぎ先の記憶部 3 の D (1 6) n e w と、論理スロット 2 の P (1 5 , 1 7) n e w との排他的論理和をとることで復元可能である。

20

【 0 2 8 5 】

30

(2 台のデータ引き継ぎ状態と 1 台の故障が同時に発生した場合)

図 7 1 は、2 台のデータ引き継ぎ状態と他の 1 台の故障が同時に発生した場合の例を示すものである。図 6 8 では、論理スロット 1 および論理スロット 2 がデータ引き継ぎ状態で、データ引き継ぎ元として記憶部 2 B 、 2 C が割り当てられ、データ引き継ぎ先として記憶部 3 B 、 3 C が割り当てられており、論理スロット 0 , 3 は通常状態である場合に、論理スロット 3 の記憶部 2 D に故障が発生した場合の例を示すものである。

【 0 2 8 6 】

S L B A = 0 からの読み出しの場合、論理スロット 0 の D (0)、論理スロット 2 の D (2) に関しては復元不要であり、論理スロット 1 の D (1) に関しては、データ引き継ぎ先の記憶部 3 B から D (1) n e w を読み出し可能である。S L B A = 2 からの読み出しの場合、論理スロット 0 の D (6) に関しては復元不要であり、論理スロット 2 の D (7) に関しては、データ引き継ぎ先の記憶部 3 C から最新のデータである D (7) n e w を読み出し可能であり、論理スロット 3 の D (8) に関しては、D (6) と、D (7) n e w と、P (6 , 8) n e w の排他的論理和をとることで復元可能である。S L B A = 5 からの読み出しの場合、論理スロット 0 の D (1 5) に関しては復元不要であり、論理スロット 1 の D (1 6) に関しては、データ引き継ぎ先の記憶部 3 B から最新のデータである D (1 6) n e w を読み出し可能であり、論理スロット 3 の D (1 7) に関しては、D (1 5) とデータ引き継ぎ先の記憶部 3 B の D (1 6) n e w と、論理スロット 2 のデータ引き継ぎ先の記憶部 3 C の P (1 5 , 1 7) n e w との排他的論理和をとることで復元可能である。

40

【 0 2 8 7 】

50

(1 台のデータ引き継ぎ状態とデータ読み出し不能エラーが同時に発生した場合)

図 7 2 は、データ引き継ぎ中に 1 台の記憶部で E C C 訂正不能エラーなどのアンコレクタブル リードエラーが発生した場合の例を示すものである。

(例 1)

制御部 2 0 0 は A L B A = 2 に対する新たなデータ D (2) の読み出し命令をクライアント 1 0 0 2 から受信すると、A L B A = 2 から S L B A = 0 および論理スロット 2 を特定する。この論理スロット 2 の記憶部 2 C からのデータ D (2) の読み出しの際に E C C 訂正不能エラーが発生した場合は、D (0) と、データ引き継ぎ先の記憶部 3 の D (1) new と、論理スロット 3 の P (0 , 2) new との排他的論理和をとることで D (2) を復元可能である。

10

(例 2)

制御部 2 0 0 は A L B A = 7 に対する新たなデータ D (7) の読み出し命令をクライアント 1 0 0 2 から受信すると、A L B A = 7 から S L B A = 2 および論理スロット 3 を特定する。この論理スロット 3 の記憶部 2 D からのデータ D (8) の読み出しの際に E C C 訂正不能エラーが発生した場合は、D (6) と、D (7) new と、記憶部 3 の P (6 , 8) new との排他的論理和をとることで D (8) を復元可能である。

(例 3)

制御部 2 0 0 は A L B A = 9 に対する新たなデータ D (9) の読み出し命令をクライアント 1 0 0 2 から受信すると、A L B A = 9 から S L B A = 3 および論理スロット 1 を特定する。この論理スロット 1 の記憶部 2 B からのデータ D (9) の読み出しの際に E C C 訂正不能エラーが発生した場合は、D (1 0) と、D (1 1) と、P (9 , 1 0) との排他的論理和をとることで D (9) を復元可能である。

20

【 0 2 8 8 】

(第 4 の実施形態)

第 4 の実施形態では、本発明をストレージネットワークからなる情報処理システム 1 に適用した場合について説明する。本実施形態においては、情報処理システム 1 は図 7 3 に示すように構成される。なお、情報処理システム 1 は図 7 4 に示すように構成されてもよい。図 7 3 に示すように、本実施形態における情報処理システム 1 は、記憶部 2、記憶部 3、1 乃至複数のその他の記憶部である記憶部 6 0 0 4、1 乃至複数のメタデータサーバ 6 0 0 3、サーバ 6 0 0 1、記憶部 2 と記憶部 3 と記憶部 6 0 0 4 とメタデータサーバ 6 0 0 3 とサーバ 6 0 0 1 を相互接続するストレージネットワーク 1 0 0 0、1 乃至複数のクライアント 6 0 0 2、サーバ 6 0 0 1 とクライアント 6 0 0 2 を相互接続するネットワーク 6 0 0 0 からなる。図 7 4 の場合は、ストレージネットワーク 1 0 0 0 は、チャンネルエクステンダ、WAN などの長距離ネットワークを介してストレージネットワーク 6 0 0 1 b に接続され、このストレージネットワーク 6 0 0 1 b に記憶部 3 が接続されている。

30

【 0 2 8 9 】

記憶部 2、3 は、ステータス格納領域 5 1 0 を具備する。記憶部 2、3 には、第 1 の実施形態と同様のデータ引き継ぎ対象記憶部アドレス格納領域 2 0 0 7、3 0 0 7 (図 3 2 参照) を具備させてもよいし、第 2 の実施形態と同様の論理ドライブ ID 格納領域 5 2 0 を具備させてもよい。ステータス格納領域 5 1 0 は、前述と同様、対応する記憶部がそれぞれ

40

0 : 初期ディスク状態

1 : 通常状態

2 : 信頼性劣化状態

3 : データ引き継ぎ元状態 (保護状態)

4 : データ引き継ぎ先状態

5 : 廃棄対象状態

の状態であることを示す。

【 0 2 9 0 】

メタデータサーバ 6 0 0 3 は、第 1 の実施形態と同様の O S 1 0 0 と制御部 2 0 0 とメ

50

タデータ300と論理ドライブステータステーブル450を格納する主メモリ6と、OS100や制御部200を実行する演算装置であるCPU5を具備し、第1の実施形態の情報処理装置111のOS100や制御部200に相当する機能を担う。メタデータサーバ6003の主メモリ6に格納されるメタデータ300はたとえば図33に示したような構造であり、主メモリ6上のメタデータやメタデータのジャーナルは、メタデータサーバ6003内の不揮発性記憶部や、ストレージネットワーク1000内の不揮発性記憶部や、記憶部2、記憶部3などに退避保存される。また、論理ドライブステータステーブル450は、例えば図40に示したような構造であり、論理ドライブが通常状態かデータ引き継ぎ状態であることを示す。メタデータサーバ6003は、たとえばサーバ6001からファイルIDを受信すると、メタデータ領域300からファイルIDを検索し、該ファイルIDに対応する論理ドライブアドレスや、記憶部アドレスであるディスク識別IDやLBAやセクタカウントなどを特定する。記憶部アドレスであるディスク識別IDにはたとえばIP(Internet Protocol)アドレスやMAC(Media Access Control)アドレスやWWN(World Wide Name)を用いてもよい。ステータステーブル510は、ストレージネットワーク1000に接続される各記憶部2、3のステータスを記憶する。

10

【0291】

なお、論理ドライブステータステーブル450を図41に示した構造とする他に、メタデータサーバ6003の主メモリ6に図75に示すようなステータステーブル650を格納し、論理ドライブステータステーブル450を格納しないようにしてもよい。このステータステーブル650においては、論理ドライブアドレスと、ディスク識別IDと、論理ドライブステータスと、ディスクステータスを管理するようにしており、ステータステーブル650は論理ドライブステータステーブル450およびステータス格納領域510の代わりに使用される。この場合は、記憶部2、3でのステータス格納領域510は不要となる。図75では、論理ドライブアドレスBは、ディスク識別IDがb1とb2の2つの記憶部を有し、これら記憶部がデータ引き継ぎ状態で、b1のディスク識別IDを有する記憶部がデータ引き継ぎ元で、b2のディスク識別IDを有する記憶部がデータ引き継ぎ先であることが示されている。

20

【0292】

ストレージネットワーク1000は、ストレージアクセスのためのネットワークであり、たとえばファイバチャネル(Fibre Channel)やイーサネット(登録商標)が用いられる。特に、ストレージネットワーク1000として、例えばSAN(Storage Area Network)やNAS(Network Attached Storage)が用いられる。SANには、たとえば、FC-SAN(Fibre Channel Storage Area Network)やIP-SAN(Internet Protocol Storage Area Network)が用いられ、その上位プロトコルとして、たとえば、SCSI(Small Computer System Interface)が用いられる。本実施形態では、ストレージネットワーク1000としてIP-SANを採用した例を示し、その上位プロトコルとしてはiSCSIが利用されるものとする。ストレージネットワーク1000にはネットワークスイッチ6009やハブ(図示せず)が含まれる。

30

【0293】

ネットワーク6000は、クライアントがサーバ6001にアクセスして種々のサービスを利用するためのネットワークであり、たとえばファイバチャネル(Fibre Channel)やイーサネット(登録商標)が用いられる。ネットワーク6000として、たとえばWANやLANなどがある。ネットワーク6000にはネットワークスイッチ(図示せず)やハブ(図示せず)が含まれる。

40

【0294】

クライアント1002は、ネットワーク6000に接続され、所望の処理を遂行するコンピュータである。クライアント1002は、典型的には、プロセッサと、メインメモリと、通信インタフェースと、ローカル入出力装置等のハードウェア資源を備え、また、デバイスドライバやオペレーティングシステム(OS)、アプリケーションプログラム400等のソフトウェア資源を備える(図示せず)。これによって、クライアント1002は

50

、プロセッサの制御の下、各種のプログラムを実行して、ハードウェア資源との協働作用により処理を実現する。例えば、クライアント6002は、プロセッサの制御の下、業務アプリケーションプログラムを実行することにより、ネットワーク6000を経由してサーバ6001にI/Oアクセスし、所望の業務システムを実現する。

【0295】

サーバ6001は、ストレージネットワーク1000とネットワーク6000に接続され、クライアント1002からの要求に応じて所望の処理を遂行するコンピュータである。サーバ6001は、典型的には、プロセッサと、メインメモリと、通信インタフェースと、ローカル入出力装置等のハードウェア資源を備え、また、デバイスドライバやオペレーティングシステム(OS)、アプリケーションプログラム等のソフトウェア資源を備える(図示せず)。これによって、サーバ6001は、プロセッサの制御の下、各種のプログラムを実行して、ハードウェア資源との協働作用により処理を実現する。例えば、サーバ6001は、クライアント1002からの要求に応じて、プロセッサの制御の下、アプリケーションサービスプログラムを実行することにより、ストレージネットワーク6000を経由して記憶部2または記憶部3または記憶部6004にI/Oアクセスし、所望のアプリケーションサービスプログラムを実現する。

【0296】

サーバ6001は、たとえば、クライアント6002からのファイルデータ読み出し要求とファイルIDを受信すると、メタデータサーバ6003にファイルIDを送信し、メタデータサーバ6003からファイルが格納されている記憶部のIPアドレスやMACアドレスやWWNなどの記憶部アドレスとLBAを受け取り、記憶部アドレスを指定したパケットをストレージネットワーク1000に送信することで記憶部に読み出し命令を送信し、記憶部から読み出しデータを受信し、クライアント1002に読み出しデータを送信する。サーバ6001は、たとえば、クライアント1002からのファイルのデータ書き込み要求とファイルIDを受信すると、メタデータサーバ6003にファイルIDを送信し、メタデータサーバ6003からファイルを格納すべき記憶部の記憶部アドレスとLBAを受け取り、クライアント6002から書き込みデータを受信し、記憶部にデータを送信することでデータを書き込む。サーバ6001は、データベース管理システム(Database management system, DBMS)が稼働しているサーバであるデータベースサーバ(DBサーバ)であってもよい。

【0297】

(寿命到達処理)

制御部200は起動されると、記憶部2についての前述した統計情報65などの信頼性情報を監視する。制御部200は、例えば図24に示すように、一定時間おき(たとえば1分おき)または一定処理ごとに記憶部2から統計情報65を取得する。制御部200は、第1の実施形態と同様、取得した信頼性情報に基づいて、記憶部2が寿命に到達したか否かを判定し、記憶部2が寿命に到達したと判定された場合、後述する寿命到達時処理を実行する。

【0298】

(寿命到達時処理)

制御部200は、接続されている記憶部2が寿命到達するかあるいは寿命到達直前になって寿命到達時処理が開始されると、記憶部2のステータス510をディスクステータス=2(信頼性劣化状態)に変更した後、新しいディスクの接続を促すメッセージをメタデータサーバ6003やサーバ6001やクライアント6002などのディスプレイ9などに表示する。新たな記憶部3が接続されると、制御部200は、ステータステーブル450の記憶部3のステータス510をディスクステータス=4(データ引き継ぎ先状態)に書き換え、さらに記憶部2のステータス510をディスクステータス=3(保護状態)に書き換える。そして、制御部200は、記憶部2+記憶部3を1つの論理ドライブ4としてOS100に認識させる。制御部200は、主メモリ6上の論理ドライブステータステーブル450またはステータステーブル650に格納された論理ドライブ4のステータス

を、ステータスが「通常状態」から「データ引き継ぎ状態」になるように書き換える。

【0299】

(データ引き継ぎ状態での論理ドライブからの読み出し)

サーバ6001は、クライアント6002からの読み出し要求とファイルIDを受信すると、メタデータサーバ6003にファイルIDを送信する。メタデータサーバ6003は、メタデータ300からファイルIDに対応する論理ドライブ4を特定し、さらに主メモリ6から論理ドライブステータステーブル450またはステータステーブル650を読み出し、特定した論理ドライブ4のステータスがデータ引き継ぎ状態であることを認識する。メタデータサーバ6003は、メタデータ300からファイルIDで指定されたファイルが格納されている記憶部2または3の記憶部アドレスとLBAを取得し、取得した記憶部アドレスとLBAをサーバ6001に送信する。サーバ6001は、受信した記憶部アドレスおよびLBAを指定したパケットをストレージネットワーク1000に送信することで記憶部2または3に読み出し命令を送信し、記憶部から読み出しデータを受信し、クライアント1002に読み出しデータを送信する。

10

【0300】

(論理ドライブへの書き込み)

サーバ6001は、たとえば、クライアント1002からのファイルのデータ書き込み要求とファイルIDを受信すると、メタデータサーバ6003にファイルIDを送信する。メタデータサーバ6003は、論理ドライブステータステーブル450またはステータステーブル650から論理ドライブ4のステータスを判定し、論理ドライブ4が通常状態であることを認識すると、主メモリ6からメタデータ300を読み出し、メタデータ300を参照してデータ書き込み用のLBAを割り当てる。メタデータサーバ6003は、LBAおよび記憶部2の記憶部アドレスとLBAをサーバ6001に送信する。サーバ6001は、受信した記憶部アドレスおよびLBAを指定したパケットをストレージネットワーク1000に送信することで記憶部2に書き込み命令を送信し、記憶部2に書き込みデータを記憶する。制御部200はメタデータ300を書き換え、書き込みデータのLBAおよびセクタカウントを、記憶部2および書き込みファイルIDを紐付ける。

20

【0301】

制御部200は、論理ドライブ4がデータ引き継ぎ状態であることを認識すると、主メモリ6からメタデータ300を読み出し、メタデータ300を参照してデータ書き込み用のLBAを割り当てる。メタデータサーバ6003は、LBAおよび記憶部3の記憶部アドレスとLBAをサーバ6001に送信する。サーバ6001は、受信した記憶部3の記憶部アドレスおよびLBAを指定したパケットをストレージネットワーク1000に送信することで記憶部2に書き込み命令を送信し、記憶部2に書き込みデータを記憶する。制御部200はメタデータ300を書き換え、書き込みデータのLBAおよびセクタカウントの記憶部2への紐付けを無効化すると共に、書き込みデータのLBAおよびセクタカウントを、記憶部3および書き込みファイルIDを紐付けることで、記憶部3への書き込みを利用して記憶部2から記憶部3へのデータ引き継ぎを実現する。

30

【0302】

(バックグラウンドバックアップ)

制御部200は、論理ドライブステータステーブル450において論理ドライブ4がデータ引き継ぎ状態のステータスである場合、クライアント6002による論理ドライブ4へのアクセスがほとんど発生しない時に(アイドル時に)、データ引き継ぎ元の記憶部2からデータ引き継ぎ先の記憶部3にバックグラウンドバックアップを行うようにしてもよい。制御部200は主メモリ6からメタデータ300を読み出し、記憶部2に紐付けられているファイルIDを検索し、記憶部2に紐付けられているファイルが存在すれば、サーバ6001を介して当該ファイルのLBAに対して記憶部2に読み出し命令を送信し、記憶部3の当該LBAに対して書き込み命令と読み出されたデータを送信して、書き込みを行い、主メモリ6上のメタデータ300を書き換えて、当該ファイルIDと記憶部2の間の紐付けを無効化すると共に、当該ファイルIDを記憶部3に紐付ける。

40

50

【 0 3 0 3 】

制御部 2 0 0 のバックグラウンドバックアップ動作としてサーバフリーバックアップを採用してもよく、その場合たとえば拡張コピーコマンドを用いてもよい。拡張コピーコマンドとして、たとえばSCSI Primary Commands- 4 (SPC- 4), INCITS T 1 0 / 1 7 3 1 -D, Revision 3 6 e (<http://www.t10.org/>)に記述されている 8 3 h EXTENDED COPYコマンドを用いてもよい。制御部 2 0 0 はバックアップ対象 L B A と記憶部 3 のアドレスとを含む拡張コピーコマンドを記憶部 2 に送信すると、記憶部 2 が該 L B A からデータを読み出し、記憶部 2 が該読み出しデータを記憶部 3 に送信し、記憶部 3 は受信データを該 L B A に書き込む。

【 0 3 0 4 】

(データ引き継ぎ完了時)

論理ドライブステータステーブル 4 5 0 において、論理ドライブ 4 のステータスが「データ引き継ぎ状態」の時、制御部 2 0 0 は、主メモリ 6 上のメタデータ 3 0 0 を定期的に読み出し、記憶部 2 に紐付けられている引き継ぎ対象ファイル I D が存在するか否かを定期的にチェックする。たとえば、制御部 2 0 0 は、論理ドライブ 4 に格納されている全ファイルのファイル I D のうち、記憶部 2 に紐付けられている引き継ぎ対象ファイル I D が存在するか否かを定期的にチェックする。存在しない場合には、制御部 2 0 0 はデータ引き継ぎ先である記憶部 3 のステータス 5 1 0 をディスクステータス = 1 (通常状態) に書き換え、データ引き継ぎ元である記憶部 2 のステータス 5 1 0 をディスクステータス = 5 (廃棄対象状態) に書き換える。制御部 2 0 0 は、記憶部 2 を論理ドライブ 4 から切り離し、記憶部 3 を論理ドライブ 4 として認識し、論理ドライブステータステーブル 4 5 0 またはステータステーブル 6 5 0 内の論理ドライブ 4 のステータスを「データ引き継ぎ状態」から「通常状態」に書き換える。

【 0 3 0 5 】

このように、データ引き継ぎ状態においては、論理ドライブ 4 へのデータ書き込みに際して、データ引き継ぎ元の記憶部 2 に対して書き込み命令を送信せず、データ引き継ぎ先の記憶部 3 に対して書き込み命令を送信する。論理ドライブ 4 からのデータ読み出しは、記憶部 2 , 3 の何れかから実行される。論理ドライブ 4 においては、クライアント 6 0 0 2 から書き込み要求を受信するごとに、記憶部 2 に格納される有効なデータが記憶部 3 に徐々に引き継がれていくことになり、新規データ書き込み動作がデータ引き継ぎ動作を兼用することになる。

【 0 3 0 6 】

(第 5 の実施形態)

第 5 の実施形態では、複数のデータセンタおよびデータセンタ間を接続する長距離ネットワークからなる情報処理システム 1 に本発明を適用した場合について説明する。本実施形態においては、情報処理システム 1 は図 7 6 に示すように構成される。本実施形態では、情報処理システム 1 は、データ引き継ぎ元のデータセンタ 3 0 0 2、データ引き継ぎ先のデータセンタ 3 0 0 3、その他のデータセンタであるデータセンタ 3 0 0 5、サーバ 3 0 0 6、データセンタ管理サーバ 3 0 0 7、これらを接続する長距離ネットワーク 3 0 0 0 からなる。本実施形態では、データセンタ管理サーバ 3 0 0 7 はデータセンタ 3 0 0 2 を論理単位としての論理データセンタ 3 0 0 4 として認識し、データセンタ 3 0 0 2 が信頼性劣化するとデータセンタ 3 0 0 2 とデータセンタ 3 0 0 3 とを論理データセンタ 3 0 0 4 として認識する。

【 0 3 0 7 】

データセンタ 3 0 0 2、3 0 0 3 は、ステータス格納領域 5 1 0 を具備する。記憶部 3 0 0 2、3 0 0 3 には、第 1 の実施形態のデータ引き継ぎ対象記憶部アドレス格納領域 2 0 0 7、3 0 0 7 (図 3 2 参照) に対応するデータ引き継ぎ対象データセンタ I D 格納領域 2 0 0 7、3 0 0 7 を具備させてもよいし、第 2 の実施形態の論理ドライブ I D 格納領域 5 2 0 に対応する論理データセンタ I D 格納領域 5 2 0 を具備させてもよい。ステータス格納領域 5 1 0 は、前述と同様、対応する記憶部がそれぞれ

10

20

30

40

50

- 0 : 初期データセンタ状態
- 1 : 通常状態
- 2 : 信頼性劣化状態
- 3 : データ引き継ぎ元状態 (保護状態)
- 4 : データ引き継ぎ先状態
- 5 : 廃棄対象状態

の状態であることを示す。

【0308】

データセンタ管理サーバ3007は、第1の実施形態と同様のOS100と制御部200とメタデータ300と論理データセンタステータステーブル450を格納する主メモリ6と、OS100や制御部200を実行する演算装置であるCPU5を具備し、第1の実施形態の情報処理装置111のOS100や制御部200に相当する機能を担う。メタデータサーバ6003の主メモリ6に格納されるメタデータ300はたとえば図33に示したような構造であり、主メモリ6上のメタデータやメタデータのジャーナルは、データセンタ管理サーバ3007内の不揮発性記憶部や、長距離ネットワーク3000内のデータセンタ3002、3003、3005などに退避保存される。また、論理データセンタステータステーブル450は、例えば図41に示したような構造であり、論理データセンタ3004が通常状態かデータ引き継ぎ状態であることを示す。データセンタ管理サーバ3007は、たとえばサーバ3006からファイルIDを受信すると、メタデータ領域300からファイルIDを検索し、該ファイルIDに対応する論理データセンタIDや、データセンタIDやLBAやセクタカウントなどを特定する。ステータステーブル510は、データセンタ3002、3003のステータスを記憶する。

【0309】

なお、論理データセンタステータステーブル450を図41に示した構造とする他に、データセンタ管理サーバ3007に図75に示したようなステータステーブル650を採用するようにしてもよい。このステータステーブル650においては、論理データセンタIDと、データセンタIDと、論理データセンタステータスと、データセンタステータスを管理するようにしており、この場合は、記憶部2、3でのステータス格納領域510は不要となる。制御部200は、メタデータ領域300の代わりに、第2の実施例のようにデータ引き継ぎ履歴格納領域を用いてデータセンタ3002からデータ3003へのデータ引き継ぎを管理するようにしてもよい。

【0310】

(寿命到達処理)

制御部200は起動されると、データセンタ3002についての前述した信頼性情報を監視する。制御部200は、例えば図24に示すように、一定時間おき(たとえば1分おき)または一定処理ごとにデータセンタ3002から統計情報65を取得する。制御部200は、第1の実施形態と同様、取得した統計情報65に基づいて、データセンタ3002が寿命に到達したか否かを判定し、データセンタ3002が寿命に到達したと判定された場合、後述する寿命到達時処理を実行する。

【0311】

(寿命到達時処理)

制御部200は、接続されているデータセンタ3002が寿命到達するかあるいは寿命到達直前になって寿命到達時処理が開始されると、データセンタ3002のステータス510をデータセンタステータス=2(信頼性劣化状態)に変更した後、新しいデータセンタの接続を促すメッセージをデータセンタ管理サーバ3007のディスプレイなどに表示する。新たなデータセンタ3003が接続されると、制御部200は、ステータステーブル450のデータセンタ3003のステータス510をデータセンタステータス=4(データ引き継ぎ先状態)に書き換え、さらにデータセンタ3002のステータス510をデータセンタステータス=3(保護状態)に書き換える。そして、制御部200は、データセンタ3002+データセンタ3003を1つの論理データセンタ3004としてOS1

10

20

30

40

50

00に認識させる。制御部200は、主メモリ6上の論理データセンタステータステーブル450を、ステータスが「通常状態」から「データ引き継ぎ状態」になるように書き換える。

【0312】

(データ引き継ぎ中での論理データセンタからの読み出し)

データセンタ管理サーバ3007は、サーバ3006からの読み出し要求とファイルIDを受信すると、メタデータ300からファイルIDに対応する論理データセンタ3004を特定し、さらに主メモリ6から論理データセンタステータステーブル450を読み出し、特定した論理データセンタ3004のステータスがデータ引き継ぎ状態であることを認識する。データセンタ管理サーバ3007は、メタデータ300からファイルIDで指定されたファイルが格納されているデータセンタ3002または3003のデータセンタアドレスとLBAを取得し、取得したデータセンタアドレスとLBAをサーバ3006に送信する。サーバ3006は、受信したデータセンタアドレスおよびLBAを指定したパケットをネットワーク3000に送信することでデータセンタ3002または3003に読み出し命令を送信し、データセンタ3002または3003から読み出しデータを受信する。

10

【0313】

(論理データセンタへの書き込み)

データセンタ管理サーバ3007は、サーバ3006からの書き込み要求とファイルIDを受信すると、論理データセンタステータステーブル450から論理データセンタ3004のステータスを判定し、論理データセンタ3004が通常状態であることを認識すると、主メモリ6からメタデータ300を読み出し、メタデータ300を参照してデータ書き込み用のLBAを割り当てる。データセンタ管理サーバ3007は、LBAおよびデータセンタ3002のデータセンタIDとLBAをサーバ3006に送信する。サーバ3006は、受信したデータセンタIDおよびLBAを指定したパケットをネットワーク3000に送信することでデータセンタ3002に書き込み命令を送信し、データセンタ3002に書き込みデータを記憶する。

20

【0314】

データセンタ管理サーバ3007は、論理データセンタ3004がデータ引き継ぎ状態であることを認識すると、主メモリ6からメタデータ300を読み出し、メタデータ300を参照してデータ書き込み用のLBAを割り当てる。データセンタ管理サーバ3007は、LBAおよびデータセンタ3003のデータセンタIDとLBAをサーバ3006に送信する。サーバ3006は、受信したデータセンタ3003のデータセンタIDおよびLBAを指定したパケットをネットワーク3000に送信することでデータセンタ3003に書き込み命令を送信し、データセンタ3003に書き込みデータを記憶する。

30

【0315】

(バックグラウンドバックアップ)

論理データセンタステータステーブル450において、論理データセンタ4がデータ引き継ぎ中のステータスである場合、サーバ3006による論理データセンタ4へのアクセスがほとんど発生しない時に(アイドル時に)、データ引き継ぎ元のデータセンタ3002からデータ引き継ぎ先のデータセンタ3003にバックグラウンドバックアップを行うようにしてもよい。制御部200は主メモリ6からメタデータ300を読み出し、データセンタ3002に紐付けられているファイルIDを検索し、データセンタ3002に紐付けられているファイルが存在すれば、サーバ3006を介して当該ファイルのLBAに対してデータセンタ3002に読み出し命令を送信してデータ3002から読み出しデータを受信し、データセンタ3003の当該LBAに対して書き込み命令と読み出されたデータを送信して、書き込みを行い、主メモリ6上のメタデータ300を書き換えて当該ファイルIDをデータセンタ3003に紐付ける。

40

【0316】

制御部200のバックグラウンドバックアップ動作としてサーバフリーバックアップを

50

採用してもよく、その場合たとえば拡張コピーコマンドを用いてもよい。拡張コピーコマンドとして、たとえばSCSI Primary Commands-4 (SPC-4), INCITS T10/1731-D, Revision 3.6e (<http://www.t10.org/>)に記述されている83h EXTENDED COPYコマンドを用いてもよい。制御部200はバックアップ対象LBAとのデータセンタ3003のIDとを含む拡張コピーコマンドをデータセンタ3002に送信すると、データセンタ3002が該LBAからデータを読み出し、データセンタ3002が該読み出しデータをデータセンタ3003に送信し、データセンタ3003は受信データを該LBAに書き込む。

【0317】

(データ引き継ぎ完了時)

論理データセンタステータステーブル450において、論理データセンタ4のステータスが「データ引き継ぎ状態」の時、制御部200は、主メモリ6上のメタデータ300を定期的に読み出し、データセンタ3002に紐付けられている引き継ぎ対象ファイルIDが存在するか否かを定期的にチェックする。たとえば、制御部200は、論理データセンタ4に格納されている全ファイルのファイルIDのうち、データセンタ3002に紐付けられている引き継ぎ対象ファイルIDが存在するか否かを定期的にチェックする。存在しない場合には、制御部200はデータ引き継ぎ先であるデータセンタ3003のステータス510をデータセンタステータス=1(通常状態)に書き換え、データ引き継ぎ元であるデータセンタ3002のステータス510をデータセンタステータス=5(廃棄対象状態)に書き換える。制御部200は、データセンタ3002を論理データセンタ4から切り離し、データセンタ3003を論理データセンタ4として認識し、論理データセンタステータステーブル450で論理データセンタ3004のステータスを「データ引き継ぎ状態」から「通常状態」に書き換える。

【0318】

このように、データ引き継ぎ中には、論理データセンタ3004へのデータ書き込みの際し、データ引き継ぎ元のデータセンタ3002に対して書き込み命令を送信せず、データ引き継ぎ先のデータセンタ3003に対して書き込み命令を送信する。論理データセンタ3004からのデータ読み出しは、データセンタ3002、3003の何れかから実行される。論理データセンタ4においては、サーバ3006から書き込み要求を受信すると、データセンタ3002に格納される有効なデータがデータセンタ3003に徐々に引き継がれていくことになり、新規データ書き込みがデータ引き継ぎ動作を兼用することになる。

【0319】

(第6の実施形態)

(中継装置)

第1及び第2の実施形態では、制御部200が、記憶部2を寿命到達したと判定した場合、寿命到達目前であると判定した場合、故障目前であると判定した場合、信頼性劣化したと判定した場合に、記憶部2に対する書き込みを行わないよう書き込み命令を制御することで、データ書き込みによる記憶部2のさらなる信頼性劣化を抑制するとともに、新規書き込みデータの損失を防止することを可能にしている。記憶部2に対するデータ書き込みを安全に制限するという観点では、情報処理装置やクライアントやサーバが記憶部2に対する書き込みを自発的に行わないようにすることが望ましい。本実施形態では、情報処理システムは中継装置5000を具備し、中継装置5000内の制御部200が記憶部2の統計情報などの信頼性情報を監視し、中継装置5000が記憶部2を寿命到達したと判定した場合、寿命到達目前であると判定した場合、故障目前であると判定した場合、信頼性劣化したと判定した場合に、記憶部2が読み出し専用デバイスであるという情報を情報処理装置111やクライアントやサーバに通知することにより、情報処理装置111やクライアントやサーバが記憶部2に対する書き込みを自発的に行わないようにさせている。本実施形態は単独で実施することが可能であるし、一方、第1の実施形態や第2の実施形態等と組み合わせることで、信頼性劣化した記憶部2に対する書き込みをより堅牢に抑制

10

20

30

40

50

することが可能になる。たとえば、本実施形態を第 1 の実施形態や第 2 の実施形態等を組み合わせることで、中継装置は記憶部 2 が読み出しと書き込みのうち読み出しのみをサポートする記憶部であるという記憶部識別情報を送信し、CPU 5 や主メモリ 6 に格納された制御部 200 は、記憶部 2 を読み出し専用記憶部として認識し、記憶部 2 を保護状態（データ引き継ぎ元状態）として認識する。

【0320】

図 77 は、デスクトップパソコンやノートパソコンなどの情報処理システム 1 に中継装置が搭載された例である。中継装置 5000 は、情報処理装置 1 内に搭載されてもよいし、情報処理装置 1 外に搭載されてもよい。中継装置 5000 は、インタフェース 19 を介して記憶部 2 と接続され、インタフェース 5001、チップセット 7 を経由して CPU 5 と接続される。なお、中継装置 5000 は、チップセット 7 を経由せず CPU 5 に直接接続されてもよい。また、中継装置 5000 は、チップセット 7 に含まれていてもよい。

10

【0321】

中継装置 5000 は制御部 200 を具備する。制御部 200 は、図 77 に示すように、全てが中継装置 5000 に含まれてもよいし、図 78 に示すように、一部が中継装置 5000 に含まれ一部が主メモリ 6 に含まれてもよいし、一部が中継装置 5000 に含まれ一部が ROM 11 など情報処理装置 111 内のその他のメモリ部分に含まれてもよい。制御部 200 は、ファームウェアやソフトウェアの形態で実装されてもよいし、ハードウェアの形態で実装されてもよい。

【0322】

20

インタフェース 19 およびインタフェース 5001 として、SATA (Serial Advanced Technology Attachment)、PCI Express (Peripheral Component Interconnect Express, PCIe)、USB (Universal Serial Bus)、SAS (Serial Attached SCSI)、Thunderbolt (登録商標)、イーサネット (登録商標)、ファイバーチャネルなどが使用可能である。インタフェース 19 とインタフェース 5001 は同じ規格のインタフェースであってもよいし、異なる規格のインタフェースであってもよい。本実施形態では、インタフェース 19 とインタフェース 5001 が SATA インタフェースである場合について説明する。

【0323】

制御部 200 は、第 1 の実施形態の図 24 や図 29 と同様にして、記憶部 2 が寿命に到達しているか否か、寿命到達目前であるか否か、あるいは故障目前であるか否かを判定し、寿命に到達しているか場合、寿命到達目前である場合、あるいは故障目前である場合に寿命到達時処理として、図 79 に示すように、通常状態から信頼性劣化モードに遷移する（ステップ S800）。通常状態と信頼性劣化モードは制御部 200 のモードであり、記憶部 2 が正常状態であるとき制御部 200 は通常状態で動作し、記憶部 2 を寿命到達したと判定した場合、寿命到達目前であると判定した場合、故障目前であると判定した場合、信頼性劣化したと判定した場合に、制御部 200 は信頼性劣化モードで動作する。寿命到達判定は、第 1 の実施形態と同様に、図 24 のような一定時間ごと、または一定処理数おき、または一定データ送受信おき、または図 29 や図 30 のように記憶部からのコマンド応答がエラーであった場合に実施される。

30

40

【0324】

制御部 200 は、インタフェース 5001 を経由して CPU 5 から受信した命令・データを、インタフェース 19 を経由して記憶部 2 に送信する。また、制御部 200 は、インタフェース 19 を経由して記憶部 2 から受信した応答・データを、インタフェース 5001 を経由して CPU 5 に送信する。インタフェース 5001 とインタフェース 19 のプロトコルが異なる場合には、制御部 200 がプロトコル変換を行った後に、変換後の命令・応答・データを CPU 5 や記憶部 2 に送信する。制御部 200 は記憶部 2 の記憶部情報を CPU 5 に送信する時、制御部 200 が通常状態であるか信頼性劣化モードであるかによって、CPU 5 に送信する記憶部 2 の記憶部情報を切り替える。すなわち、制御部 200 は、通常状態では記憶部 2 が読み出しおよび書き込み可能な記憶部であるという記憶部情

50

報をCPU5に送信し、信頼性劣化モードでは記憶部2が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報をCPU5に送信する。

【0325】

制御部200は、CPU5から記憶部2の記憶部情報要求を受信すると、記憶部情報要求に対する応答として、通常状態では記憶部2が読み出しおよび書き込み可能な記憶部であるという記憶部情報をCPU5に送信し、信頼性劣化モードでは記憶部2が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報をCPU5に送信する。記憶部情報要求としては、たとえば、ACS-3記載のECh IDENTIFY DEVICEコマンドや、SCSI Primary Commands-4 (SPC-4)記載のA3h REPORT IDENTIFYING INFORMATIONコマンドや、INFORMATION TECHNOLOGY Multi-Media Commands-6 (MMC-6)(<http://www.t10.org/>)に記載の46h GET CONFIGURATIONコマンドや、NVM Express Revision 1.1記載の06h Identifyコマンドなどが用いられる。

10

【0326】

図80は、制御部200がCPU5から記憶部2の記憶部情報要求を受信した場合の処理手順の例を示すものである。制御部200がCPU5から記憶部2の記憶部情報要求を受信すると(ステップS801)、制御部200は、記憶部2が信頼性劣化モードであるか否かを、制御部200自身が通常状態にあるか信頼性劣化モードにあるかに基づいて判定する(ステップS802)。通常状態である場合、制御部200は、記憶部2に記憶部情報要求を送信し(ステップS803)、記憶部2から記憶部情報を受信し、CPU5に受信した記憶部情報を送信する(ステップS804)。信頼性劣化モードである場合、制御部200は、記憶部2に記憶部情報要求を送信し(ステップS805)、記憶部2から記憶部情報を受信し、受信データを書き換えて記憶部2がROMデバイスであるかのように記憶部情報を変更し、CPU5に変更した記憶部情報を送信する(ステップS806)。なお、制御部200は、ステップS801のCPU5からの記憶部情報要求を受信すること無しに、制御部200が自発的にS802～S806の処理を実施してもよい。また、制御部200は、記憶部2に対する記憶部情報要求送信をステップS801とステップS802の間に行い、S803およびS805での記憶部2に対する記憶部情報要求送信を行わないようにしてもよい。

20

【0327】

図81に示すように、制御部200は、信頼性劣化モードのとき、記憶部2に記憶部情報要求を送信することなく記憶部2がROMデバイスであるかのような記憶部情報を生成し、CPU5に送信するようにしてもよい。すなわち、制御部200は、記憶部2の記憶部情報要求をCPU5から受信すると(ステップS810)、記憶部2が信頼性劣化モードであるか否かを、制御部200自身が通常状態にあるか信頼性劣化モードにあるかに基づいて判定する(ステップS811)。通常状態である場合、制御部200は、記憶部2に記憶部情報要求を送信し(ステップS812)、CPU5に受信した記憶部情報を送信する(ステップS813)。信頼性劣化モードである場合、制御部200は、記憶部2に記憶部情報要求を送信することなく、記憶部2がROMデバイスであるとCPU5に見せかける記憶部情報を生成し、CPU5に生成した記憶部情報を送信する(ステップS814)。

30

40

【0328】

通常状態での動作時に、記憶部2が読み出しおよび書き込み可能な記憶部であるという記憶部情報として、制御部200はたとえば記憶部2がATAデバイスであることをCPU5に明示的に通知することが望ましい。たとえば、ATA/ATAPI Command Set-3 (ACS-3)で記述されているDevice Signatureにおいて、LBA(7:0)を01hに、LBA(15:8)を00hに、LBA(23:16)を00hとしてCPU5に出力することで、記憶部2がATAデバイスであることをCPU5に通知することができる。

【0329】

信頼性劣化モードでの動作時に、記憶部2が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報として、制御部200はたとえば記憶部2がATA P

50

I デバイスであることを C P U 5 に明示的に通知する。たとえば、ACS-3で記述されている Device Signature において、L B A (7:0) を 01h に、L B A (15:8) を 14h に、L B A (23:16) を EBh として C P U 5 に出力することで、記憶部 2 が A T A P I デバイスであることを C P U 5 に通知することができる。さらに、信頼性劣化モードでの動作時に、記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報として、制御部 2 0 0 は、たとえば C P U 5 から INCITS Multi-Media Commands-6 (MMC-6) 記載の 46h GET CONFIGURATION コマンドを受信したとき、Random Writable (Feature Number = 0020h)、Incremental Streaming Writable (Feature Number = 0021h)、Write Once (Feature Number = 0025h) などの Features において、書き込み系機能は全て非サポートであることを C P U 5 に返信する。これにより、O S 1 0 0 が Windows (登録商標) などである場合でも、記憶部 2 を読み出し可能なデバイスとして O S 1 0 0 に認識させることが可能となり、O S 1 0 0 や、O S 1 0 0 より上位のレイヤのアプリケーションプログラム 4 0 0 に対しては記憶部 2 があたかもリードオンリーデバイスであるかのように見え、O S 1 0 0 やアプリケーションプログラム 4 0 0 が記憶部 2 に対して誤って書き込み命令を送信してしまうことを防止できる。

10

【0330】

あるいは、信頼性劣化モードでの動作時に、記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報として、制御部 2 0 0 はたとえば記憶部 2 が A T A デバイスであることを C P U 5 に明示的に通知し、かつ C P U 5 から ACS-3 記載のコマンド ECh IDENTIFY DEVICE を受信したとき、書き込み系機能は全て非サポートであるという情報を C P U 5 に通知してもよい。

20

【0331】

記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であるか否かの通知方法はこれ以外の様々な形態をとってもよい。C P U 5 は、記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという情報を受信することにより、O S 1 0 0 は、記憶部 2 に適用されるドライバソフトウェアとして読み出し専用の例えば A T A P I の読み出し専用記憶部用ドライバを適用し、O S 1 0 0 上で記憶部 2 はたとえば C D - R O M や D V D - R O M や B D - R O M などのような読み出し専用記憶部として認識される。たとえば、図 8 2 や、図 8 3 に示すように、情報処理装置 1 1 1 のユーザは、O S 1 0 0 により、ディスプレイ装置 9 経由で記憶部 2 がたとえば C D - R O M や D V D - R O M や B D - R O M などのような読み出し専用デバイスであるとアイコンのグラフィック等により明示的に通知されることが望ましい。さらに望ましくは、ユーザが記憶部 2 にファイル等を書き込もうと試みた場合、図 8 4 のように O S 1 0 0 がディスプレイ装置 9 を経由して記憶部 2 がライトプロテクトされていることを通知することが望ましい。

30

【0332】

O S 1 0 0 が記憶部 2 に読み出し専用記憶部用ドライバを適用した場合など、C P U 5 や O S 1 0 0 が記憶部 2 を読み出しと書き込みのうち読み出しのみ可能な記憶部として認識した場合であっても、C P U 5 は記憶部 2 に対して読み出し命令を送信可能である。C P U 5 が記憶部 2 に読み出し命令を送信した時、制御部 2 0 0 は記憶部 2 に読み出し命令を送信し、制御部 2 0 0 は記憶部 2 から読み出しデータを受信し、制御部 2 0 0 は読み出しデータを C P U 5 に送信する。このようにして、C P U 5 は、制御部 2 0 0 が通常状態と信頼性劣化モードのどちらの場合であっても、記憶部 2 からデータを読み出すことができる。

40

【0333】

記憶部 2 のデータ破壊や故障によるユーザデータ損失を防止する観点では、信頼性劣化モードにおける制御部 2 0 0 は、記憶部 2 に対する書き込みコマンドを一切送信しないよう構成されることが望ましい。しかし、制御部 2 0 0 は、たとえば O S 1 0 0 のシステム情報など一部のデータを記憶部 2 に書き込む必要がある場合は、例外として記憶部 2 に当該データを書き込むことを許可してもよいが、当該書き込みデータはデータ量が記憶部 2 の容量に対して十分に小さいことが望ましい。さらに望ましくは、ユーザが誤って書き込

50

みコマンドを送信してしまい記憶部 2 に対してデータを書き込んでしまうことを防止するために、記憶部 2 に対する通常の書き込みコマンドを C P U 5 から受信しても記憶部 2 に書き込みコマンドを送信せず、記憶部 2 に対して例外的にデータを書き込む必要がある場合として、特殊なコマンドを用いた書き込みコマンドを C P U 5 から受信した場合に限り、記憶部 2 に対して書き込みコマンドを送信するようにすることが望ましい。たとえば、C P U 5 の記憶部 2 に対する通常の書き込みコマンドとして、たとえば ACS-3 記載の 35h W R I T E D M A E X T や 61h W R I T E F P D M A Q U E U E D のような書き込みコマンドが使われる場合、通常状態の制御部 2 0 0 が 35h W R I T E D M A E X T コマンドや 61h W R I T E F P D M A Q U E U E D コマンドを通常書き込みコマンドとして C P U 5 から受信すると、通常状態の制御部 2 0 0 は該コマンドを記憶部 2 に転送し、信頼性劣化モードの制御部 2 0 0 が 35h W R I T E D M A E X T コマンドや 61h W R I T E F P D M A Q U E U E D コマンドを通常書き込みコマンドとして C P U 5 から受信すると、信頼性劣化モードの制御部 2 0 0 は書き込みコマンドを記憶部 2 に送信しない。一方、信頼性劣化モードの制御部 2 0 0 がたとえば INCITS ACS-3 に記述されている 30h W r i t e S e c t o r s コマンドや 3Fh W r i t e L o g E x t コマンドや SCT Command Transport やその他ベンダー独自のコマンドなどを特殊書き込みコマンドとして C P U 5 から受信すると、信頼性劣化モードの制御部 2 0 0 は該コマンドを記憶部 2 に転送する。

10

【 0 3 3 4 】

以上、情報処理システム 1 がデスクトップパソコンやノートパソコンなどのコンピュータシステムである場合について説明したが、情報処理システム 1 はたとえば図 8 5 のようなストレージネットワークを有する情報処理システムであってもよい。図 8 5 では、インタフェース 1 9 やインタフェース 5 0 0 1 としてのストレージネットワーク 1 0 0 0 が用いられ、中継装置 5 0 0 0 としてのネットワークスイッチ 6 0 0 9 が用いられる。

20

【 0 3 3 5 】

ストレージネットワーク 1 0 0 0 として、例えば S A A N (Storage Area Network) や N A S (Network Attached Storage) が用いられる。S A N には、たとえば、F C - S A N (Fibre Channel Storage Area Network) や I P - S A N (Internet Protocol Storage Area Network) が用いられ、その上位プロトコルとして、たとえば、S C S I (Small Computer System Interface) が用いられる。たとえば、ストレージネットワーク 1 0 0 0 として I P - S A N が採用可能であり、その上位プロトコルとしては i S C S I (Internet Small Computer System Interface) が採用可能である。

30

【 0 3 3 6 】

ネットワークスイッチ 6 0 0 9 はストレージネットワーク 1 0 0 0 上においてクライアントやサーバや記憶部など複数のネットワーク機器間の接続を行うネットワークデバイスであり、ネットワーク機器からパケットを受信すると、受信パケットの宛先アドレスをもとに宛先となるネットワーク機器に受信パケットを送信する。

【 0 3 3 7 】

ネットワークスイッチ 6 0 0 9 は制御部 2 0 0 を具備する。制御部 2 0 0 は、図 8 5 に示すように全てがネットワークスイッチ 6 0 0 9 に含まれても良いし、一部がネットワークスイッチ 6 0 0 9 に含まれ一部がたとえばクライアント 7 0 0 0 A の主メモリ 6 に含まれても良いし、一部がネットワークスイッチ 6 0 0 9 に含まれ一部がたとえばクライアント 7 0 0 0 A の R O M などシステム 1 内のその他の部分に含まれても良い。制御部 2 0 0 は、ファームウェアやソフトウェアの形態で実装されてもよいし、ハードウェアの形態で実装されてもよい。

40

【 0 3 3 8 】

制御部 2 0 0 は、第 1 の実施形態の図 2 4 や図 2 9 と同様にして、ストレージネットワーク 1 0 0 0 に接続された 1 乃至複数の記憶部が寿命に到達しているか否か、寿命到達目前であるか否か、あるいは故障目前であるか否かを判定し、寿命に到達しているか場合、寿命到達目前である場合、あるいは故障目前である場合に寿命到達時処理として、寿命到達時処理対象の記憶部のみを対象として通常状態から信頼性劣化モードに遷移する。通常状態と信頼性劣化モードは、ストレージネットワーク 1 0 0 0 に接続された 1 乃至複数の

50

記憶部それぞれに対応して存在する制御部 200 のモードであり、たとえば記憶部 2A が正常状態であるとき制御部 200 の記憶部 2A 用モードは通常状態で動作し、たとえば記憶部 2A を寿命到達したと判定した場合、寿命到達目前であると判定した場合、故障目前であると判定した場合、信頼性劣化したと判定した場合に、制御部 200 は記憶部 2A 用モードとして信頼性劣化モードで動作する。制御部 200 は、寿命到達時処理対象の記憶部を対象として通常状態から信頼性劣化モードに遷移した場合であっても、寿命到達時処理対象ではない記憶部に対しては通常状態で動作する。寿命到達判定は、第 1 の実施形態と同様に、図 24 のような一定時間ごと、または一定処理数おき、または一定データ送受信おき、または図 29 や図 30 のように記憶部からのコマンド応答がエラーであった場合に実施される。

10

【0339】

本実施形態では例としてストレージネットワーク 1000 に、2 つのクライアント 7000A と 7000B と、2 つの記憶部 2A と 2B が接続された場合について説明するが、システム 1 のネットワーク機器の構成は他の任意の構成を採用しても良い。また、クライアント 7000A の代わりにサーバ 7000A が用いられてもよいし、クライアント 7000B の代わりにサーバ 7000B が用いられてもよい。ネットワーク機器であるクライアント 7000A、クライアント 7000B、記憶部 2A、記憶部 2B にはそれぞれストレージネットワーク上でアドレスが割り当てられる。アドレスとして、たとえば IP (Internet Protocol) アドレスや MAC (Media Access Control) アドレスを使用することができる。たとえば、クライアント 7000A の CPU 5 が、記憶部 2A のアドレスのみを指定する命令およびデータをストレージネットワーク 1000 に送信すると、ネットワークスイッチ 6009 がこの命令やデータを構成するパケット内のアドレスから記憶部 2A を特定し、ネットワークスイッチ 6009 はパケットを記憶部 2A にのみ送信する。たとえば、記憶部 2A が、クライアント 7000A のアドレスのみを指定する応答およびデータをストレージネットワーク 1000 に送信すると、ネットワークスイッチ 6009 が前記応答やデータを構成するパケット内のアドレスからクライアント 7000A を特定し、ネットワークスイッチ 6009 はそのパケットをクライアント 7000A にのみ送信する。このアドレス指定は、単一のネットワーク機器のみでなく、複数のネットワーク機器を指定することもできる。

20

【0340】

制御部 200 はたとえば記憶部 2A の記憶部情報をたとえばクライアント 7000A に送信する時、制御部 200 の記憶部 2A 用モードが通常状態であるか信頼性劣化モードであるかによって、クライアント 7000A に送信する記憶部 2A の記憶部情報を切り替える。すなわち、制御部 200 は、制御部 200 の記憶部 2A 用モードが通常状態のとき記憶部 2A が読み出しおよび書き込み可能な記憶部であるという記憶部情報をクライアント 7000A に送信し、制御部 200 の記憶部 2A 用モードが信頼性劣化モードのとき記憶部 2A が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという記憶部情報をクライアント 7000A に送信する。記憶部 2A の記憶部情報は同様にしてクライアント 7000B にも送信される。また、記憶部 2B の記憶部情報も同様にしてクライアント 7000A やクライアント 7000B にも送信される。

30

40

【0341】

このように本実施形態では、制御部 200 は、記憶部 2 を寿命到達したと判定した場合、寿命到達目前であると判定した場合、故障目前であると判定した場合、信頼性劣化したと判定した場合に、CPU 5 に送信される記憶部 2 の記憶部情報を加工または生成することにより、記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であると CPU 5 や OS 100 に認識させることで、記憶部 2 に対するデータ書き込みを防止することができる。また、記憶部 2 が読み出しと書き込みのうち読み出しのみ可能な記憶部であるという認識処理は OS 100 など下位のソフトウェアレイヤで行われるため、アプリケーションプログラム 400 など上位のソフトウェアレイヤやユーザは記憶部 2 への特別な制御を必要としない。

50

【 0 3 4 2 】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

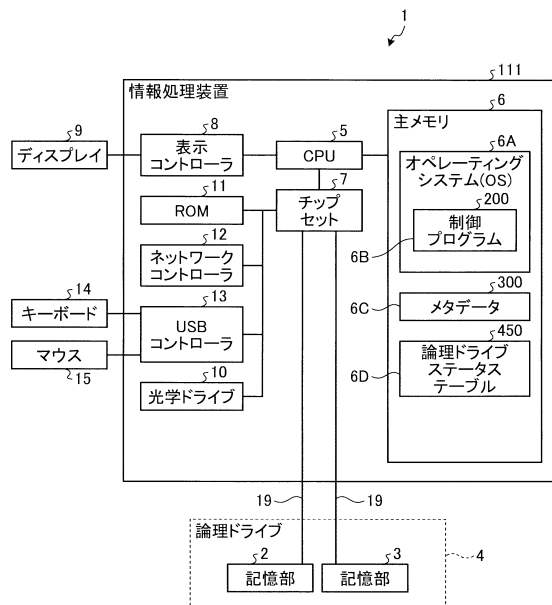
【 符号の説明 】

【 0 3 4 3 】

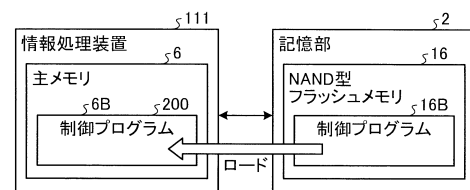
1 情報処理システム、2, 2A-2D 記憶部(記憶装置)、3 記憶部(記憶装置)、4 論理ドライブ、6 主メモリ、16 NAND型フラッシュメモリ、19 インタフェース、111 情報処理装置、200 制御プログラム(制御部)、300 メタデータ、400 アプリケーションプログラム、450 論理ドライブステータステーブル、510 ステータス記憶領域(ステータステーブル)、520 論理ドライブID記憶領域、530 スロット番号記憶領域、550 引き継ぎ履歴(引き継ぎ履歴記憶領域)、1000 ストレージネットワーク、1001 ネットワークスイッチ、1003 ストレージアレイ装置。

10

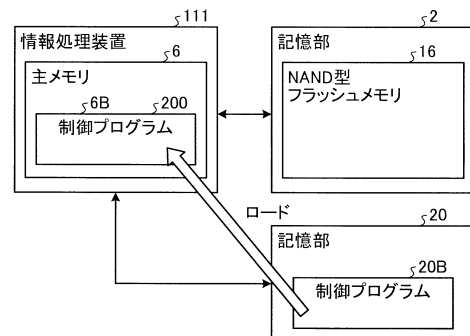
【 図 1 】



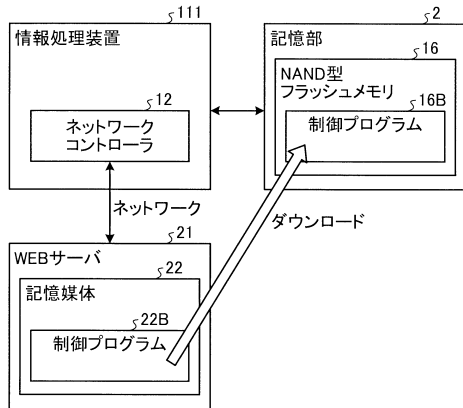
【 図 2 】



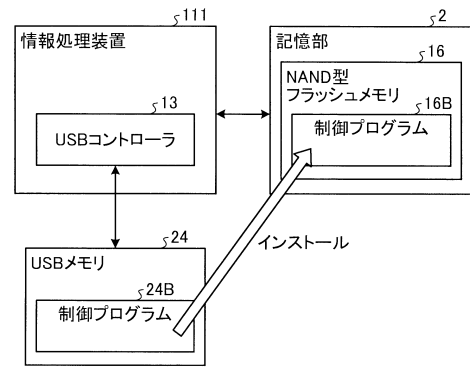
【 図 3 】



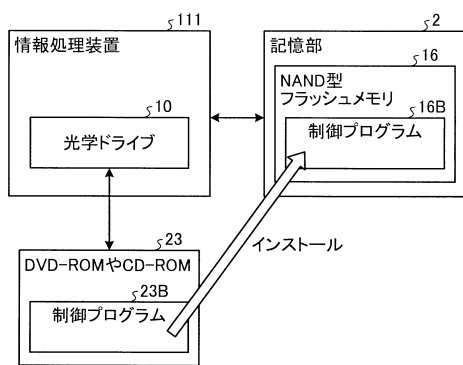
【図 4】



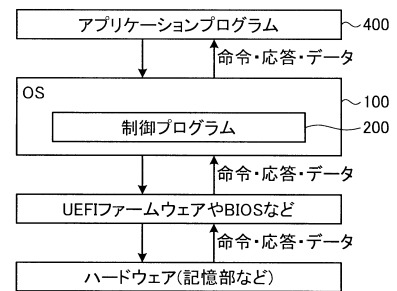
【図 6】



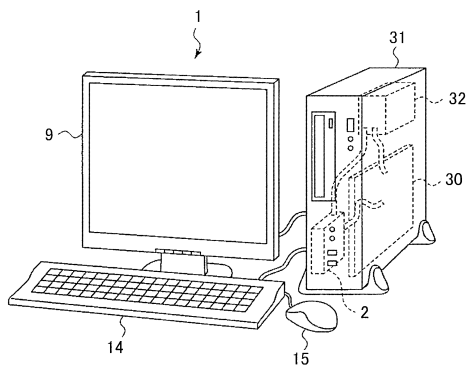
【図 5】



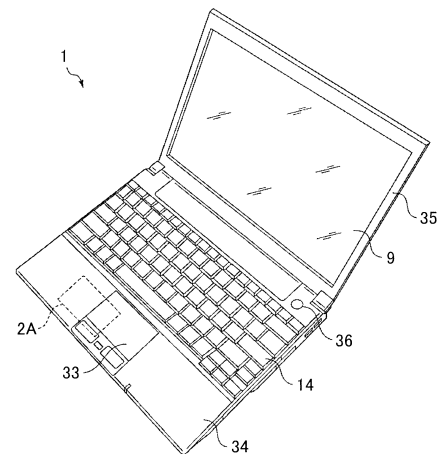
【図 7】



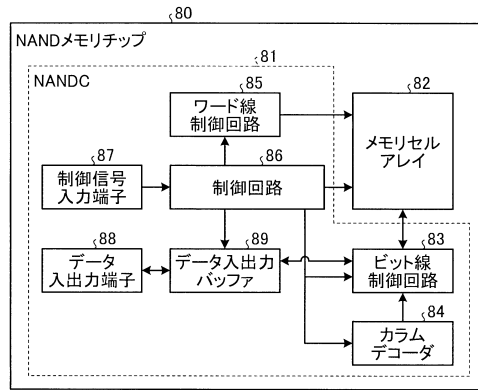
【図 8】



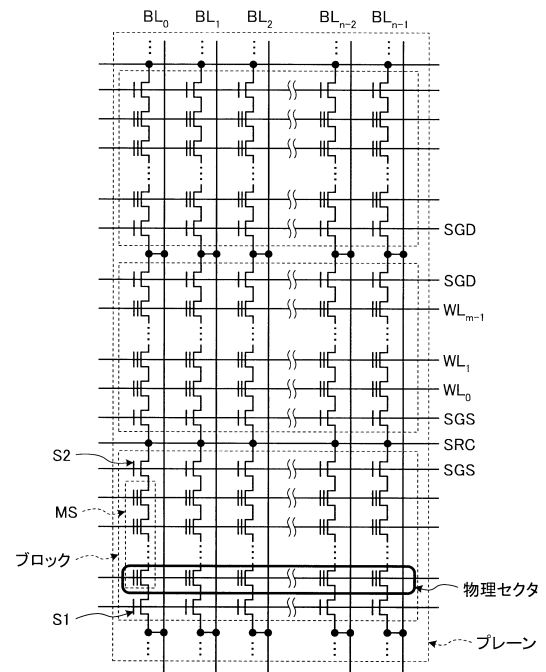
【図 9】



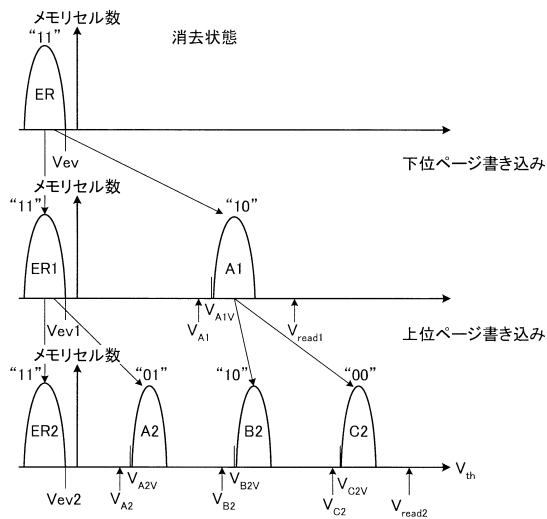
【図 10】



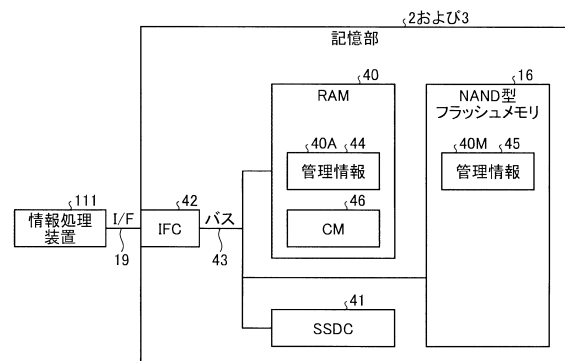
【図 11】



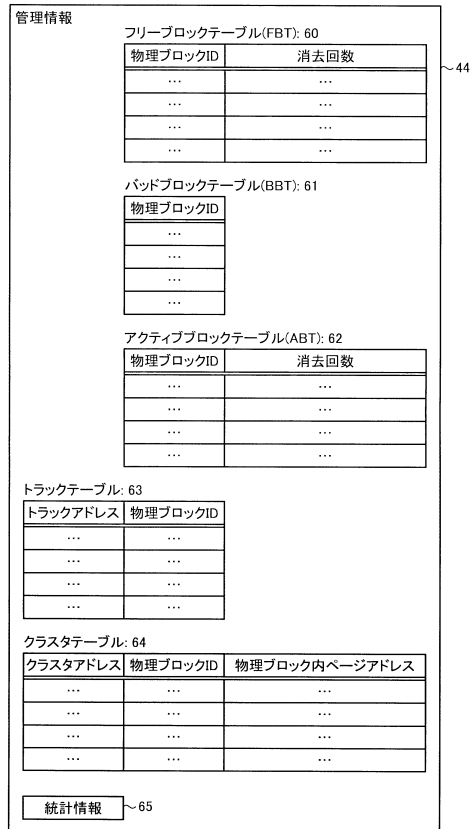
【図 12】



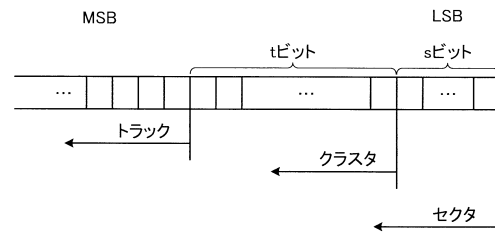
【図 13】



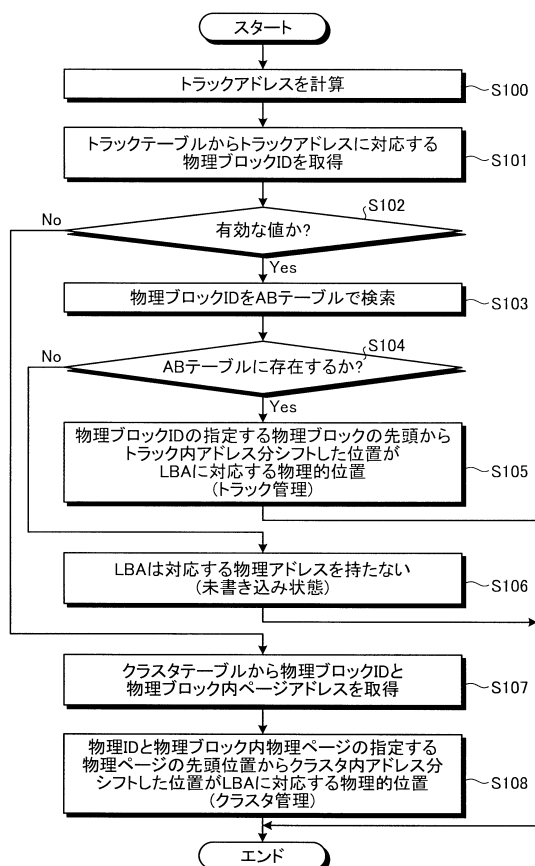
【図 14】



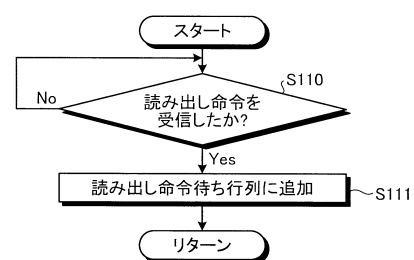
【図 15】



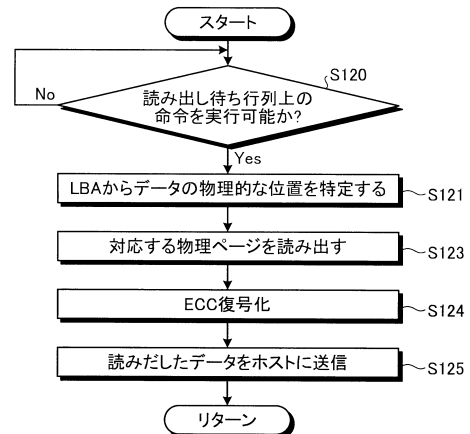
【図 16】



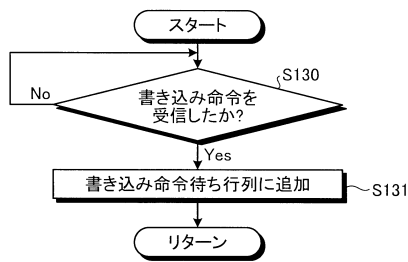
【図 17】



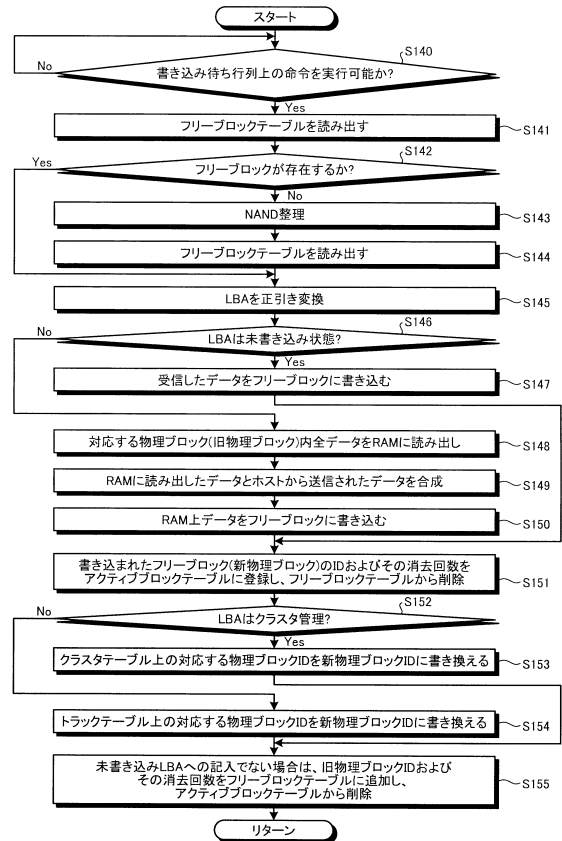
【図 18】



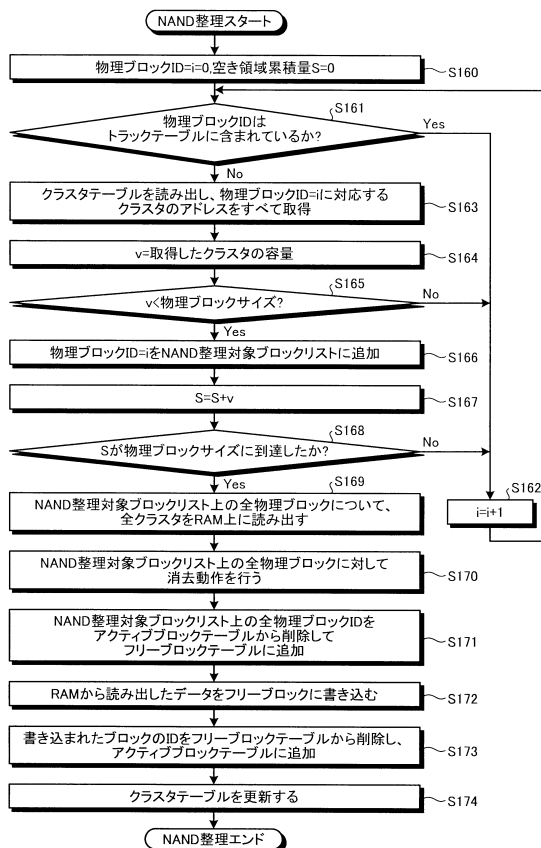
【図 19】



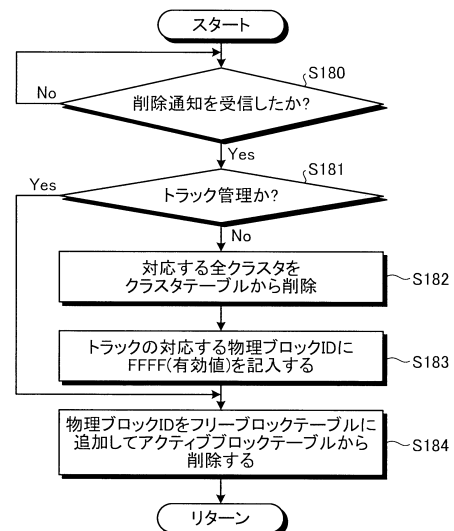
【図 20】



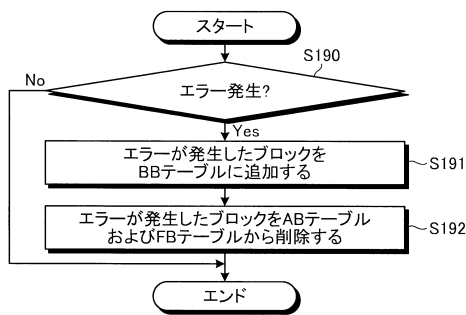
【図 21】



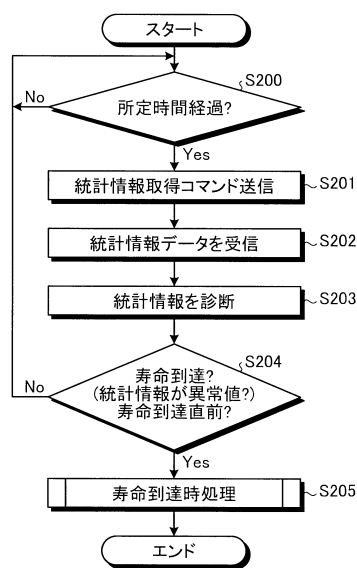
【図 22】



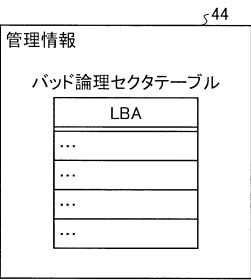
【図 2 3】



【図 2 4】



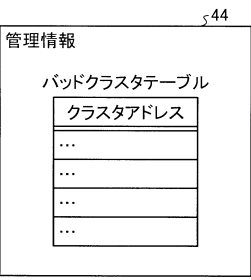
【図 2 5】



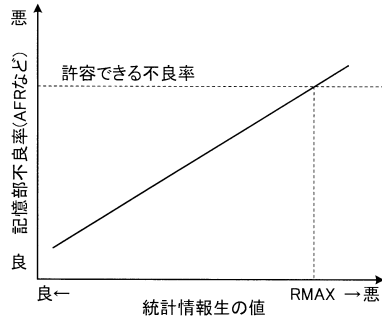
【図 2 7】

attribute ID	attribute 名	Value	Threshold	Worst	Raw Data
...	パッドブロック数総計 (統計情報X01)
...	パッド論理セクタ数総計 (統計情報X02)
...	消去回数総計 (統計情報X03)
...	消去回数平均値 (統計情報X04)
...	NANDメモリの書き込みエラー発生回数累積値 (統計情報X05)
...	NANDメモリの読み込みエラー発生回数累積値 (統計情報X06)
...	読み出しエラーセクタ数総計 (統計情報X07)
...	書き込みエラーセクタ数総計 (統計情報X08)
...	ECG訂正不能回数総計 (統計情報X09)
...	リトライ総計 (統計情報X10)
...	ペットポイントECG訂正単位総計 (統計情報X11)
...	インターフェース19のデータ化エラー発生回数 (統計情報X12)
...	インターフェース19の通信速度ダウン回数 (統計情報X13)
...	インターフェース19のレーン数ダウン回数 (統計情報X14)
...	インターフェース19のエラー発生回数 (統計情報X15)
...	RAM40のエラー発生回数 (統計情報X16)
...	記憶部2の使用時間総計 (統計情報X17)
...	起動回数 (統計情報X18)
...	不正電源発生回数 (統計情報X19)
...	温度が推奨動作温度の最高値を上回った時間累計 (統計情報X20)
...	温度が推奨動作温度の最低値を下回った時間累計 (統計情報X21)
...	温度が推奨動作温度の最大値 (統計情報X22)
...	コマンドの応答時間平均値 (統計情報X23)
...	NANDの応答時間最大値 (統計情報X24)
...	NANDの応答時間平均値 (統計情報X25)
...	発生温度 (統計情報X26)
...	最高温度 (統計情報X27)
...	最低温度 (統計情報X28)
...	管理情報冗長度 (統計情報X29)
...	RAM40への書き込みデータ量合計 (統計情報X30)
...	統計情報増加率 (統計情報X31)
...	NAND管理失敗フラグ (統計情報X32)

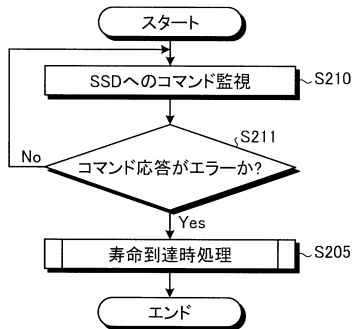
【図 2 6】



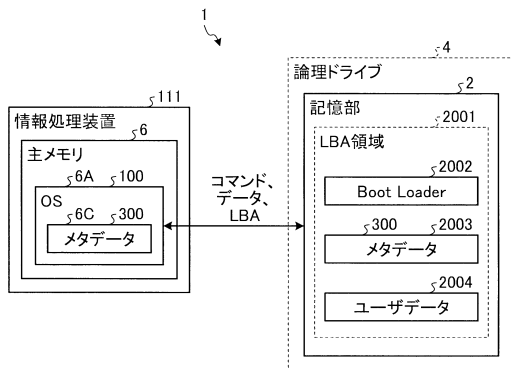
【図 28】



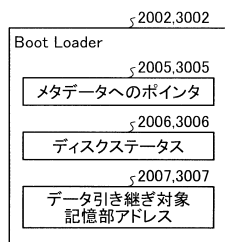
【図 29】



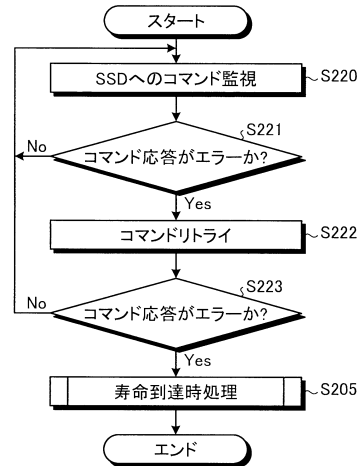
【図 31】



【図 32】



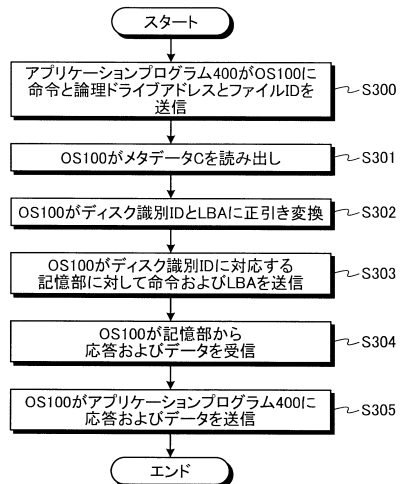
【図 30】



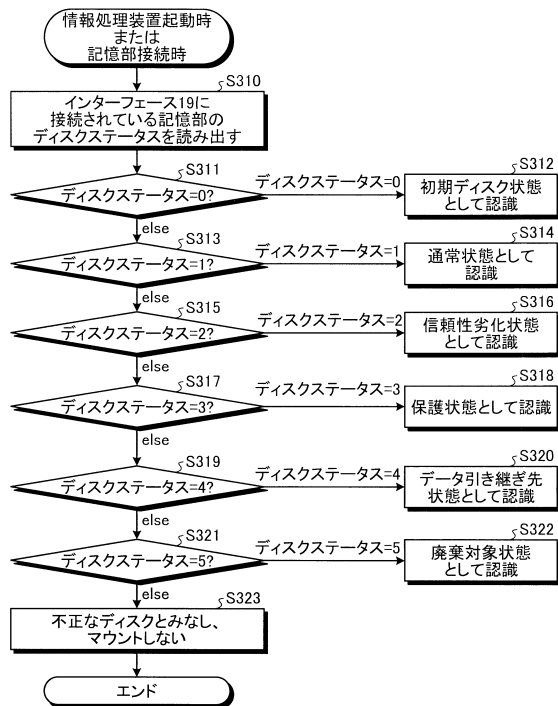
【図 33】

メタデータ				
ファイル識別ID	論理ドライブアドレス	ディスク識別ID	LBA	セクタカウント
...	論理ドライブ4	記憶部2
...	論理ドライブ4	記憶部2
...	論理ドライブ4	記憶部2
...	論理ドライブ4	記憶部2
...	論理ドライブ4	記憶部2

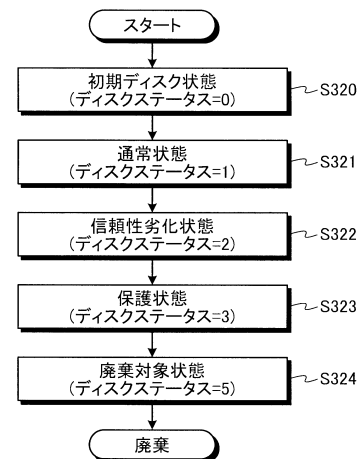
【図 34】



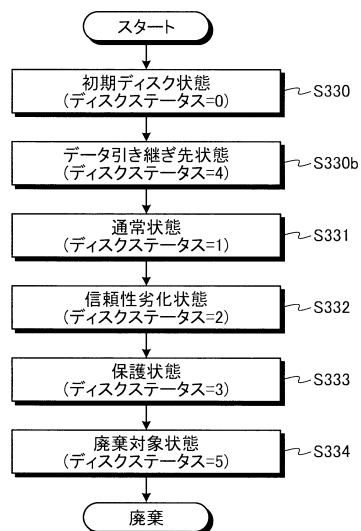
【図 35】



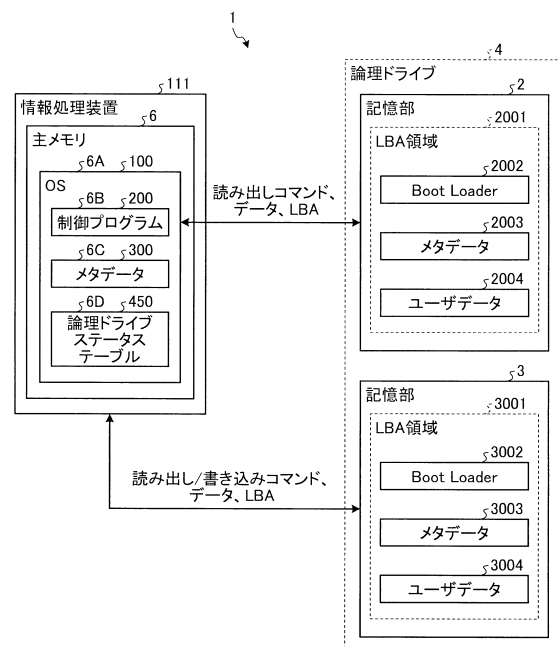
【図 36】



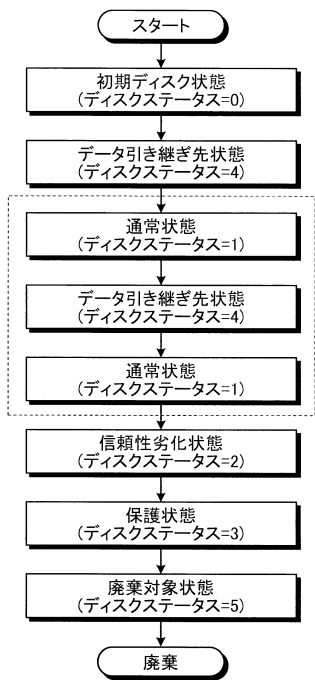
【図 37】



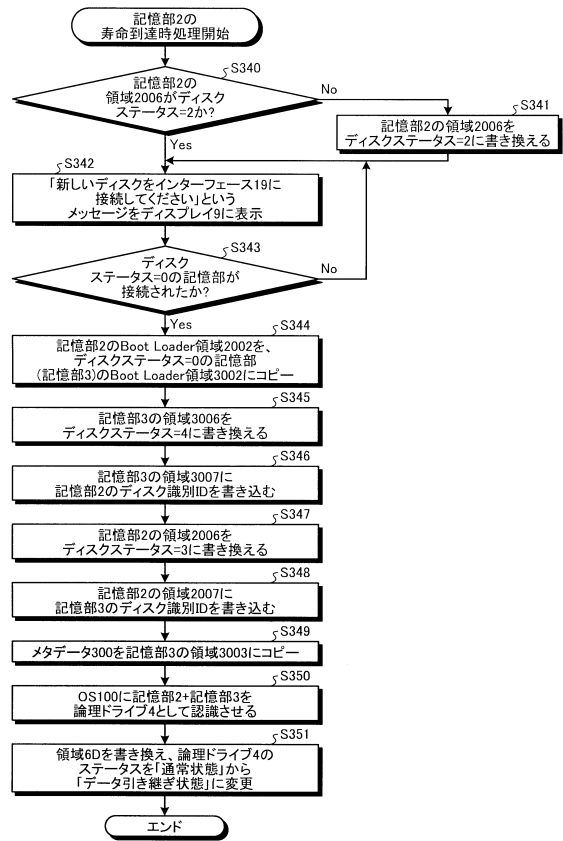
【図 38】



【 図 3 9 】



【 図 4 0 】

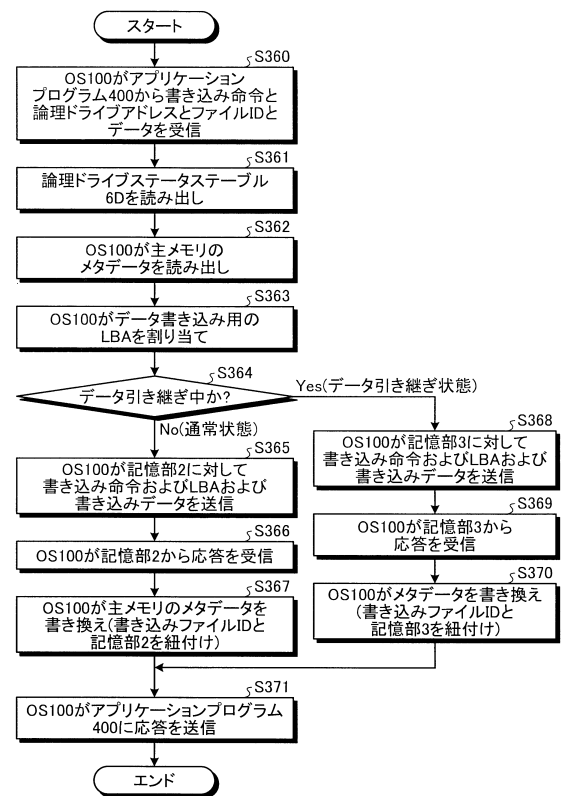


【 図 4 1 】

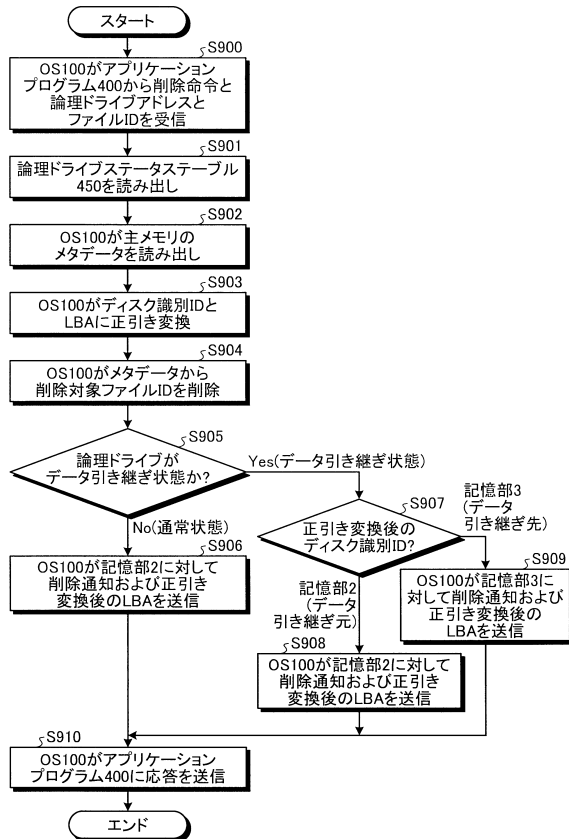
論理ドライブステータステーブル § 6D § 450

論理ドライブアドレス	ステータス
...	通常状態
...	データ引き継ぎ状態
...	通常状態

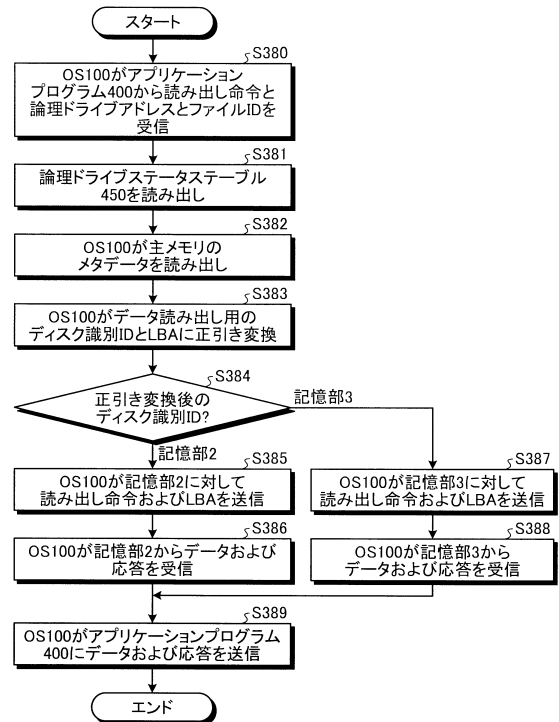
【 図 4 2 】



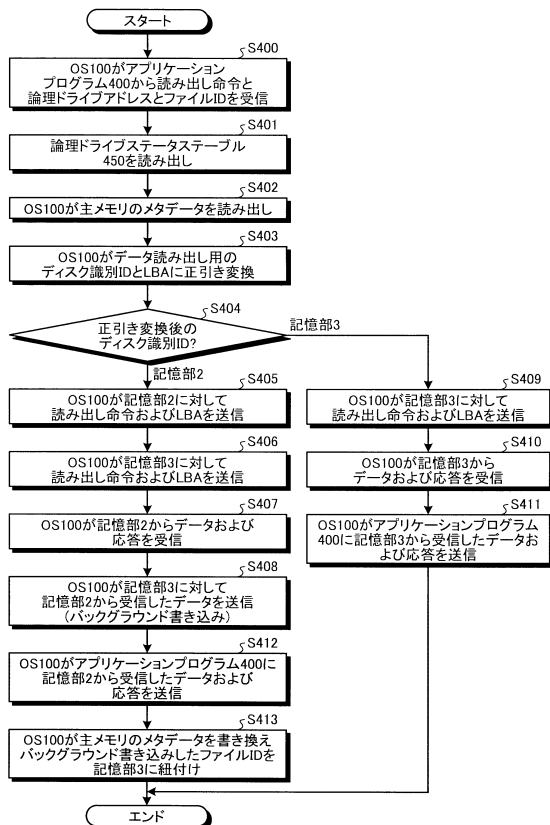
【図 4 3】



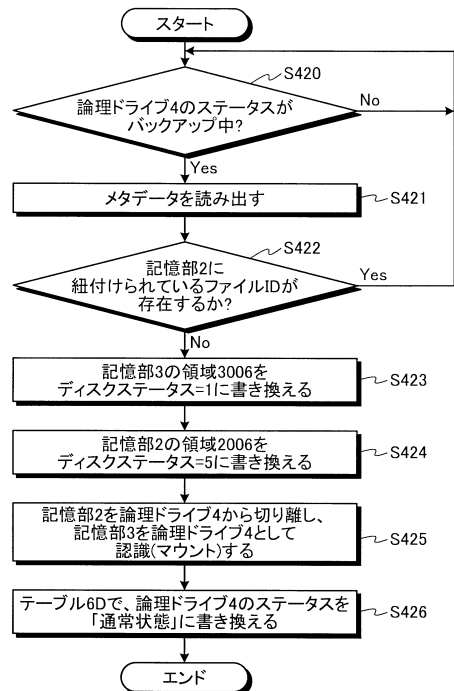
【図 4 4】



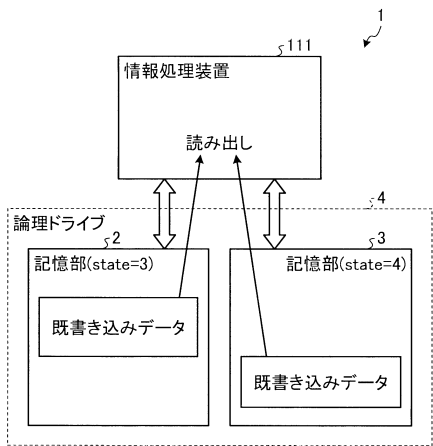
【図 4 5】



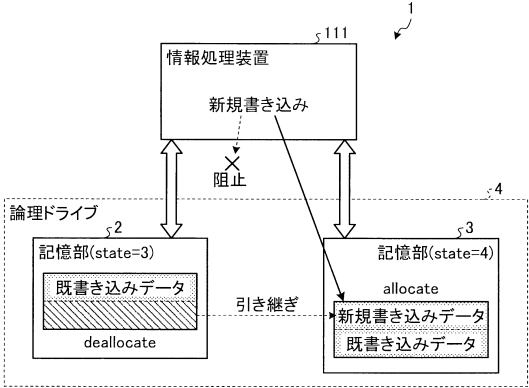
【図 4 6】



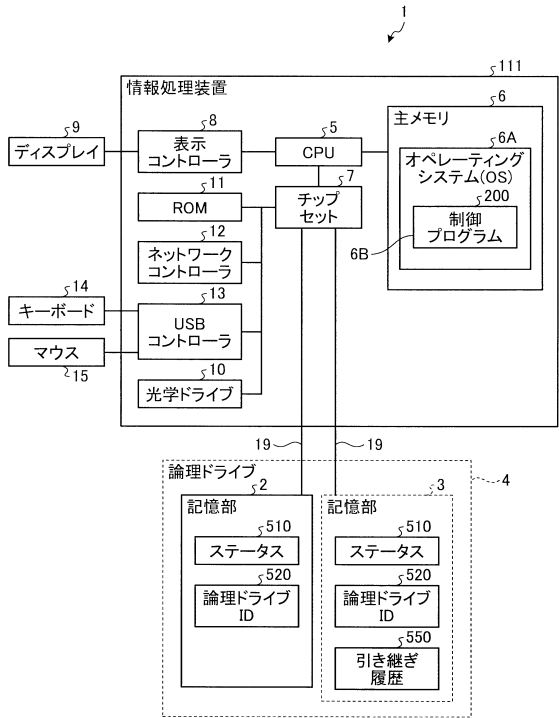
【図 47】



【図 48】



【図 49】

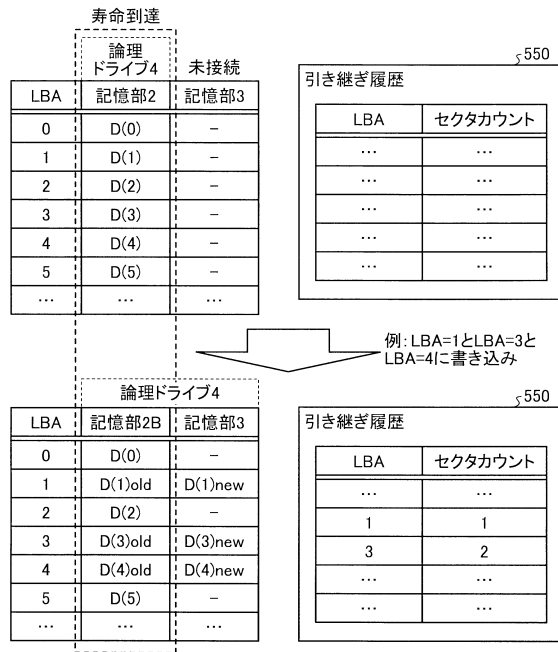


【図 50】

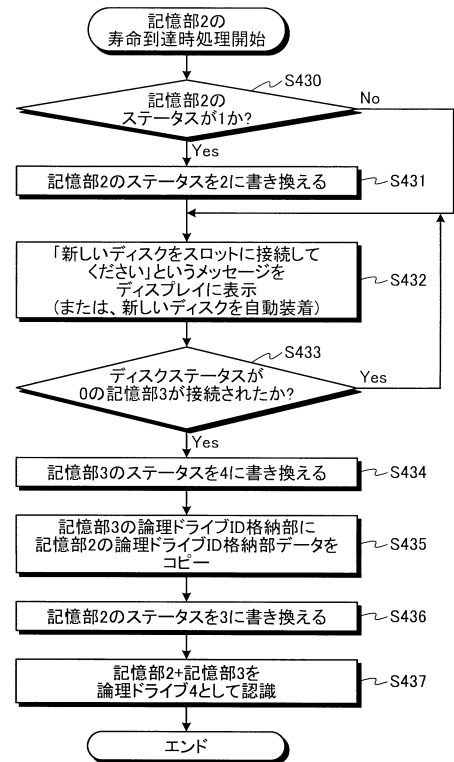
引き継ぎ履歴

LBA	セクタカウント
...	...
...	...
...	...
...	...

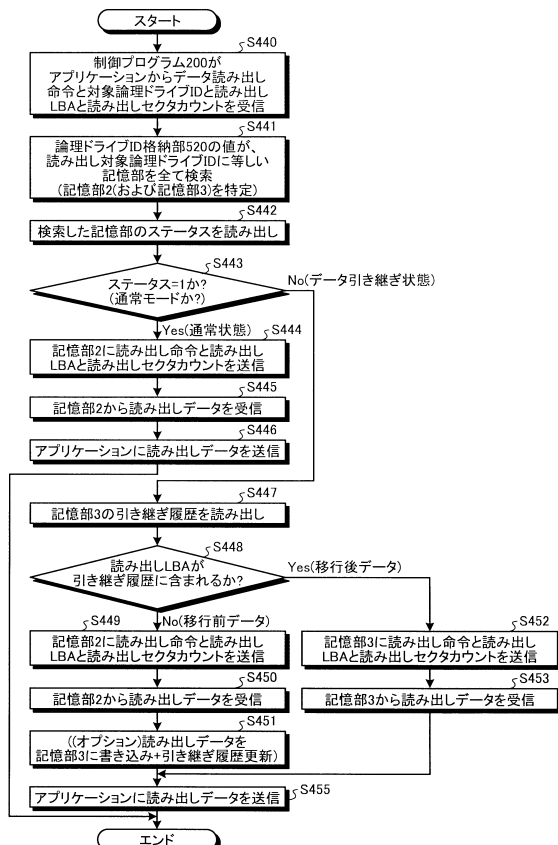
【図 5 1】



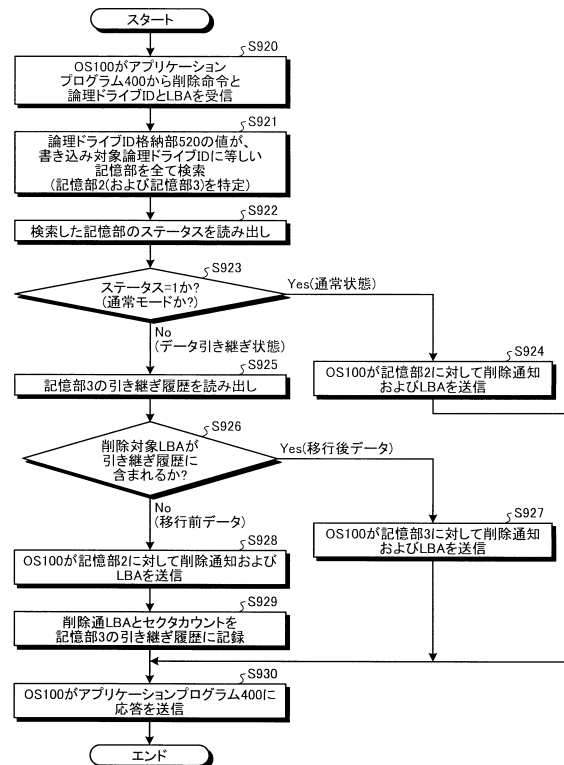
【図 5 2】



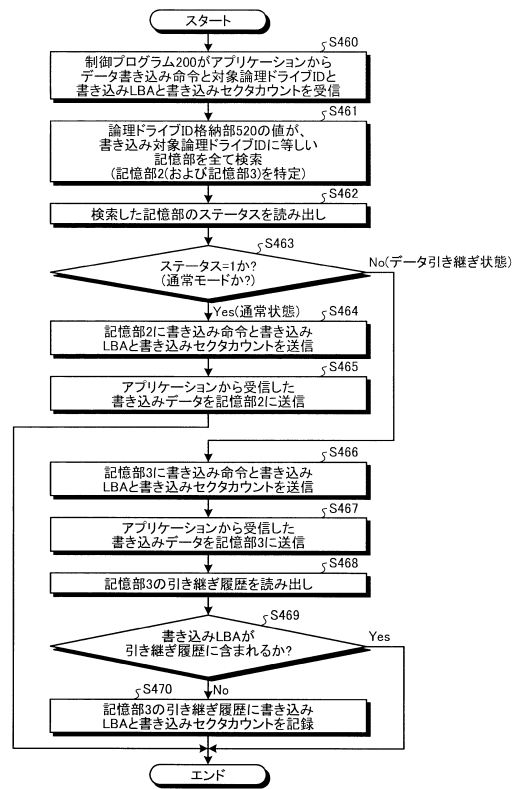
【図 5 3】



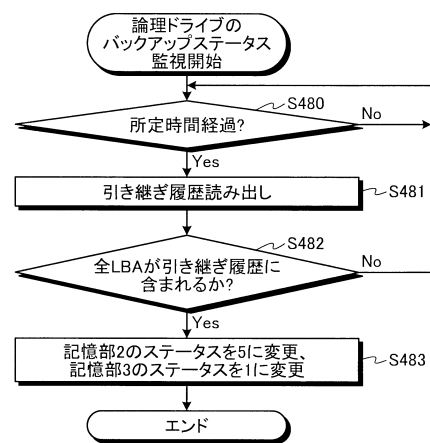
【図 5 4】



【図 5 5】



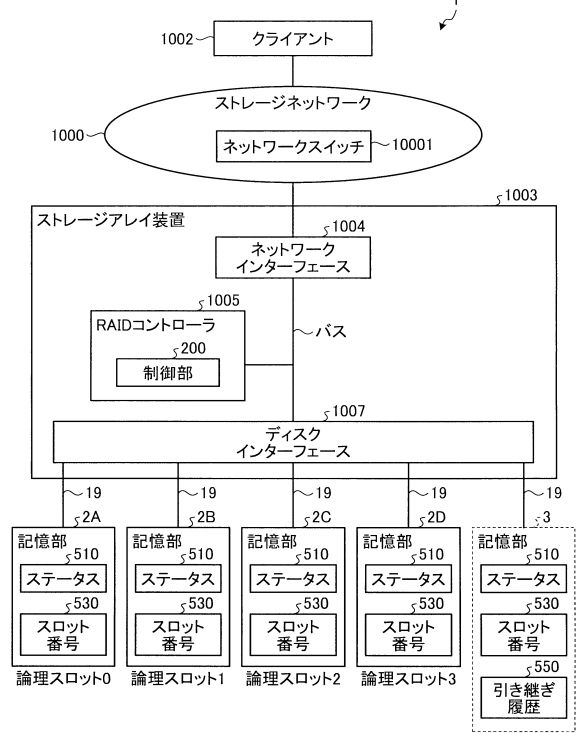
【図 5 6】



【図 5 7】

論理ドライブ4			
ステップ	記憶部2	記憶部3	状態
1	1	未接続	信頼性劣化前
2	2	未接続	記憶部3を接続するよう要求
3	3	0	記憶部3を接続直後
4	4	4	データ引き継ぎ中
5	5	1	データ引き継ぎ完了

【図 5 8】



【図 59】

	論理 スロット0	論理 スロット1	論理 スロット2	論理 スロット3
記憶部LBA	記憶部2A	記憶部2B	記憶部2C	記憶部2D
0	D(0)	D(1)	D(2)	P(0,2)
1	D(3)	D(4)	P(3,5)	D(5)
2	D(6)	P(6,8)	D(7)	D(8)
3	P(9,11)	D(9)	D(10)	D(11)
4	D(12)	D(13)	D(14)	P(12,14)
5	D(15)	D(16)	P(15,17)	D(17)
...

【図 61】

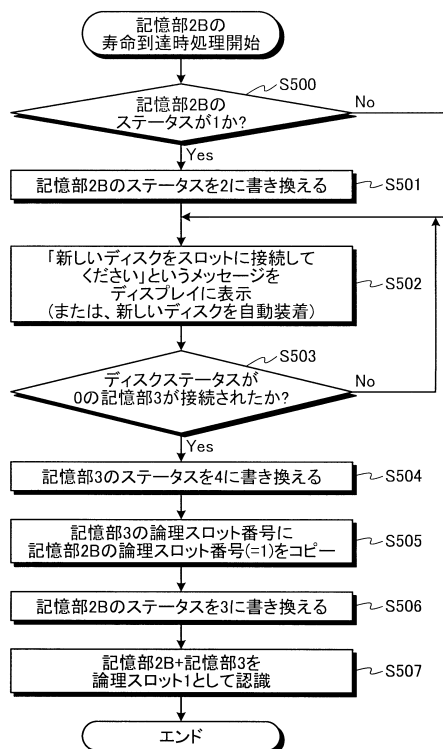
§550

引き継ぎ履歴	
記憶部LBA	セクタカウント
...	...
...	...
...	...
...	...

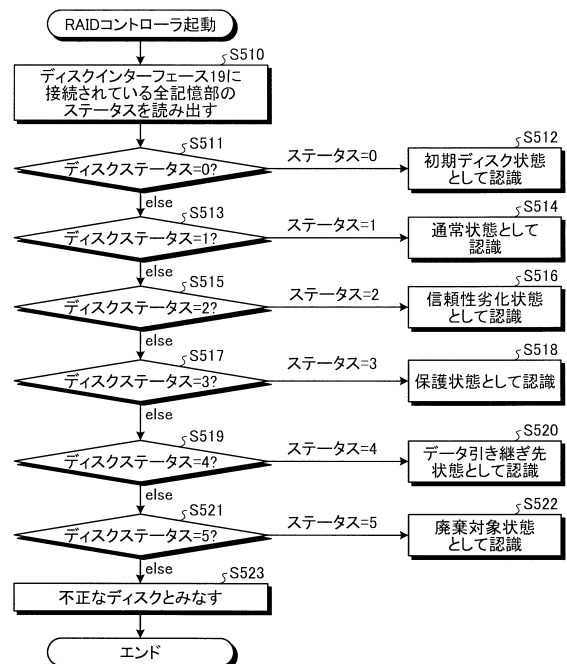
【図 60】

	論理 スロット0	論理 スロット1	論理 スロット2	論理 スロット3	論理 スロット1
記憶部LBA	記憶部2A	記憶部2B	記憶部2C	記憶部2D	記憶部3
0	D(0)	D(1)	D(2)	P(0,2)	-
1	D(3)	D(4)	P(3,5)	D(5)	-
2	D(6)	P(6,8)	D(7)	D(8)	-
3	P(9,11)	D(9)	D(10)	D(11)	-
4	D(12)	D(13)	D(14)	P(12,14)	-
...

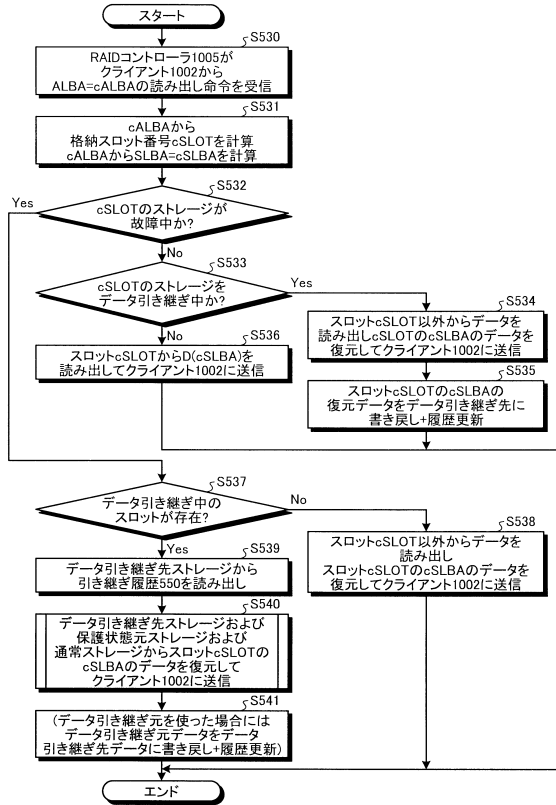
【図 62】



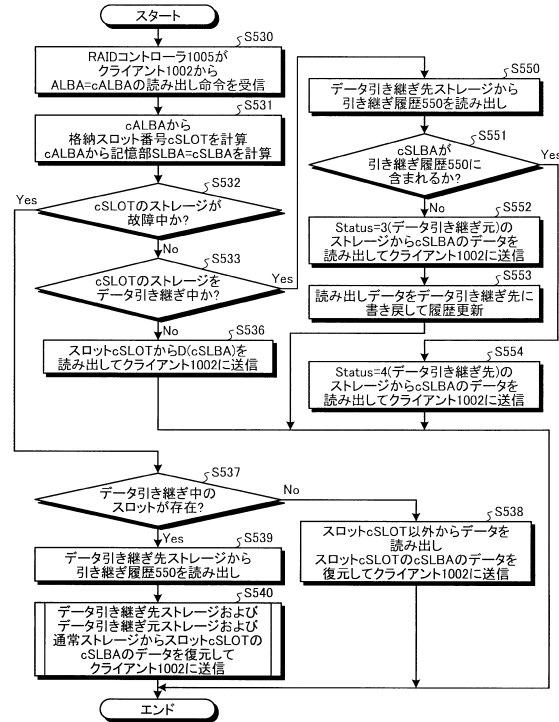
【図 63】



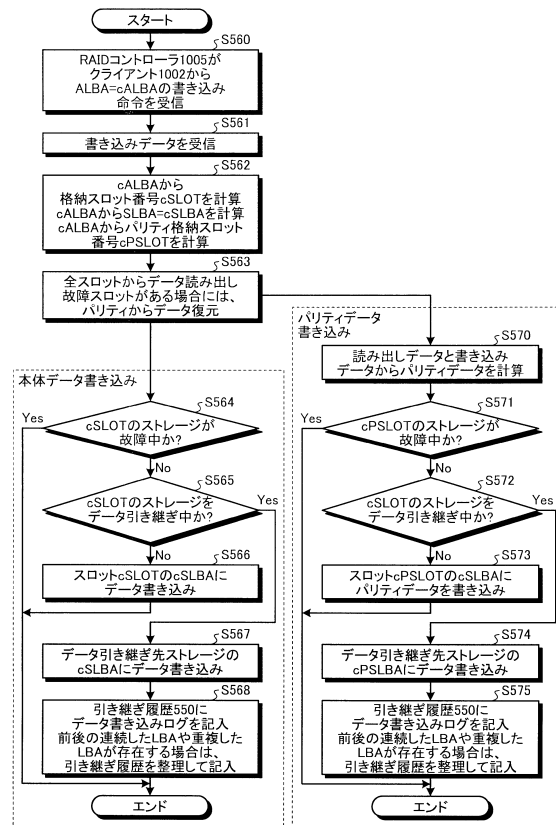
【図 6 4】



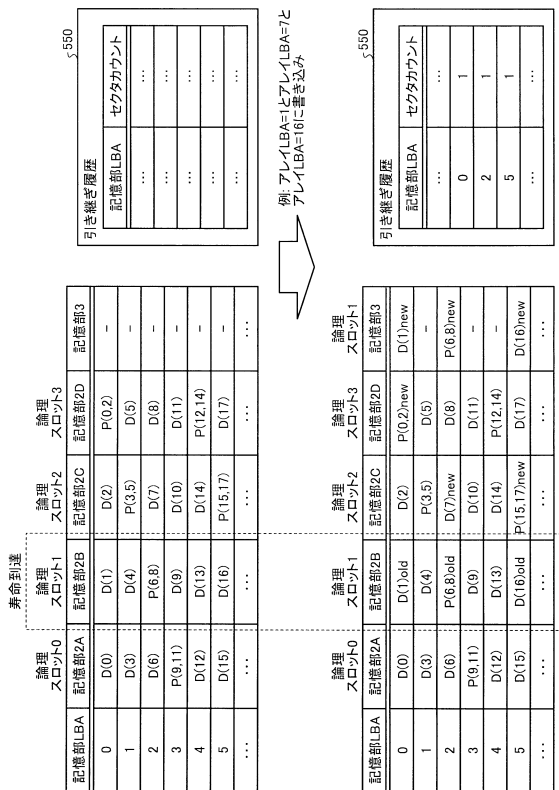
【図 6 5】



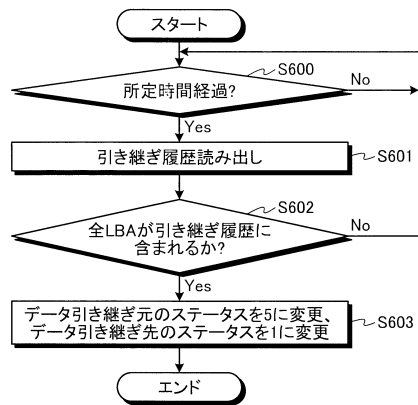
【図 6 6】



【図 6 7】



【図 68】



【図 69】

ステップ	信頼性劣化 データ引き継ぎ先 論理スロット1		信頼性劣化 データ引き継ぎ先 論理スロット2		信頼性劣化 データ引き継ぎ先 論理スロット3		状態
	記憶部2A スロット0	記憶部2B スロット1	記憶部2C スロット2	記憶部2D スロット3	記憶部3 スロット1	記憶部3 スロット2	
1	1	1	1	1	未接続	未接続	信頼性劣化前
2	1	2	1	1	未接続	未接続	記憶部3を接続するよう要求
3	1	2	1	1	0	0	記憶部3接続直後
4	1	3	1	1	4	4	データ引き継ぎ中 読み出しは記憶部2A/2C/2Dを 優先して行い、記憶部2Bの LBAデータはバリエーションから適宜復元
5	1	5	1	1	1	1	データ引き継ぎ完了

【図 70】

記憶部 LBA	信頼性劣化 データ引き継ぎ先 論理スロット1		信頼性劣化 データ引き継ぎ先 論理スロット2		信頼性劣化 データ引き継ぎ先 論理スロット3		読み出し時の挙動
	記憶部2A スロット0	記憶部2B スロット1	記憶部2C スロット2	記憶部2D スロット3	記憶部3 スロット1	記憶部3 スロット2	
0	D(0)	D(1)old	D(2)	P(0,2)new (故障)	D(1)new	D(0), D(2)復元不要 D(1)はデータ引き継ぎ先である記憶部3から読み出し	読み出し時の挙動
1	D(3)	D(4)	P(3,5)	D(5) (故障)	-	D(3), D(4)は復元不要 D(5)はD(3)とD(4)とP(3,5)をXORして復元	
2	D(6)	P(6,8)old	D(7)new	D(8) (故障)	P(6,8)new	D(6), D(7)は復元不要 D(8)はD(6)とD(7)とP(6,8)newをXORして復元	
3	P(9,11)	D(9)	D(10)	D(11) (故障)	-	D(9), D(10)は復元不要 D(11)はD(9)とD(10)とP(9,11)をXORして復元	
4	D(12)	D(13)	D(14)	P(12,14) (故障)	-	復元不要	
5	D(15)	D(16)old	P(15,17)new	D(17) (故障)	D(16)new	D(15)は復元不要 D(16)はデータ引き継ぎ先である記憶部3から読み出し D(17)はD(15)とD(16)とP(15,17)newをXORして復元	
...

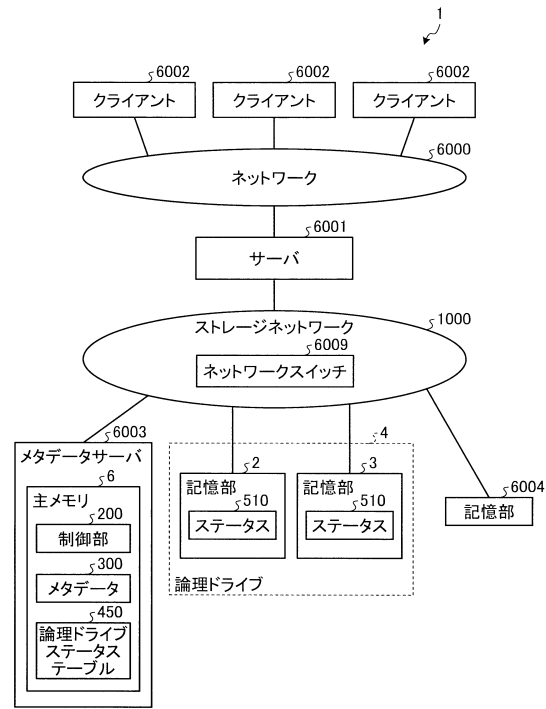
【図 71】

記憶部 LBA	信頼性劣化 データ引き継ぎ先 論理スロット1		信頼性劣化 データ引き継ぎ先 論理スロット2		信頼性劣化 データ引き継ぎ先 論理スロット3		読み出し時の挙動
	記憶部2A スロット0	記憶部2B スロット1	記憶部2C スロット2	記憶部2D スロット3	記憶部3B スロット2	記憶部3C スロット3	
0	D(0)	D(1)old	D(2)	P(0,2)new (故障)	D(1)new	-	読み出し時の挙動
1	D(3)	D(4)	P(3,5)	D(5) (故障)	-	-	
2	D(6)	P(6,8)old	D(7)old	D(8) (故障)	P(6,8)new	D(7)new	
3	P(9,11)	D(9)	D(10)	D(11) (故障)	-	-	
4	D(12)	D(13)	D(14)	P(12,14) (故障)	-	-	
5	D(15)	D(16)old	P(15,17)old	D(17) (故障)	D(16)new	P(15,17) new	
...

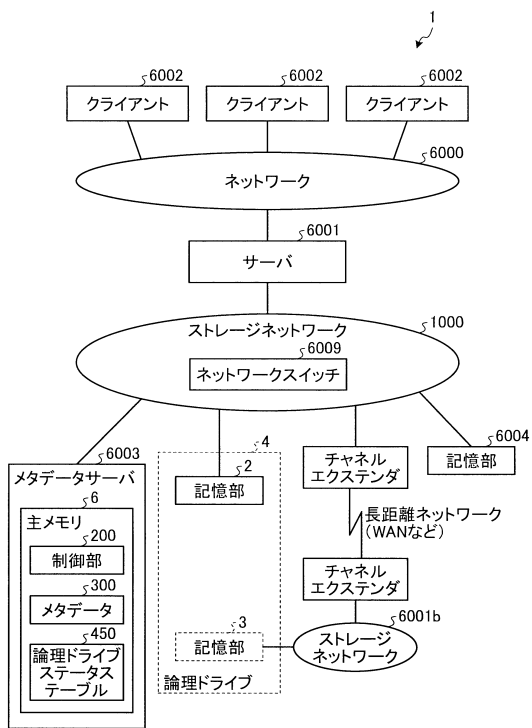
【図 7 2】

論理データ引き継ぎ元		論理データ引き継ぎ先		読み出し時の挙動	
スロット0	記憶部2A	記憶部2B	記憶部2C	スロット3	記憶部3
記憶部LBA					
0	D(0)	D(1)old	D(2)UECC	D(1)new	ALBA=2の読み出しでECC訂正不能エラー →D(0)とD(1)newとP(0,2)newをXORして復元
1	D(3)	D(4)	P(3,5)	D(5)	-
2	D(6)	P(6,8)old	D(7)new	P(6,8)new	ALBA=7の読み出しでECC訂正不能エラー →D(6)とD(7)newとP(6,8)newをXORして復元
3	P(9,11)	D(9)UECC	D(10)	D(11)	ALBA=9の読み出しでECC訂正不能エラー →D(10)とD(11)とP(9,11)をXORして復元
4	D(12)	D(13)	D(14)	P(12,14)	-
5	D(15)	D(16)old	P(15,17)new	D(17)	-
...

【図 7 3】



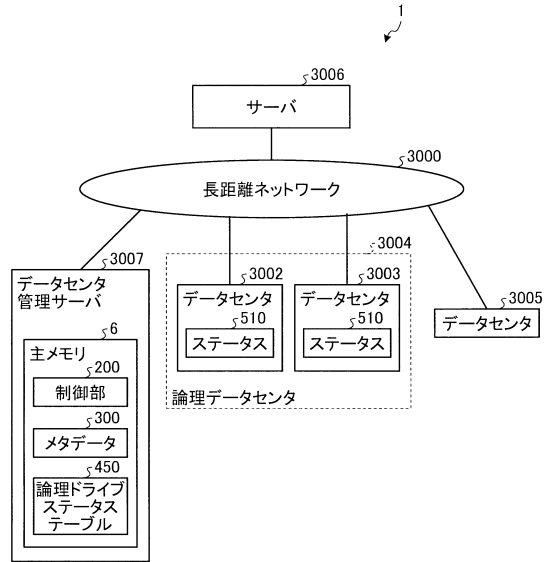
【図 7 4】



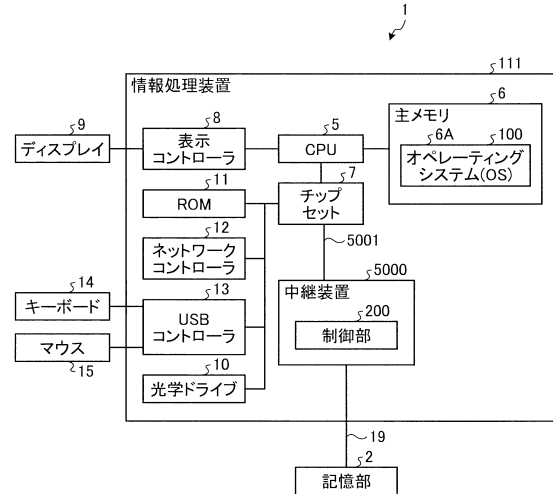
【図 7 5】

ステータステーブル			
論理ドライブアドレス	ディスク認識ID	論理ドライブステータス	ディスクステータス
A	a	通常状態	1(通常状態)
B	b1	データ引き継ぎ中	3(データ引き継ぎ元)
B	b2	データ引き継ぎ中	4(データ引き継ぎ先)
C	c	通常状態	1(通常状態)

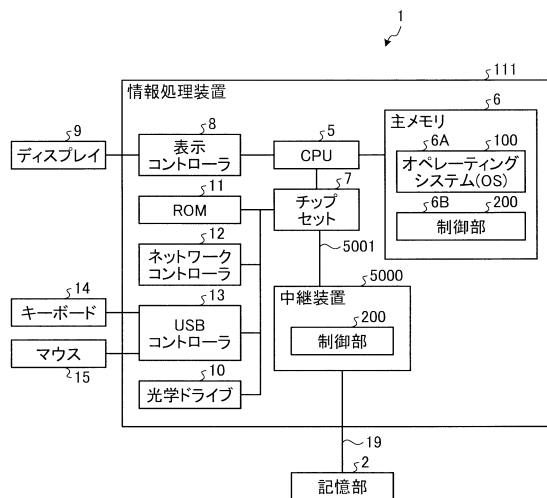
【図 76】



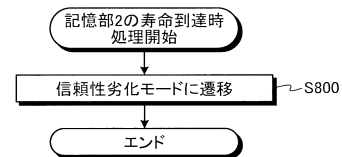
【図 77】



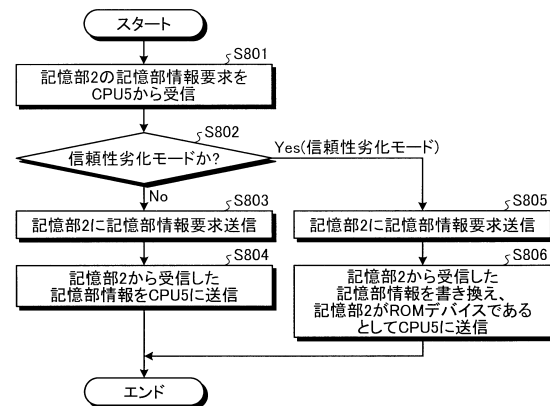
【図 78】



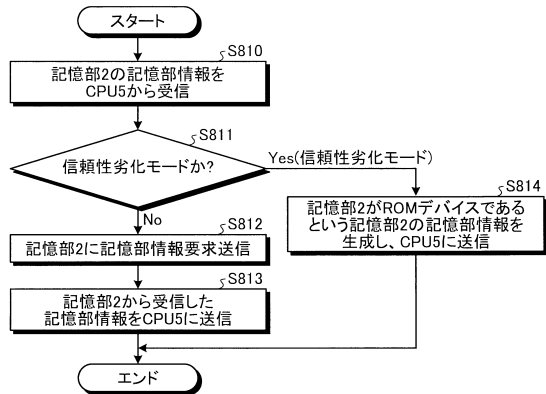
【図 79】



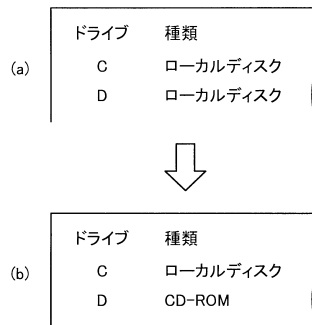
【図 80】



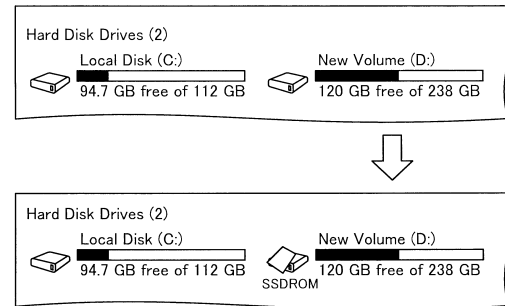
【図 8 1】



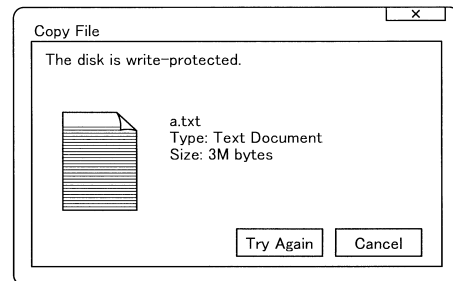
【図 8 2】



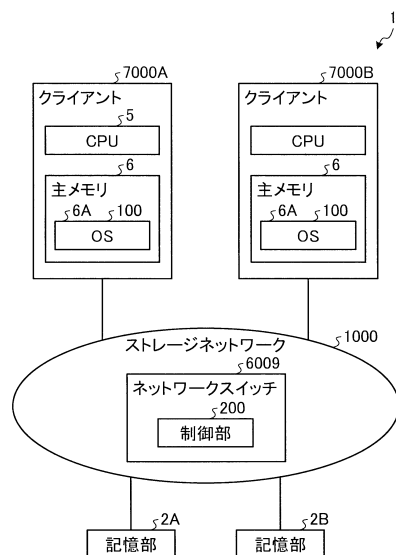
【図 8 3】



【図 8 4】



【図 8 5】



フロントページの続き

審査官 滝谷 亮一

(56)参考文献 米国特許出願公開第2010/0306577(US, A1)
特開2012-022619(JP, A)
特表2012-509521(JP, A)

(58)調査した分野(Int.Cl., DB名)
G06F 12/16
G06F 3/06
G06F 13/10